# International Journal of

# **Communications, Network** and System Sciences

ISSN: 1913-3715 Vol. 1, No. 2, June 2008





Scientific Research Publishing

# **JOURNAL EDITORIAL BOARD**

ISSN 1913-3715 (Print) ISSN 1913-3723 (Online) Http://www.srpublishing.org/journal/ijcns/

Editors-in-Chief	
Prof. Tom Hou	Department of Electrical and Computer Engineering, Virginia Tech., USA
Prof. Huaibei Zhou	Advanced Research Center for Sci. & Tech., Wuhan University, China
Editorial Board	
Prof. Ji Chen	Department of Electrical Engineering, University of Houston, USA
Prof. Jong-Wha Chong	College of Information & Communications, Hanyang University, Korea (South)
Prof. Laurie Cuthbert	Department of Computer Science, University of London at Queen Mary, UK
Prof. Klaus Doppler	Radio Communications CTC, Nokia Research Center, Nokia Corporation, Finland
Prof. Thorsten Herfet	Telecommunications Lab, Saarland University, Germany
Dr. Li Huang	Holst Centre, Stiching IMEC Nederland, Netherlands
Dr. Yi Huang	Department of Electrical Engineering and Electronics, University of Liverpool, UK
Prof. Chun Chi Lee	Department of Computer and Communication, Shu-Te University, Taiwan (China)
Prof. Myoung-Seob Lim	Faculty of Electronics & Information Engineering, Chonbuk National University, Korea
	(South)
Prof. Zhihui Lv	School of Information Science and Engineering, Fudan University, China
Dr. Kosai Raoof	Laboratoire des Images et Signaux, France
Prof. Heung-Gyoon Ryu	School of Electrical, Electronic and Computer Engineering, Chungbuk National
	University, Korea (South)
Prof. Rainer Schoenen	RWTH Aachen University, Germany
Dr. Lingyang Song	University Graduate Center, Norway
Prof. Dash Wu	Center for Risk Analysis, University of Toronto, Canada
Dr. Hassan Yaghoobi	Mobile Wireless Group, Intel Corporation, USA
Editorial Assistant	
Ting Chen	
Ruoshan Kong	

#### **Guest Reviewers**

Loren Schwiebert Wayne State University A. N. Skodras Hellenic Open University Zhiguo Wan K.U.Leuven, Belgium Damla Turgut University of Central Florida

Qingshan Shan University of Strathclyde Minghua Xia ETRI Beijing R&D Center Valerio Bellandi University of Milan Gerard Rowe University of Auckland Prabu D PTG Lab (R&D) Mugen Peng Beijing University of Posts & Telecommunication Linfeng Yuan Huazhong University of Science and Technology Chong Shen Cork Institute of Technology Ireland

*I. J. Communications, Network and System Sciences*, 2008, 2, 105-206 Published Online May 2008 in SciRes (http://www.SRPublishing.org/journal/ijcns/).

### **TABLE OF CONTENTS**

#### Volume 1

May 2008

Novel Joint Chip Sampling and Phase Synchronization Algorithm for Multistandard UMTS Systems	
Y. SERRESTOU, K. RAOOF, J. LIÉNARD	105
Performance Analysis of an AMC System with an Iterative V-BLAST Decoding Algorithm	
S.J. RYOO, K.W. LEE, I. HWANG	119
Impact of Depolarization Phenomena on Polarized MIMO Channel Performances N. PRAYONGPUN, K. RAOOF	124
A Quadratic Constraint Total Least-squares Algorithm for Hyperbolic Location K. YANG, J.P. AN, Z. XU	130
Analysis of Lifetime of Large Wireless Sensor Networks Based on Multiple Battery Levels R.H. ZHANG, Z.P JIA1, D.F. YUAN	136
Beacon-driven Leader Based Protocol over a GE Channel for MAC Layer Multicast Error Control Z. LI, T. HERFET	144
Routing and Wavelength Assignment in GMPLS-based 10 Gb/s Ethernet Long Haul Optical Networks with and without Linear Dispersion Constraints L.N. BINH	154
On Modeling and Accuracy Analysis of the Available Bandwidth Measurement Based-on Packet-pair Sampling	
J. LIU, D.F. ZHANG, J.H. JIN	168
Streaming Multimedia over Wireless Mesh Networks D.Q. LIU, J. BAKER	177
A Novel Adaptive Hybrid Error Correction Scheme for Wireless DVB Services G.P. TAN, T. HERFET	187
A Co-verification Method Based on TWCNP-OS for Two-way Cable Network SOC C. LI, X.T. ZHANG, Y.D WAN, Q. WANG	199

# International Journal of Communications, Network and System Sciences (IJCNS)

#### **Journal Information**

#### **SUBSCRIPTIONS**

The International Journal of Communications, Network and System Sciences (Online at Scientific Research Publishing, www.SRPublishing.org) is published quarterly by Scientific Research Publishing, Inc. 3306 Apple Grove CT, Herndon, VA 20171, USA.

E-mail: jcnss@srpublishing.org

#### Subscription rates: Volume 1 2008

Print: \$50 per copy. Electronic: free, available on www.SRPublishing.org. To subscribe, please contact Journals Subscriptions Department, E-mail: jcnss@srpublishing.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: jcnss@srpublishing.org

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

Reprints Co-ordinator, Scientific Research Publishing, Inc. 3306 Apple Grove CT, Herndon, VA 20171, USA. E-mail: jcnss@srpublishing.org

#### COPYRIGHT

Copyright© 2008 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: jcnss@srpublishing.org



## Novel Joint Chip Sampling and Phase Synchronization Algorithm for Multistandard UMTS Systems

Youssef SERRESTOU<sup>1</sup>, Kosai RAOOF<sup>2</sup>, Joël LIÉNARD<sup>2</sup>

 <sup>1</sup>LCIS-INPG, 50 Rue Barthélémy de Laffemas, 26902 cedex, Valence, France
 <sup>2</sup>GIPSA-LAB, CNRS UMR 5216, 961 Rue de Houille Blanche, 38402 St. Martin d'Hères, France E-mail: <sup>1</sup>youssef.serrestou@esisar.inpg.fr, <sup>2</sup>kosai.raoof@gipsa-lab.inpg.fr

#### Abstract

CDMA Timing and phase offsets tracking remain as one of considerable factors that influence the performances of communication systems. Many algorithms are proposed to solve this problem. In general, these solutions process separately the chip sampling offset and phase rotation. In addition, most of proposed solutions can not assure a compromise between robustness criteria and low complexity for implementation in real time applications. In this paper we present an efficient algorithm for chip sampling instant and phase errors. The robustness and the low complexity of this algorithm are evaluated, firstly by simulation and then tested by real experimentation for UMTS standard. Simulation results show that the proposed algorithm provides very efficient compensation for sampling clock offset and phase rotation. A real time implementation is achieved, based on TigerSharc DSP, while using a complete UMTS transmission-reception chain. Experimental results show robustness in real conditions.

Keywords: Synchronization, DS-CDMA, UMTS, Joint Estimation, Early-late Loop, Phase Locked Loop

#### 1. Introduction

Code Division Multiple Access (CDMA) system offers several advantages such as robustness to multi-path fading channel and multiple access properties [1]. Thanks to these advantages, CDMA are adopted in high data rate applications such as UMTS standard (Universal Mobile Telecommunication System), and it is a good candidate for future applications, such as the 4<sup>th</sup> generation phone mobile and wireless LAN [2]. However, the performance of such systems is affected by the time sampling offsets and phase rotation, due to many factors such as channel fading, multipath problem, clock imperfection, etc.

Recent research works are interested in solving this problem of synchronization for different applications [3– 7]. In general the synchronization techniques of DS-CDMA systems can achieve synchronization in two distinct steps: acquisition step and tracking step. The acquisition step, or the initial synchronization step, involves in determining the timing/phase offsets of the incoming signal to within a specified range known as the pull-in region of tracking loop [8-11]. Upon the successful completion of the acquisition step, the tracking step refines and maintains synchronization. When the problem of tracking is treated, generally, phase offset tracking is omitted or treated separately with a more complex computation than delay tracking [3-7, 12, 13]. The tracking step is based on recursive estimation of timing/phase error and feedback correction. So the performance of synchronization depends greatly on the properties of estimators and on the precision of correctors. To estimate time delay the conventional delay lock loop (DLL) has been extensively used. In a DLL the received signal is cross-correlated with an early and a late copy of a pilot signal. The proposed tracking algorithm is inspired from the concept of lock loop to establish estimation [14]. So the complex correlation function of the received signal and pilot signal is computed in three points (which corresponds to original, early and late copy of the pilot signal). This computation allows estimating jointly the time delay and the phase rotation. A second order filters are used to update the estimation. For the feedback corrector, an interpolator [15] is used to find the sample at the estimated instant;

the used interpolator gives good compromises between precision and real time application. The correction of phase rotation consists of multiplying the received signal by the complex exponential of estimated phase [14,15].

In this study we are interested in implementing the tracking algorithm in a complete transmission-reception chain. The aim of this chain is to demonstrate the feasibility of a reconfigurable software radio communication system based on UMTS standard. The tracking algorithm was adapted to UMTS specifications [16–20]. In simulation context, the performances of these algorithms were evaluated in the presence of a fading and additive white Gaussian noise channel (AWGN). In the case of a multipath channel, acquisition and tracking blocks synchronize the signal of the high power path. One block within the complete chain contains an algorithm for channel estimation [21,22].

DS-CDMA was adopted in UMTS standard as an access scheme. And for radio access there are two modes, FDD (Frequency Division Duplex) and TDD (Time Division Duplex), for operating with paired and unpaired bands respectively. The possibility to operate in either FDD or TDD mode allows an efficient use of the available spectrum according to the frequency allocation in different regions. The modulation scheme is QPSK, and because of the CDMA nature the spreading (and scrambling) process is closely associated with modulation. To spread transmitted signal different families of spreading codes are used. For separating channels from the same source, orthogonal codes are derived with code tree structure. And for separating different cells in FDD mode Gold codes with 10 ms period (38400 chips at 3.84 Mcps) are used and for TDD mode we implement codes with period of 16 chips and midamble sequences of different length depending on the environment [16-20]. Our tracking algorithm has adopted for both UMTS modes TDD and FDD. Simulation has allowed evaluating the performances and determining suitable parameters for implementation. The final real time implementation was accomplished with a 32-bit TigerSharc DSP.

The rest of this paper is organised as follows. Section 2 presents the algorithm context. The basic aim of this section is to present the theoretical models of transmitted/received signals in UMTS and the communication channel that we consider. So the different phases of signal forming (spreading, scrambling, channel effects, noise) are described. In section 3, the mathematical model of the tracking algorithm is presented, as well as it's adaptation for UMTS-FDD and UMTS-TDD standard. We focus on the estimators' performances with theoretical demonstration ended by simulation to show the efficiencies of this algorithm. Section 4 describes simulation platform and result analysis. The implementation and experiment results are detailed in section 5. Some conclusions are given in section 6.

#### 2. Base-Band Signal Model in UMTS

#### 2.1. Introduction

Our tracking algorithm tries to maintain synchronization of the received signal by a UMTS terminal. So we are interested in the downlink mode (base station to terminal). We introduce some useful notions for modeling the base-band signal in the context of UMTS standard. These notions are specified and described in detail in the 3rd Generation Partnership Project (3GPP); see [16–20].

The access scheme is Direct-Sequence Code Division Multiple Access (DS-CDMA). There are two radio access modes, FDD (Frequency Division Duplex) and TDD (Time Division Duplex), which allows operating with paired and unpaired bands respectively. The modulation scheme is QPSK, and different families of spreading codes are used. These codes are OVSF (orthogonal variable spreading factor). And for separating different cells in FDD mode Gold codes (scrambling codes) with 10 ms period (38400 chips at 3.84 Mega chip per second Mcps) are used and for TDD mode we use codes with period of 16 chips and midamble sequences of different length depending on the environment [16-20]. In FDD mode the uplink and downlink transmissions use two different radio frequencies. In TDD mode the uplink and downlink transmissions are modulated over the same frequency by using synchronized time intervals.

In both modes the transmitted signal is decomposed into radio frames. A radio frame has duration of 10 ms and contains 38400 chips at the rate of 3.84 Mcps. It is divided into 15 slots, and every slot contains 2560 chips.

A physical channel is therefore defined as a simple code (or number of codes used to spread this channel), additionally in TDD mode the sequence of time slots completes the definition of a physical channel.

The information rate, in terms of symbols/s, varies with the spreading factor (number of chips by symbol). Spreading factors (noted SF) are from 512 to 4 for FDD downlink, and from 16 to 1 for TDD downlink. Thus the respective modulation symbol rates vary from 960 k symbols/s to 15 k 7.5 k symbols/s for FDD downlink, and for TDD the momentary modulation symbol rates shall vary from 38400 symbols/s to 240 k symbols/s.

#### 2.2. Signal Model in UMTS-FDD Mode

In FDD mode a practical fixed rate channel with predefined symbol (1+j) sequence is transmitted within each frame, this common pilot channel (CPICH), use the first spreading code in tree with spreading factor equal to 256. The proprieties of this channel allows receiver to identify cell's location. In our work we have used this channel as signal pilot to estimate the time delay and

phase rotation.

Let N physical channels be transmitted with QPSK symbols  $\{d_i^{(n)} = \pm 1 \pm j\}$ , n denotes the channel number and l represents the symbol number.

Every symbol is transformed to a sequence of chips by spreading operation. SF "spreading factor" denotes the spreading code length, *i.e.* the number of chip per symbol. The combination of all channels is scrambled by the same code that we denoted SC. The following combined signal is obtained:

$$d(t) = \sum_{n}^{N} G_{n} \sum_{l}^{N_{n}} d_{l}^{(n)} \sum_{k}^{SF_{n}} C_{k,SF_{n}}^{n} SC_{k+l \cdot SF_{n}} \delta(t - (k + l \cdot SF_{k}) \cdot T_{c})$$
(1)

where  $G_n$  is the gain of channel n,  $C_{k,SF_n}^n \in \{-1,1\}$  is the k<sup>th</sup> chips in spreading code used to spread the n<sup>th</sup> channel,  $SF_n$  is the spreading factor used for the n<sup>th</sup> channel and  $SC_k \in \{-1,1\}$  denotes the k<sup>th</sup> chip in scrambling code. T<sub>c</sub> is chip duration and  $\delta(.)$  is the Dirac delta function.

A root raised cosine impulse filter, denoted  $h_e(t)$ , is used to shape this signal for transmission.

For propagation channel, we adopt a fading channel with AWGN noise. Its impulse response is modeled as follow:

$$h_{c}(t) = \alpha(t)e^{j\varphi(t)}\delta(t - \tau(t))$$
<sup>(2)</sup>

where  $\alpha(t)$  is the fading of propagation channel given from Jake's Doppler spectrum,  $\tau(t)$  is the time delay, and  $\varphi(t)$  represents the random phase shift introduced by propagation, and it will be added to another phase shift which is caused by the imperfections of receivertransmitter oscillators.

The received signal is filtered by a root raised cosine with impulse response  $h_r(t)$ . Then it is sampled with a period  $T_e$ . We denote by OSF= $T_c/T_e$  the over sampling factor, i.e. number of samples per chip. Let g(t) be the convolution product  $Cr_0(t)*h_c(t)*h_r(t)$ , the received signal may be written as:

$$y_{r,FDD}(t) = e^{\varphi(t)} \left\{ \alpha(t) \sum_{n}^{N} G_{n} \sum_{l}^{N_{n}} a_{l,n} \sum_{k}^{SF_{n}} \sum_{m}^{OSF} C_{k,SF_{k}}^{n} SC_{k+lSF_{k}} g(t-\tau(t) - OSF(k+lSF_{k})T_{e} - mT_{e}) + \eta(t) \right\}$$
(3)

$$y_{r,FDD}(t) = \underbrace{y_{r0}(t)}_{other \ channel} + e^{\varphi(t)}\alpha(t)G_{cpich}\sum_{l}^{N_{n}}(1+j)\sum_{k}^{256}\sum_{m}^{OSF}SC_{k+256l}g(t-\tau(t)-OSF(k+l256)T_{e}-mT_{e})$$

where

$$y_{r0}(t) = e^{\varphi(t)} \left\{ \alpha(t) \sum_{n}^{N} G_{n} \sum_{l}^{N_{n}} a_{l,n} \sum_{k}^{SF_{n}} \sum_{m}^{OSF} C_{k,SF_{k}}^{n} SC_{k+lSF_{k}} g(t-\tau(t) - OSF(k+lSF_{k})T_{e} - mT_{e}) + \eta(t) \right\}$$

$$y_{r,TDD}^{k}(t) = \alpha(t)e^{j\varphi(t)} \left[ \sum_{\substack{n=0\\ n=0}}^{N_{k}} d_{n}^{(k,1)}w^{k} \sum_{i=1}^{16} s_{i}^{k} \cdot \sum_{l=0}^{OSF^{-1}} g(t - (i - 1)OSFT_{e} - 16OSFnT_{e} - lT_{e} - \tau(t)) + \frac{1}{data} \right]$$

$$\sum_{\substack{n=1\\ n=1}}^{L_{m}} m_{n}^{k} \sum_{l=0}^{OSF^{-1}} g(t - (i - 1)OSFT_{e} - 16OSFnT_{e} - lT_{e} - \tau(t)) + \frac{1}{midamble}$$

$$\sum_{\substack{n=0\\ n=0}}^{N_{k}} d_{n}^{(k,2)}w^{k} \sum_{i=1}^{16} s_{i}^{k} \cdot \sum_{l=0}^{OSF^{-1}} g(t - (i - 1)OSFT_{e} - 16OSFnT_{e} - lT_{e} - \tau(t)) + \frac{1}{Noise} \right] + \eta(t) + \eta(t)$$

where  $\eta(t)$  is an additive Gaussian noise,  $T_e$  is the sampling time ( $T_c$ =OSF. $T_e$ ),  $T_c$  is the chip duration, and  $\phi(t)$  is the global phase shift introduced by propagation channel and oscillators imperfections. In FDD the received signal contains the CPICH symbols, so we can separate the signal in two parts; one part contains CPICH

symbols and the other part contains the other transmitted channels, which allows us to write equation (4).

Equation (4) gives us the received signal model in baseband for UMTS-FDD, the calculus and demonstrations in the next sections are based on this model.

Copyright © 2008 SciRes.

(4)

#### 2.3. Signal Model in UMTS-TDD Mode

The UMTS TDD physical channel is a combination of two data fields, a midamble, and a guard period. There are two burst types proposed in [20], namely burst type 1 and 2. As illustrated in Figure 1, both types have the same length of 2560 chips and are terminated by a guard period of 96 chips in order to avoid overlapping with consecutive time slots. Burst type 1 has a longer midamble (512 chips) suitable for cases where long training periods are required for adaptation and tracking. The midamble code is used as signal pilot for tracking algorithm.



Figure 1. Two burst types

With the same notation as the previous sub-section B, the received signal in TDD modes will be written as described in (4). In equation (4),  $d^{(k,1)}$  denotes the data symbols transmitted before the midamble and  $d^{(k,2)}$  denotes the data symbols after the midamble. The combination of the user specific channelization and cell specific scrambling codes can be seen as a user and cell specific spreading code, with chip values, are denoted by  $S_i^k$  in (4).

#### 3. Digital Joint Synchronization Algorithm

#### 3.1. Introduction

The sampling instant of the received signal is fluctuated, and the phase is shifted because of the channel effects, the reception clock imperfections, and receivertransmitter oscillators' imperfections. The received signal needs synchronization to extract the transmitted sequences. So that synchronization is divided into two phases, acquisition phase and tracking phase. Acquisition is the first step that the mobile phone does, known as "cell search" in UMTS standard. It consists finding the beginning of signal (top frame and top slot) within a half chip precision. This operation allows also finding the cell's parameters. The tracking step updates the phase and the sampling instant continuously in time. In this section we will detail the algorithm used in tracking phase. The developed algorithm consists of four steps, as shown on Figure1, assuming that the baseband signal is obtained by radio frequency down conversion followed by a linear filtering system, and the acquisition is done correctly.

#### **3.2. Timing Correction**

We look to extract from the received samples the correct transmitted samples. To perform this operation the received signal will be interpolated at the estimated correct instant. The following formula presents the fundamental equation for digital interpolation [23]:

$$y(kT_i) = y((m_k + \mu_k)T_e) = \sum_{i=0}^N x[(m_k - i)T_e] \cdot h[(i + \mu_k)T_e]$$
$$kT_i = \left[\frac{kT_i}{T_e}\right]T_e + \mu_k T_e \Longrightarrow m_k = \left[\frac{kT_i}{T_e}\right]$$

where x is the received signal,  $T_e$  is the theoretical sampling period,  $T_i$  represents the fluctuated sampling period, and h(t) is the impulse response of interpolated filter. The interpolation coefficients are modified in function of the delay  $\mu_k$ .

In our work we estimate the delay time  $\mu_k$  then we adapt the interpolator coefficients. N is the number of interpolator coefficients, and [] denotes integer part function.

In literature there are many methods to interpolate digital signal: linear interpolation, polynomial interpolation, piecewise-parabolic interpolation and some others [23,24]. The choice of an interpolator depends on two criterions: precision and complexity. In our study we have chosen a truncated cardinal sine wave. In UMTS the over-sampling factor is reconfigurable, and to assure real time criteria the number of interpolators' coefficients depends on the over-sampling factor.

#### 3.3. Phase Correction

The second step is to correct the phase shift. The interpolated complex samples (I and Q) are multiplied by complex exponential of the estimate phase.

$$z(kT_i) = y(kT_i) e^{-j\phi_k}$$

where  $y(kT_i)$  is the interpolated complex samples and  $\hat{\varphi}_{i}$ 

is the estimated phase. For phase correction, in order to find the phase reference, the sign of the real and imaginary parts allows to find the reference region in the QPSK plan.

#### 3.4. Joint Estimation of Delay and Phase

Now we are arrived in one of the principal point in our work, the correction of the sampling instant and phase require the knowledge of time delay and phase shift. We use a procedure that allows us estimating jointly these parameters. These measurements are filtered by a loop filters, then they are used by interpolation and phase correction blocks. We use two different methods for estimation in the TDD and FDD modes. In the following

#### NOVEL JOINT CHIP SAMPLING AND PHASE SYNCHRONIZATION ALGORITHM FOR MULTISTANDARD UMTS SYSTEMS

sub-sections we explain these methods.

#### 3.4.1. Estimation for UMTS-FDD Mode

The sampling instant delay and the phase shift are estimated over one symbol duration,  $T_s=SF^*T_c$ , in order to correct the next symbol timing. The sampling instant is corrected by interpolation and the phase is corrected by complex multiplication between received signal and estimated phase.

We use as a pilot signal the corresponding scrambling code, because the CPICH channel is a copy of the scrambling code. So the received signal is divided into successive symbols (256 chips i.e. 256.OSF samples), and we calculate the correlation between each symbol and the original, early, and late copy of the corresponding symbol.

Let  $y_{r,s}$  be the s<sup>th</sup> symbol of received signal (described in equation (4)) and SC<sub>s</sub> represents the sequence of scrambling code, which was aligned in transmission with this part of the signal. The time delay and phase shift are supposed to be constant during one symbol period, then without correction of signal we obtain the following correlation function at instant  $\varepsilon$ :

$$\gamma(\varepsilon) = \gamma_{y_{r_{0,s}}, SC_s}(\varepsilon) + \alpha(s) \underbrace{e^{\phi(s)}}_{fading} \underbrace{G_{cpich}(1+j)}_{shift} \gamma_{SC_s}(\varepsilon - \tau(s)) + \underbrace{N(\varepsilon)}_{Noise} + \underbrace{N(\varepsilon)}_{Noise}$$
(6)

The spreading code is orthogonal, that means the correlation function of two codes is equal to zero. The first term is negligible. The scrambling code local auto-correlation function can be approximated by a triangular function:

$$\gamma_{sc}(\varepsilon) \approx 1 - \left| \frac{\varepsilon}{T_c} \right|$$
  
For  $\varepsilon = -\frac{T_c}{2}, 0, \frac{T_c}{2}$  we have:  
 $r(\varepsilon = 0) \approx e^{\varphi(s)} \alpha(s) G_{cpich}(1+j) r_{SC_r}(-\tau(s)) + N(\varepsilon = 0)$   
 $\approx e^{\varphi(s)} \alpha(s) G_{cpich}(1+j) (1 - \left| \frac{\tau(s)}{T_c} \right|) + N(\varepsilon = 0)$   
 $r(\varepsilon = -\frac{T_c}{2}) \approx e^{\varphi(s)} \alpha(s) G_{cpich}(1+j) r_{SC_r}(-\frac{T_c}{2} - \tau(s)) + N(-T_c/2)$   
 $\approx e^{\varphi(s)} \alpha(s) G_{cpich}(1+j) (1 - \left| \frac{1}{2} + \frac{\tau(s)}{T_c} \right|) + N(-T_c/2)$   
 $r(\varepsilon = \frac{T_c}{2}) \approx e^{\varphi(s)} \alpha(s) G_{cpich}(1+j) r_{SC_r}(\frac{T_c}{2} - \tau(s)) + N(T_c/2)$   
 $\approx e^{\varphi(s)} \alpha(s) G_{cpich}(1+j) (1 - \left| \frac{1}{2} - \frac{\tau(s)}{T_c} \right|) + N(-T_c/2)$   
 $(7)$ 

We observe that these measurements contain information about the time delay and the phase shift; our

idea is to extract this information from these measurements. Real and imaginary parts of one-time correlation function allow estimating the phase rotation  $\phi(s)$ . The early and late correlation functions give estimation of the time delay  $\tau(s)$ . Let define the measurements  $\Delta_{\tau}$  and  $\Delta_{\phi}$ , as the sampling instant delay and phase shift. And Re(.), Im(.) are the real and imaginary part function of a complex number.

$$\begin{split} \Delta_{\varphi}(s) &= \operatorname{Re}(r_{y_{r,s},SC_{s}}\left(\varepsilon=0\right)) - \operatorname{Im}(r_{y_{r,s},SC_{s}}\left(\varepsilon=0\right)) \\ &= 2\alpha(s)G_{cpich}\left(1 - \left|\frac{\tau(s)}{T_{c}}\right|\right) \sin[\varphi(s)] + \eta(\varphi,\tau) \\ \Delta_{1}(s) &= \operatorname{Re}(r_{y_{r,s},SC_{s}}\left(\varepsilon=\frac{T_{c}}{2}\right)) + \operatorname{Im}(r_{y_{r,s},SC_{s}}\left(\varepsilon=\frac{T_{c}}{2}\right)) \\ &= 2\alpha(s)G_{cpich}\left(1 - \left|\frac{1}{2} - \frac{\tau(s)}{T_{c}}\right|\right) \cos[\varphi(s)] + \eta_{1}(\varphi,\tau) \\ \Delta_{2}(s) &= \operatorname{Re}(r_{y_{r,s},SC_{s}}\left(\varepsilon=-\frac{T_{c}}{2}\right)) + \operatorname{Im}(r_{y_{r,s},SC_{s}}\left(\varepsilon=-\frac{T_{c}}{2}\right)) \\ &= 2\alpha(s)G_{cpich}\left(1 - \left|\frac{1}{2} - \frac{\tau(s)}{T_{c}}\right|\right) \cos[\varphi(s)] + \eta_{2}(\varphi,\tau) \\ &= 2\alpha(s)G_{cpich}\left(1 - \left|\frac{\tau(s)}{T_{c}}\right|\right) \sin\varphi(s) + \eta_{1}(\varphi,\tau) \\ \Rightarrow \begin{cases} \Delta_{\varphi}(s) &\equiv 2\alpha(s)G_{cpich}\left(1 - \left|\frac{\tau(s)}{T_{c}}\right|\right) \sin\varphi(s) + \eta_{1}(\varphi,\tau) \\ &\equiv 2\alpha(s)G_{cpich}\left(\left|\frac{1}{2} + \frac{\tau(s)}{T_{c}}\right| - \left|\frac{1}{2} - \frac{\tau(s)}{T_{c}}\right|\right) \cos[\varphi(s)] + \eta_{2}(\varphi,\tau) \end{cases}$$

Let  $\hat{\tau}(s)$  be the estimation of the sampling instant delay and  $\hat{\varphi}(s)$  is the estimation of phase shift, with a correction loop. Then the previous equations become:

$$\begin{cases} \Delta_{\varphi}(s) \cong 2\alpha(s)G_{cpich}\left(1 - \left|\frac{\hat{\tau}(s) - \tau(s)}{T_c}\right|\right) \sin[\varphi(s) - \hat{\varphi}(s)] \\ + \eta_1(\varphi, \tau) \\ \Delta_x(s) \cong 2\alpha(s)G_{cpich}\left(\left|\frac{1}{2} + \frac{\hat{\tau}(s) - \tau(s)}{T_c}\right| - \left|\frac{1}{2} - \frac{\hat{\tau}(s) - \tau(s)}{T_c}\right|\right) \\ \cos[\varphi(s) - \hat{\varphi}(s)] + \eta_2(\varphi, \tau) \end{cases}$$

$$\tag{9}$$

If the system of estimation and correction converge correctly, the time delay and phase error must converge to zero. In this condition:

$$\sin[\varphi(s) - \hat{\varphi}(s)] \approx \varphi(s) - \hat{\varphi}(s)$$
$$\cos(\varphi(s) - \hat{\varphi}(s)) \approx 1$$

With this condition, equation (8) gives the estimation of delay and phase as follows:

$$\Delta_{\varphi}(s) \cong 2\alpha(s)G_{cpich}\left[\varphi(s) - \hat{\varphi}(s)\right] + \eta_{1}(\varphi,\tau)$$
$$\Delta_{\tau}(s) \cong \frac{4\alpha(s)G_{cpich}}{OSF}\frac{\hat{\tau}(s) - \tau(s)}{T_{e}} + \eta_{2}(\varphi,\tau)$$
(10)

Copyright © 2008 SciRes.

where OSF is the over sampling factor, i.e. the number of samples per chip,  $T_c$  is the chip duration and  $T_e$  is the sampling period. Hence,  $T_c$ =OSF.T<sub>e</sub>.  $G_{cpich}$  is the channel gain of CPICH [17,18], and  $\alpha(s)$  is the fading of the propagation channel. Always, we suppose that the fading is invariant for symbol duration ( $T_s$ =SF.T<sub>c</sub>=SF.OSF.T<sub>e</sub>). These measurements are filtered by a second order filter that we discuss it in the next section. Now we will establish these measurements in the case of TDD mode.

#### 3.4.2. Estimation for UMTS-TDD Mode

We are going here to explain how we estimate sampling instant delay and phase shift in TDD mode. In this mode we process slot by slot the 2560 chips i.e. 2560.OSF samples. This is different from FDD mode where we process symbol by symbol. We have chosen this way because in the TDD mode we don't have a permanent channel with a known spreading code. In addition the TDD mode will be used when we have a good propagation channel conditions. In this case we can consider that the delay time and the phase shift are constant over one slot duration. For these raisons we use midamble code as a signal pilot to achieve synchronisation. So we calculate correlation function between the received midamble and the transmitted midamble in three points (codes aligned, the midamble code shifted to right and to left by OSF/2 samples). These measurements are filtered to estimate time delay and phase shift. Estimation of the sampling instant delay and the phase shift over actual slot is used to correct the next slot. The sampling instant is corrected by interpolation and the phase is corrected by complex multiplication between received signal and complex exponential of the estimated phase.

With the same notations and calculus as in the previous sub-section concerning FDD, we can derive the measurements of phase and delay errors as follows:

$$\Delta_{\varphi}(s) \cong \alpha(s)(\varphi(s) - \hat{\varphi}(s)) + \eta_1(\varphi, \tau)$$

$$\Delta_{\tau}(s) \cong \frac{2\alpha(s)}{OSF} \frac{\hat{\tau}(s) - \tau(s)}{T_e} + \eta_2(\varphi, \tau)$$
(11)

These parameters are filtered by a second order filters in order to estimate phase and delay. The structure of these filters is the same as in FDD mode. But filters' coefficients are different. This difference comes from the difference in the gain, sampling frequency and updating time period. We study in the following section these loop filters.

#### 3.5. Closed Loop Estimation

A Second-order loops should be able to track delay and phase offset as long as the values of the coefficients is carefully chosen.

So, the  $\Delta_{\tau}$  can be filtered by a second order filter with

 $\lambda_1$  and  $\lambda_2$  as coefficients. The output of the filter is used to determine the interpolator impulse response to find the correct sample. In FDD mode the interpolator response is updated every symbol (256 chips), but in TDD mode is updated every slot (2560 chips). The interpolation filter is updated for each symbol or slot (s) which is used to correct the sampling instant error for the next symbol or slot (s+1). Here we describe how a time loop filter updates time delay:

$$\mu_{s+1} = \mu_s + \varepsilon_k + \lambda_1 \Delta_\tau(s)$$

$$\varepsilon_{s+1} = \varepsilon_s + \lambda_2 \Delta_\tau(s)$$
(12)

where  $\mu$  is the estimated instant to be used by interpolator and  $\epsilon$  is the error of estimation.

With a similar structure the measurement  $\Delta_{\phi}$  is filtered by a second-order filter with two coefficients  $\gamma_1$ and  $\gamma_2$ . The output of the filter is used to correct the phase by complex multiplication. The phase is updated every symbol in FDD mode and every slot in TDD mode. So the calculated phase in symbol or slot (s) is used to correct the phase of the next symbol or slot (s+1). The following equation shows how a phase loop filter updates the phase:

$$\phi_{s+1} = \phi_s + \delta_s + \gamma_1 \Delta_{\varphi}(s)$$

$$\delta_{s+1} = \delta_s + \gamma_2 \Delta_{\varphi}(s)$$
(13)

where  $\phi$  is the estimated instant to be used by phase corrector and  $\delta$  is the estimation error.

From equations (10) and (12), we can derive the expression of sampling instant estimation:

$$\mu_{s+1} = \mu_s + F_{\tau}(z) \left( A_{\tau} \alpha(s) \frac{\tau(s) - \hat{\tau}(s)}{T_e} + \eta(s) \right)$$

where  $A_{\tau}$  depends on the mode (TDD or FDD):

$$A_{\tau} = \begin{cases} \frac{2}{OSF} & for \ TDD \ mode \\ \frac{4G_{CPICH}}{OSF} & for \ FDD \ mode \end{cases}$$

and  $F_{\tau}(z)$  is the loop filter transfer function:

$$F_{\tau}(z) = \lambda_{1} + \frac{\lambda_{2} z^{-1}}{1 - z^{-1}}$$

For the phase estimation, we find the same equation, but filter coefficients are different:

$$\phi_{s+1} = \phi_s + F_{\varphi}(z) \left( A_{\varphi} \alpha(s) \frac{\varphi(s) - \hat{\varphi}(s)}{T_e} + \eta(s) \right)$$

again  $A_{\varphi}$  depends on the mode (TDD or FDD) :

$$A_{\varphi} = \begin{cases} 1 & for \ TDD \ mode \\ \\ 2G_{CPICH} & for \ FDD \ mode \end{cases}$$

and  $F_{\varphi}(z)$  is the loop filter transfer function:

$$F_{\varphi}(z) = \delta_1 + \frac{\delta_2 z^{-1}}{1 - z^{-1}}$$

Estimator properties must be studied to determine filter coefficients. We observe that the time loop filter and phase loop filter have the same structure; in addition they have the same structure for both of modes FDD and TDD. So, we are studying the closed loop estimation in the general case using this filter structure. In what follows, we illustrate estimator properties in order to demonstrate how filters' coefficients are determined.

Figure 2 shows the diagram of a closed loop. We denote by m the parameter that we want to estimate and correct.



Figure 2. Estimation loop

The open loop transfer function of this system can be written as:

$$H_m(z) = \frac{A_{m1}z + (A_{m2} - A_{m1})}{z^2 + (A_{m1} - 2)z + (A_{m2} - A_{m1} + 1)}$$

where  $A_{m1}$  and  $A_{m2}$  are defined as :

$$A_{m1} = A_m \beta_1$$
$$A_{m2} = A_m \beta_2$$

We calculate the coefficients  $A_{m1}$  and  $A_{m2}$  for the case of a stable system. In order to have two complex conjugate poles, we should have  $(A_{m1}^2 - 4A_{m2}) < 0$ , and

$$\begin{cases} A_{m1} < 2\sqrt{A_{m2}} \\ A_{m2} > 0 \end{cases}$$
(14)

If the first condition is verified then the poles of  $H_m(z)$  are:

$$\begin{cases} p_m = \frac{2 - A_{m1} + j\sqrt{-A_{m1}^2 + 4A_{m2}}}{2} \\ p_m^* = \frac{2 - A_{m1} - j\sqrt{-A_{m1}^2 + 4A_{m2}}}{2} \end{cases}$$
(15)

We have a second order system; this system can be stable if the modules of its poles are smaller than one. These conditions imply that

$$(2 - A_{m1})^2 - A_{m1}^2 + 4A_{m2} < 4 \Longrightarrow A_{m2} - A_{m1} < 0$$

and with the first condition the criteria of stability becomes:

$$\begin{cases} A_{m2} < A_{m1} < 2\sqrt{A_{m2}} \\ 0 < A_{m2} < 4 \end{cases}$$
(16)

In order to achieve good estimation, we must choose filters' coefficients which minimize the estimation error (the residual error after convergence). With this criterion we can calculate the appropriate coefficients. From the structure of the closed loop estimation (Figure 4), the variance of the estimation error can be calculated as follows:

$$\operatorname{var}(m-\hat{m}) = \operatorname{var}(\frac{1}{A}h_m(t)*\eta(t))$$

where  $h_m(t)$  is the impulse response of the open loop filter, and  $\eta(t)$  denotes an additive Gaussian noise. Let  $\eta'(t) = h_m(t) * \eta(t)$  defines the filtered noise and N<sub>0</sub> is the power spectral density of the noise  $\eta(t)$ . So the power spectral density of the noise  $\eta'(t)$  is:

$$\Gamma_{\eta'}(f) = \frac{N_0}{A^2} \left| H_m \left( z = e^{2j\pi f} \right) \right|^2$$

The variance of the noise  $\eta'(t)$  is relied to the power spectral density by the following relation:

$$\operatorname{var}(\eta'(t)) = \int_{R} \Gamma_{\eta'}(f) \, df = \int_{R} \frac{N_{0}}{A^{2}} \left| H_{m}(z = e^{2j\pi f}) \right|^{2} \, df$$

By proceeding to a variable change, the integral can be calculated on the unit circle in Z plane.

$$\operatorname{var}(\varepsilon) = \frac{N_0 A_{m1}^2}{2j\pi A^2} \oint \frac{z - z_m}{(z - p_m)(z - p_m^*)} \frac{1 - z_m z}{(1 - p_m z)(1 - p_m^* z)} dz$$
(17)

where  $z_m$  is the zero and  $p_m$  is the pole the transfer function. The function to be integrated have four poles,

Copyright © 2008 SciRes.

I. J. Communications, Network and System Sciences, 2008, 2, 105-206

two of them are out of the unit circle in Z plane, so by using the residue theorem we can calculate this integral then we obtain:

$$\operatorname{var}(\varepsilon) = \frac{N_0 A_m^2}{A^2} \frac{1}{\left|1 - p_m^2\right|^2 (1 - \left|p_m\right|^2)} \left[ \left(1 + \left|p_m\right|^2\right) \left(1 + z_m^2\right) - 4 \Re e(p_m) z_m \right]$$
(18)

To find filters' coefficients with efficient estimator, the variance of the estimator must be optimized. There is a region of coefficients that allows to an accepted variance relative to Cramer-Rao bound. Also in simulation we have found that for these values the variance of estimation error is minimal. From these values we can calculate the different filters' coefficients.

In the FDD mode, and for the time loop filter (TLF) we fix the coefficient as follows:

$$A_{\tau} = \frac{2}{OSF}, \ \lambda_1 = \frac{0.15 \bullet OSF}{G_{CPICH}} \text{ and } \ \lambda_2 = \frac{0.05 \bullet OSF}{G_{CPICH}}$$

and for the phase loop filter (PLF):

$$A_{\varphi} = 2G_{CPICH}, \ \delta_1 = \frac{0.3}{G_{CPICH}} \text{ and } \delta_2 = \frac{0.1}{G_{CPICH}}$$

If we don't know the gain of the channel CPICH (which is less than one) we fix it to one. Even for this case we are sure that our system is stable and the estimation error stay minimal.

In the TDD mode, for the time loop filter (TLF) we fix the coefficients to the following values:

$$A_r = \frac{2}{OSF}, \ \lambda_1 = 0.3OSF \text{ and } \lambda_2 = 0.1 \bullet OSF$$

and for the phase loop filter (PLF):

$$A_{\alpha} = 1, \ \delta_1 = 0.6 \ \text{and} \ \delta_2 = 0.2$$

#### 4. Simulation Platform

In order to test the proposed algorithm and to determine all necessary parameters for implementation, we have carried out an environment of simulation. The UMTS base band signal in different modes is defined as well as the propagation channel.

Figure 3 shows the modeling of the propagation channel, that we consider in simulation context. For the Fading phenomena we use a classic Jake's spectrum noised with a complex Gaussian noise to model the Doppler's effect and the frequency selectivity. The Doppler frequency is considered as variable. The simulated noise is an AWGN. The variance of this noise determines the signal to noise ratio (SNR). In order to introduce delay time in transmitted signal, the samples are delayed by interpolation. We have simulated different sampling drift functions; linear function and triangular function. To shift the phase, the transmitted signal was multiplied by complex exponential of a variable phase. Many forms of shift functions are experimented to evaluate the performance of our algorithm.



Figure 3. Propagation channel modeling

We try here to show the performance of our algorithm. So we show the obtained results for many configurations. In the first configuration, we considered the FDD mode, and three physical channels, i.e. three spreading codes, to be transmitted. Two channels contained random data and the third is the particular channel CPICH. The spreading factors are SF1=SF2=16 and SF of CPICH is 256. To introduce time delay in the transmitted signal we use a drift with 100 ppm  $(10^{-4} T_e \text{ of drift per sample})$ duration). And we shift the phase linearly by a coefficient of 194 radian/s. and we fix Doppler frequency to 220 Hz i.e. 120 km/h. and one propagation path is considered. The SNR is variable. The following Figures 4 and 5 show both real and imaginary parts of the received CPICH channel for two values of SNR (SNR=5 dB and 15 dB), where the transmitted symbol value of this channel is (1+j). We obtain this figure by correlation between the corrected received signals and scrambling code as it described in the section "estimation in FDD mode". So these results show the stability, convergence and efficiency of our algorithm.

10 frames are simulated, which correspond to 1500 symbols of CPICH channel. The real and imaginary parts converge to 1. The transient phase takes about 10 symbols duration, so the convergence is established after ten to fifteen symbols. That means the synchronization is quickly established and maintained. We observe that the variance of synchronized symbol depends on the SNR.

Figure 6 shows the constellation of CPICH for the same values of SNR.

As it has described before our algorithm is based on joint estimation of time delay and phase shift. So to study our estimator we search for unbiased and optimal estimators. For the first criteria we observe the convergence of estimation and for the second we look to minimize variance.

#### NOVEL JOINT CHIP SAMPLING AND PHASE SYNCHRONIZATION ALGORITHM FOR MULTISTANDARD UMTS SYSTEMS



Figure 4. Real part of CPICH channel after correction



Figure 5. Real part of CPICH channel after correction



Figure 6. CPICH constellation after synchronization

In the following figures we show the error of delay and phase estimation for two different SNR values. Figure 7 shows delay estimation error and Figure 8 shows phase estimation error.



Figure 8. Phase estimation error

We observe that the mean of our estimation error is equal to zero, which allows us to conclude that our estimator is not biased. To show estimator's efficiency we have observed the variances of estimators in function of SNR.

To demonstrate the functionality and the performance of our algorithm, we have transmitted two physical channels of data, and after synchronization symbols are extracted by de-spreading, which allow comparing with the original transmitted symbols. In the first figure we show the constellation of data. We have repeated this operation for many different SNR and then bit error rates are calculated in function of signal to noise ratio.



Figure 9. Variance of delay time estimation in function of SNR



Figure 10. Variance of phase estimation



Figure 11. Constellation of decoding data after synchronization



Figure 12. BER in function of SNR

Copyright © 2008 SciRes.

Figure 11 shows QPSK data constellation after synchronization. Figure 12 shows bit error rate in function of SNR.

The TDD mode is simulated in similar conditions and we have obtained the same performance as in FDD mode. Even for low SNR=15, the symbol error rate is acceptable.

# 5. Implementation in DSP TigerSharc TS101

After validation of the approach by simulation, the proposed algorithm is then implemented in a DSP TigerSharc TS101 as well as a complete UMTS transmission-reception chain is carried out. This section discusses the strategy of implementation, the complexity, and the real time processing [25,26].

Algorithm functionality is divided into 4 blocks as it described in Figure 1; interpolation, phase correction, correlation and estimation.

The interpolation is based on a look up table (LUT), which contains the interpolator filter coefficients. The number of coefficients depends on the over sampling factor while the total number of stored coefficients depends on the processing precision. For phase correction, samples of the trigonometry functions (SIN and COS) are stored in a look up table. The size of these tables depends on the precision of correction.

The scrambling code in FDD mode and the midamble code in TDD mode are determined by the acquisition block and stocked in memory to be used for tracking, channel estimation, and equalization. These parameters depend on the cell where the mobile is localized. To assure real time processing, received data are processed frame by frame; the size of each frame is one slot. In the next section we give some details about implementation and the obtained results.

#### 5.1. Interpolation Process in DSP

After estimation of the time delay  $\mu_k$  in slot k, the next slot is corrected by interpolation. The used interpolator is given by its impulse response function:

$$h_k(n) = \frac{\sin \pi (n + \mu_k)}{\pi (n + \mu_k)}$$

where  $-1 < \mu_k < 1$ , -N <= n <= N and 2N+1 is the number of coefficients. This will accelerate the data fetch functions. In order to fill the LUT we have calculated the coefficients of interpolators.

As the memory space is limited, it is not possible to provide coefficients for each time delay. For this reason estimated time delay is quantified with a step  $\delta$ . The choice of the step of quantification is based on the required precision. In our case, the samples of data are

#### NOVEL JOINT CHIP SAMPLING AND PHASE SYNCHRONIZATION ALGORITHM FOR MULTISTANDARD UMTS SYSTEMS

coded with 12 bits. So the step size must be more than  $2^{-11}$ . Let  $\overline{\mu}_k$  denotes the quantification value of  $\mu_k$  so we have:

$$\forall \mu_k, -1 < \mu_k < -1 \ \exists m, m\delta < \mu_k < (m+1)\delta \Rightarrow \overline{\mu}_k = m\delta + \frac{\delta}{2}$$

Then the impulse response function of the interpolator k can be written as:

$$h_k(n) = \frac{\sin \pi (n + m\delta + \frac{\delta}{2})}{\pi (n + m\delta + \frac{\delta}{2})}$$

The data are coded with fixed point with fractional format (Q<sub>12</sub>: 12 bits are used). So it is not useful to choose precision for quantification of time delay more than the precision used in data, because the time delay is estimated by correlation of this data. So that time delay can be quantified with a step of  $\delta = 10^{-7}$ . The LUT for interpolation is organized as described in Figure 13. The size of this table is:

$$S = \left(2N+1\right)\frac{2}{\delta}$$

where 2N+1 is the number of interpolation coefficients, this number depends on the oversampling factor (OSF). For OSF=2 we need 5 coefficients to assure good precision and for OSF=4 or 8 we use only 3 coefficients. The number of coefficients corresponds to the number of samples used to extract the correct sample.



Figure 13. Interpolation coefficients LUT

#### 5.2. Phase Correction & Processing in DSP

To store samples of trigonometric functions, complex number representation is used, and the phase is quantified with a step of:

$$\delta = \frac{\pi}{512}$$

The following table shows the arranged functions in the DSP memory:



Figure 14. Trigonometric functions LUT

The size of this table is 1024 words of 16 bits. SIN and COS tables are stored in the same 32-bit word, where the SIN samples are stocked in the more signified 16-bit (MSB) and COS samples in the last significant 16-bit (LSB). This method allows the DSP to fetch and execute complex multiplication in one clock cycle. So we get benefit from the wide memory DSP structure.

#### 5.3. Complexity of Tracking Algorithm

We have evaluated the number of cycles required for every function and the total number of cycles required for tracking. The criterion is to assure real time. We start by giving all parameters, so we process slot by slot. The time duration of one slot is  $T_{slot}=666.667 \ \mu s$ . the DSP clock period is  $T_{cycle}=4$  ns. We detail here the required number of cycles, which depends on the mode used (TDD or FDD) and the over sampling factor. Table 1 shows that for all values of OSF and for both modes, the processing in real time is assured. The functions that need an important number of cycles are interpolation and correlation.

#### 5.4. Implementation Results

For the same simulation conditions, the algorithm is evaluated in a real experimental platform. The obtained results demonstrate the efficiency of the proposed tracking algorithm. For a signal to noise ratio SNR=5dB the symbol error rate is lower than 0.1%.

#### Y. SERRESTOU ET AL.

#### Table 1. Implementation complexity

		Interpolation	Phase correction	Correlation	Estimation	Total cycles	Time
FDD	OSF=2	51420 cycles	10480 cycles	31030 cycles	147 cycles	93077 cycles	372.31 µs
mode	OSF=4	51460 cycles	20740 cycles	31030 cycles	169 cycles	103399 cycles	413.60 µs
	OSF=8	102620 cycles	41140 cycles	15660 cycles	169 cycles	159589 cycles	638.36 µs
TDD	OSF=2	51224 cycles	<b>10265</b> cycles	618 cycles	132 cycles	67807 cycles	271.23 µs
mode	OSF=4	51226 cycles	20505 cycles	15390 cycles	132 cycles	87253 cycles	349.01 µs
	OSF=8	102450 cycles	41010 cycles	15390 cycles	132 cycles	158982 cycles	635.93 µS







Figure 16. CPICH constellation after synchronization

# NOVEL JOINT CHIP SAMPLING AND PHASE SYNCHRONIZATION ALGORITHM FOR MULTISTANDARD UMTS SYSTEMS



Figure 17. Constellation of decoded data after synchronization for SNR=5dB and 15dB

For this algorithm, the acquisition step is required only at the beginning of the transmission, or when the system is out-off synchronization. When the acquisition is achieved with the required precision at the beginning, tracking algorithm maintains continuously the synchronization of the received signal. This synchronization is maintained even under critical propagation conditions. Experiments shows that for time drift of 100ppm, i.e. of  $10^{-4}$  T<sub>e</sub> of drift per sample duration, and phase drift of 194 rad/s and at SNR = 5 dB, tracking algorithm maintains synchronization and corrects the received signal with a high precision, so the BER is lower than 0.1%. In addition, during simulation we have considered just the case of propagation channel with single path, but in experimentation the multipath channel were taken into consideration. So this algorithm tracks the principal path, which has the higher power at reception. The implementation allows to change configuration, so we can change the mode (FDD or TDD) as well as the over sampling factor (OSF) without loss of synchronization.

#### 6. Conclusions

In this paper, digital joint chip timing and phase tracking for reconfigurable standard UMTS was introduced. The mathematical model of this digital tracking algorithm has been presented and used to determine the optimal parameters for implementation. The algorithm was validated by simulation and implemented in 16-bit fixed point TigerSharc DSP. Simulation was carried out in many configurations in order to confirm theoretical parameters and to evaluate the performances and robustness of this algorithm. Experiences demonstrate the feasibility of real implementation within reconfigurable UMTS communication receiver. It shows high performance and ability to assure real time constraints and reconfiguration ability. We note that this algorithm can be adapted for other applications that use different techniques for multiple access schemes, such as chaotic communication systems. We think that our algorithm can be easily adapted for such application, and can show more performance in terms of bit error rate.

Our future work will be focalised on joint synchronization for MIMO systems.

#### 7. References

- R. Prasad and T. Ojanperä, "An overview of CDMA evolution toward wideband CDMA," IEEE Commun. Lett, vol. 1, no. 1, pp. 2–29, January 1998.
- [2] T. Samanchuen and S. Tantaratana, "symbol synchronization for MC-CDMA using timing estimator and delay locked loop," Proceeding of IEEE/ ISCIT'2005, pp. 366–369.
- [3] L. Tomba and W.A. Krzymien, "Sensitivity of the MCCDMA Access Scheme to Carrier Phase Noise and Frequency Offset," IEEE Trans. Veh. Technol., vol. 48, no.5, pp. 1657–1665, September 1999.
- [4] R.A. Iltis, "A DS-CDMA tracking mode receiver with joint channel/delay estimation and MMSE detection," IEEE Trans. Communications, vol. 49, no. 10, pp. 1770– 1779, 2001.
- [5] M. Latva-aho and J. Lilleberg, "Parallel interference cancellation based delay tracker for CDMA receivers," in Proceeding of the 30<sup>th</sup> Conference on Information Sciences and Systems.
- [6] W. Zha and S.D. Blostein, "Multiuser Delay-Tracking CDMA Receiver," EURASIP Journal on Applied Signal Processing, pp. 1355–1364, December 2002.
- [7] B. Jovic et al., "A robust sequence synchronization unit

for multi-user DS-CDMA chaos-based communication systems," in Signal Processing, vol. 87, pp. 1692–1708, 2007.

- [8] R.L. Peterson, R.E. Ziemer, and D.E. Borth, "Introduction to Spread Spectrum Communications," Prentice Hall, Inc., New Jersey, pp. 149–318 (Chapters 4 and 5), 1995.
- [9] R.E. Ziemer and R.L. Peterson, "Introduction to Digital Communication," second edition, Prentice Hall, Inc., New Jersey, pp. 611–615 (Chapter 9), 2001.
- [10] J.S. Lee and L.E. Miller, "CDMA Systems Engineering Handbook," Artech House Publishers, Boston, pp. 774– 837 (Chapter 7), 1998.
- [11] S. Parkvall, E.G. Ström, and B. Ottersten, "The impact of timing errors on the performance of linear DS-CDMA receivers," IEEE Journal on Selected Areas in Communications, vol. 14, no. 8, pp. 1660–1668, 1996.
- [12] J.C. Lin, "Low-complexity code tracking loop with chiplevel differential detection for DS/SS receivers," IEEE Electronic Letters, vol. 36, no. 24, 23rd November 2000.
- [13] E. G. Strom, S. Parkvall, S. L. Miller, and B. Ottersten, "Propagation delay estimation in asynchronous directsequence code-division multiple access systems," IEEE Trans. Communications, vol. 44, no. 1, pp. 84–93, 1996.
- [14] Y. Serrestou, K. Raoof, J. Liénard, "Digital joint sampling instant and phase synchronization for UMTS standard," IEEE/ACES International Conference on Wireless Communications and Applied Computational Electromagnetics PN: 905565, Hawaii, 2005.
- [15] F.M. Gardner, "Interpolation in digital Modems," IEEE Trans. on comm., vol.41, no.3, March 1999.
- [16] 3GPP TS 25.211, "Physical layer-General description," Technical Specifications of the 3rd Generation Partnership Project. http://www.3gpp.org/ftp/Specs
- [17] 3GPP TS 25.211, "Physical channels and mapping of transport channels onto physical channel FDD,"

Technical Specifications of the 3rd Generation Partnership Project. http://www.3gpp.org/ftp/Specs

- [18] 3GPP TS 25.212, "Multiplexing and channel coding FDD," Technical Specifications of the 3rd Generation Partnership Project. http://www.3gpp.org/ftp/Specs
- [19] 3GPP TS 25.221, "Physical channels and mapping of transport channels onto physical channel TDD," Technical Specifications of the 3rd Generation Partnership Project. http://www.3gpp.org/ftp/Specs
- [20] 3GPP TS 25.222, "Multiplexing and channel coding TDD," Technical Specifications of the 3rd Generation Partnership Project. http://www.3gpp.org/ftp/Specs
- [21] L. Krikidis, J.L. Danger, and L. Naviner, "A DS-CDMA multi-stage inter-path interference canceller for high bit rates," 2004 IEEE Eighth International Symposium on Spread Spectrum Techniques and Applications, pp. 405– 408, August 30 – September 2, 2004.
- [22] L. Krikidis, J.L. Danger, and L. Naviner, "Flexible and reconfigurable receiver architecture for WCDMA systems with low spreading factors," Electronics Letters, vol. 41, no. 1, pp 22–24, January 6, 2005.
- [23] F.M. Gardner, "Interpolation in digital Modems," IEEE Trans. on comm., vol. 41, no. 3, March 1999.
- [24] R.E. Crochiere, "Optimum FIR digital filter implementations for decimation, interpolation, and narrow-band filtering," IEEE Trans. on acoustics speech and signal processing, vol. ASSP-23, no. 5, October 1975.
- [25] A. Devices, "Tuning C source code for the TigerSharc DSP compiler," Engineer to engineer note, EE-147, January 2001.
- [26] A. Devices, "16-bit FIR filters ADSP-TS20x TigerSharc processors," Engineer to engineer note, EE-211, January 2004.



## Performance Analysis of an AMC System with an Iterative V-BLAST Decoding Algorithm

Sangjin RYOO<sup>1</sup>, Kyunghwan LEE<sup>2</sup>, Intae HWANG<sup>3</sup>

 <sup>1</sup> Department of Electronics Engineering, Chonnam National University, Korea
 <sup>2</sup> T&M Algorithm Research, Innowireless, Inc., Korea
 <sup>3</sup> Department of Electronics & Computer Engineering, Chonnam National University, Korea E-mail: <sup>1</sup>sjryoo@empal.com, <sup>2</sup>signalds@innowireless.co.kr, <sup>3</sup>hit@chonnam.ac.kr

#### Abstract

In this paper, the iterative Vertical-Bell-lab Layered Space-Time (V-BLAST) decoding algorithm of an Adaptive Modulation and Coding (AMC) system is proposed, and the corresponding MIMO scheme is analyzed. The proposed decoding algorithm adopts iteratively extrinsic information from a Maximum A Posteriori (MAP) decoder as an a priori probability in the two decoding procedures of the V-BLAST scheme of ordering and slicing in an AMC system. Furthermore, the performance of the proposed decoding algorithm is compared with that of a conventional V-BLAST decoding algorithm and a Maximum Likelihood (ML) decoding algorithm in the combined system of an AMC scheme and a V-BLAST scheme. In this analysis, each MIMO schemes are assumed to be parts of the system for performance improvement.

Keywords: Iterative V-BLAST Decoding, MAP Decoder, AMC, STD, Turbo Code

#### 1. Introduction

To improve the throughput performance together with the development of the MIMO scheme, the AMC scheme has attracted considerable attention as the forerunner of next-generation mobile communication systems [1]. The AMC scheme adapts a coding rate and modulation scheme to the channel condition [2], resulting in improved throughput performance. Consequently, the combination of a MIMO scheme and an AMC scheme can potentially improve the throughput performance. V-BLAST [3,4] was selected as the MIMO multiplexing scheme [5] and the turbo-coding [6] was chosen as the channel coding scheme of the AMC due to the complexity of the aforementioned combined system. The turbo-coding scheme with iterative decoding implies the use of parallel concatenated recursive systematic convolutional codes. Such a scheme is iteratively decoded with a Posteriori Probabilities (APP) algorithms for the constituent codes [7,8]. In addition, the turbo decoding algorithms used with MIMO is currently an area that is actively researched [9,10].

A performance analysis is offered here of AMC systems with several V-BLAST decoding algorithms

including the turbo decoding algorithm used with MIMO. For greater performance improvement, the proposed system utilizes a  $2\times 2$  MIMO channel using two transmitter antennas and two receiver antennas, a  $4-2\times 2$  MIMO channel applying a Selection Transmit Diversity (STD) scheme [11] that selects two antennas from four transmitter antennas, a  $4\times 4$  MIMO channel using four transmitter antennas and four receiver antennas, and a  $8\times 8$  MIMO channel using eight transmitter antennas and eight receiver antennas.

#### 2. The AMC System with the Proposed Iterative V-BLAST Decoding Algorithm

Figure 1 shows the structure of the AMC system used with the proposed iterative V-BLAST decoding algorithm. An AMC system that uses a conventional V-BLAST decoding algorithm combines a V-BLAST scheme with a turbo-coded AMC system. The proposed decoding algorithm of an AMC system differs from the conventional V-BLAST decoding algorithm insofar as the extrinsic information from a MAP decoder is used as an a priori probability in the ordering and slicing decoding procedures of the V-BLAST scheme [12].



Figure 1. Transmitter-receiver structure of an AMC system with the proposed iterative V-BLAST decoding algorithm

This scheme operates iteratively and is defined as the main MAP iteration. Furthermore, whenever the scheme operates internally, iterative decoding of the MAP decoder is performed. This method is defined as sub MAP iteration. For the proposed system, a system equipped with M transmitter antennas and N receiver antennas is considered. It is assumed that each transmission channel is modeled as a flat Rayleigh fading channel. The received signal in the V-BLAST receiver is denoted by

$$X = Hs + n \tag{1}$$

where  $X = [x_1, ..., x_N]^T$  is the received signal vector,  $s = [s_1, \dots, s_M]^T$  is the transmitted symbol vector, H is the  $N \times M$  channel matrix,  $n = [n_1, \dots, n_N]^T$  is the noise vector. The superscript T signifies the transpose matrix, and the noise vector, n, is modeled as a complex Gaussian random process. In addition,  $s_m$  is the  $2^{Q}$ -ary modulated symbol; that is  $s_m = f(d_1^m, \dots, d_Q^m) \in \Phi = \{\varphi_1, \dots, \varphi_{2^o}\}$ , where Q denotes the bit number per symbol,  $f(\cdot)$  denotes the symbol modulation function,  $\{d_q^m\}_{q=1,...,Q}$  represents the q-th information bits that correspond to  $s_m$ , and  $\{\varphi_i\}_{i=1,\dots,2^{\mathcal{O}}}$  represents the *i*-th symbol. The proposed slicing algorithm does not make a hard decision with the received signal but makes a decision with the extrinsic information from the MAP decoder [13]. This extrinsic information from the MAP decoder is the log-likelihood function, which can be described as

$$L_{m,q} = \log \frac{p\left(d_q^m = 1\right)}{p\left(d_q^m = 0\right)}$$
(2)

where  $L_{m,q}$  is the extrinsic information that corresponds to  $d_q^m$  [14]. Specifically,  $\{d_q^m\}_{q=1,...,Q}$  is determined by  $\{L_{m,q}\}_{q=1,...,Q}$ , respectively. (e.g., if  $L_{m,q}$  is greater than 0,  $d_q^m$  is determined to be 1. Otherwise,  $d_q^m$  is determined to be 0.) The proposed slicing algorithm then performs the quantization operation appropriate to the constellation in use corresponding to  $\{d_q^m\}_{q=1,...,Q}$ . In a conventional V-BLAST ordering procedure, the decoding order is determined by the SNR of the corresponding layer. The conventional V-BLAST ordering is described as

$$l_k = \arg\min_{m} ||(H_k^{\dagger})_m||^2$$
(3)

where k denotes the decoding stage and the superscript  $\dagger$  represents the pseudo-inverse matrix. The SNR is a function of the channel power, and the layer with the largest channel power is the first layer that is decoded. A high SNR signifies a low symbol error rate. From this fact, it follows that the maximum SNR criterion can be considered to be a specific version of the minimum symbol error criterion. The proposed ordering algorithm is a function not only of the SNR but also of the extrinsic information. It can be modified accordingly to

$$l_k = \arg\min_m P_m(e|X_k, H_k, L_m^{(i)})$$
(4)

where  $P_m(e|X_k, H_k, L_m^{(i)})$  is the symbol error probability of the *m*-th layer and  $L_m^{(i)} = [L_{m,1}^{(i)}, \cdots, L_{m,Q}^{(i)}]^T$  is the extrinsic information vector of the  $l_k$ -th layer at the *i*-th main MAP iteration. The symbol error probability,  $P_m$ , can be calculated from

$$Pm(e|Xk, Hk, Lm(i))$$

$$= \frac{1}{2^{\varrho}} \sum_{q=1}^{2^{\varrho}} \sum_{p=l, p\neq q}^{2^{\varrho}} P(\varphi_q | L_m^{(i)}) P(\varphi_q \to \varphi_p | X_k, H_k, L_m^{(i)})$$
(5)

where  $\varphi_q$  is the original transmitted symbol,  $\varphi_p$  is the possible symbol excluding the original transmitted symbol ( $\varphi_q$ ), and  $P(\varphi_q \rightarrow \varphi_p | X_k, H_k, L_m^{(i)})$  is the pair-wise symbol error probability, which can be obtained from

$$P(\varphi_q \rightarrow \varphi_p | X_k, H_k, L_m^{(i)})$$
  
=  $P[p(\varphi_q | y_m) < p(\varphi_p | y_m)]$   
=  $P[log p(\varphi_q | y_m) < log p(\varphi_p | y_m)]$  (6)

where  $y_m$  is the desired symbol that deletes the interference of other symbols after the nulling process of the V-BLAST decoding in the received symbol of the *m*-th layer,  $x_m$ . With the assumption that the variance of noise corresponding to the *m*-th layer is  $\sigma_m^2/2$ , in Eq. (6), the log posteriori function of  $\varphi_p$  is described by

$$\log p(\varphi_p|\mathbf{y}_m) \tag{7}$$

$$= log [ p(\varphi_p | L_m^{(i)}) p(y_m | \varphi_p) / p(y_m) ]$$
  
= log p(\varphi\_p | L\_m^{(i)}) + [Re(\varphi\_p - \varphi\_q)(2y\_m - (\varphi\_p + \varphi\_q))^\*] / 2\sigma\_m^2

where the superscript \* signifies a complex conjugate.

#### 3. Simulation Results

#### **3.1. MCS Level and Parameters for Simulation**

Table 1 shows the Modulation and Coding Scheme

(MCS) level selection thresholds, and Table 2 shows the simulation parameters. The detailed parameters in Table 1 are established on the basis of the 1X EV-DO standards [15]. There are many references in the selection of the MCS level selection threshold. For example, the threshold can be selected to satisfy the required Bit Error Rate (BER) and the required Frame Error Rate (FER). As more emphasis is placed here on the data transmission rate, the threshold that maximizes the throughput performance was selected. That is, the threshold of the selected MCS level is derived from the MCS level transmission rate performance intersection in each system. One frame is set up with one transmission slot with a frame length of 2,048 symbols. If one bit error occurs in one frame, it is assumed to be a frame error. When a frame error does not occur, the transmission rate is calculated in accordance with the V-BLAST technique in the order of (bit length  $\times$  data rate  $\times$  number of transmit antenna). The performance of the transmission rate closely corresponds to the capacity of the FER. Thus, in accordance with the transmission rate, a performance analysis is obtained by the error probability.

Table 1. MCS levels

MCS level	Data rate (Kbps)	Number of bits per frame	Code rate	Modulation
1	614.4	1,024	1/3	QPSK
2	1,228.8	2,048	2/3	QPSK
3	1,843.2	3,072	2/3	8PSK
4	2,457.6	2,096	2/3	16QAM

**Table 2. Simulation parameters** 

Parameter	Value
Turbo-coding scheme	PCCC
MAP iteration of the AMC system with a conventional V-BLAST technique	4
Main MAP iteration of the AMC system with the proposed V-BLAST technique	4
Sub MAP iteration of the AMC system with the proposed V-BLAST technique	2
Channel	Flat fading

#### 3.2. Complexity of Each Decoding Algorithm

This section outlines the complexity of the proposed decoding algorithm, the conventional V-BLAST decoding algorithm, and the ML decoding algorithm in the combined system of an AMC scheme and a V-BLAST scheme. The multiplication operation contributes to the complexity of implementing the system. Except for the procedures of a channel deinterleaver and the MAP decoder in the receiver, each decoding algorithm was compared to the number of

multiplication operations, as shown in Table 3 [16]. In this table, C is the number of symbols, S is the number of sub MAP iterations, L is the number of main MAP iterations, and B is the number of bits per symbol. Some examples of the table show that the proposed decoding algorithm is more complex than a conventional V-BLAST decoding algorithm but less complex than an ML decoding algorithm. In particular, when used with a higher-order modulation, the proposed decoding algorithm is less complex than the ML decoding algorithm. According to the table, as the modulation changes from QPSK to 16QAM in the case of M=N=4, the computational complexity of the proposed decoding algorithm ranges from approximately 24% to 0.1% of the complexity of the ML decoding algorithm. Furthermore, comparing with the complexity of the proposed scheme in [16], the complexity of the proposed scheme is relatively less complex for M=N=4, QPSK, and L=3 or 4.

 Table 3. Complexity of each decoding algorithm

 (L=4, S=2, M=N=4)

	ML decoding	Conventional decoding	Proposed decoding
Required multiplications	C <sup>M</sup> (M+1)N	(M+1)N <sup>3</sup> + (3/2)M <sup>2</sup> N+ [(7/2)M-1]N-1	$(M+1)N^{3}+$ $L[M^{2}N(B+1)+$ (3M-1)N-1]
QPSK	5,120	467	1,260
8PSK	81,920	467	1,516
16QAM	1,310,720	467	1,772

# **3.3.** Performance of the AMC Systems with Several V-BLAST Decoding Algorithms

Figure 2 shows the throughputs of the AMC systems with several V-BLAST decoding algorithms in a  $2\times2$  MIMO scheme. It is clear that the proposed decoding algorithm achieves a better throughput performance compared to the conventional V-BLAST decoding algorithm over the entire SNR range. Additionally, the proposed decoding algorithm is close to the existing ML decoding algorithm in terms of the throughput performance.

Figure 3 shows the throughputs of the AMC systems with several V-BLAST decoding algorithms in a  $2\times 2$  and  $4-2\times 2$  MIMO channel. It is demonstrated that the systems in a  $4-2\times 2$  MIMO channel achieve superior throughput performance relative to the others. The systems in the  $4-2\times 2$  MIMO channel that utilize a STD scheme show improvements in the SNR through the selection diversity gain. These improvements lead to a reduced error rate and an increase in the probability of selecting the MCS level with a higher data rate. These systems therefore achieve greater throughput

performance compared to the other systems. In addition, the proposed decoding algorithm achieves superior throughput performance relative to the conventional V-BLAST decoding algorithm in  $4-2\times 2$  MIMO channel using a STD scheme. It can be inferred that the proposed decoding algorithm achieves this effect as well in conjunction with a STD and a MIMO diversity scheme.

Figure 4 shows the throughputs of the AMC systems with several V- BLAST decoding algorithms in a  $2\times2$ , a  $4\times4$ , and an  $8\times8$  MIMO scheme. The results show that the approximate maximum throughput improvement for these three MIMO schemes is 421 Kbps, 545 Kbps, and 880 Kbps, respectively. Accordingly, it can be inferred that the effect of the proposed decoding algorithm increases as the number of system antennas increases. Moreover, when each MIMO scheme is applied, the performance is enhanced significantly.



Figure 2. Throughputs of the AMC systems with several V-BLAST decoding algorithms in a 2×2 MIMO scheme



Figure 3. Throughputs of the AMC systems with several V-BLAST decoding algorithms in a  $2\times 2$  and  $4-2\times 2$  MIMO scheme



Figure 4. Throughputs of the AMC systems with several V-

BLAST decoding algorithms in a 2×2, 4×4, and 8×8 MIMO scheme

#### 4. Conclusions

In this paper, to improve the throughput performance in a downlink, AMC systems with several V-BLAST decoding algorithms were implemented and compared. It was found that the performance can be improved through application of the STD as a MIMO diversity scheme. Through the SNR improvement of the receiver of the systems that utilized a STD scheme, the error probability was decreased in the range of a relatively low SNR and, ultimately, the throughput performance was improved. The throughput performance can also be enhanced by increasing the number of antennas in the MIMO channel.

The proposed decoding algorithm achieves a better throughput performance than the conventional V-BLAST decoding algorithm over the entire SNR range. For the example of M=N=4 and QPSK, it was demonstrated that the proposed decoding algorithm has nearly 24% lower complexity than the existing ML decoding algorithm while it provides an approximate increase of 8.3% in capacity compared to the conventional V-BLAST decoding algorithm.

In addition, the simulation results show that the maximum throughput improvement in each MIMO channel is nearly 421 kbps (a 17.7% increase in capacity) for a  $2\times2$  MIMO, 545 kbps (an 8.3% increase in capacity) for a  $4\times4$  MIMO, and 880 kbps (a 5.5% increase in capacity) for an  $8\times8$  MIMO. Thus, the effect of the proposed decoding algorithm increases while the number of system antennas increases. Accordingly, if the MIMO schemes or the MIMO channel can be applied in each case for a higher throughput performance, the proposed decoding algorithm will then be a practical candidate for next-generation mobile communication systems.

#### 5. References

- A.J. Goldsmith and S.G. Chua, "Adaptive Coded Modulation for Fading Channels," IEEE Trans. on Comm., vol. 46, no. 5, pp. 595–602, May 1998.
- [2] A. Bhargave and R.J.P. de Fegueiredo, "A new MIMO detector for iterative decoding with multiple antenna systems," Military Communications Conference, MILCOM 2005. IEEE, vol. 3, pp. 1428–1432, October 2005.
- [3] G.J. Foschini, "Layered Space-Time Architecture for Wireless Communication in a Fading Environment When Using Multi-Element Antennas," Bell Labs Technical Journal, pp. 41–59, 1996.
- [4] A. Bhargave, R.J.P. de Figueiredo, and T. Eltoft, "A Detection Algorithm for the V-BLAST System," GLOBECOM' 01. IEEE, vol. 1, pp. 494–498, November

2001.

- [5] A. Alamouti, "A simple transmit diversity technique for wireless communications," IEEE JSA on Comm., vol. 16, pp. 1451–1458, October 1998.
- [6] S. Benedetto and G. Montorsi, "Unveling Turbo Codes: some results on parallel concatenated coding schemes," IEEE Trans. on Inform. Theory, vol. 42, pp. 409–429, March 1996.
- [7] Q.F. Chen, H.F. Wang, M. Chen, and S.X. Cheng, "An Improved Turbo-BLAST System for Quasi-static Channel," The 15th IEEE International Symposium on Personal, Indoor and Mobile Radio Communications, 2004. PIMRC 2004, vol. 3, no. 5–8, pp. 1588–1591, September 2004.
- [8] Y. Li, Y. Yang, and H.S. Yan, "Using Turbo Code in BLAST System," Proceedings of the 2003 International Conference on Neural Networks and Signal Processing, 2003, vol. 2, pp. 1477–1480, December 2003.
- [9] T. Matsumoto and R.S. Thoma, "Turbo Transceivers for MIMO Wireless Communications and Their Performance Verification via Multi-Dimensional Channel Sounding," IEICE Trans. Commun. vol. E88-B, no. 6, pp. 2239– 2251, June 2005.
- [10] S. Haykin, McMaster University, M. Sellathurai, Y.D. Jong, and T. Willink, "Turbo-MIMO for Wireless

Communications," IEEE Communications Magazine, pp. 48–53, October 2004.

- [11] M. Sandell, "Analytical analysis of transmit diversity in WCDMA on fading multipath channels," PIMRC99, vol.2, pp. 946–950, September 1999.
- [12] Z.W. Catherine, H. Sweatman, J.S. Thompson, B. Mulgrew, and P.M. Grant, "Comparison of Detection Algorithm including BLAST for Wireless Communication using Multiple Antennas," PIMRC'00, vol. 1, pp. 698–703, 2000.
- [13] A. Elkhazin, N. Plataniotis, and S. Pasupathy, "Reduced-Dimension MAP Turbo-BLAST Detection," IEEE Transactions on Communications, vol. 54, no. 1, pp. 108–118, January 2006.
- [14] R. Wang, H. Wang, C. Fan, X. Zhang, and D.C. Yang, "Research on Modified Structure of Turbo-Blast System," The 17th Annual IEEE International Symposium on Personal, Indoor and Mobile Radio Communications, PIMRC'06, pp. 1–5, September 2006.
- [15] 3GPP2 C.P9010, "Draft baseline text for the physical layer portion of the 1X EV specification," pp. 9–78, August 2000.
- [16] H.Z. Sung, J.W. Kang, and K.B. Lee, "A Simplified Maximum Likehood Detection for MIMO Systems," IEICE Trans. Commun., vol. E89-B, no. 8, pp. 2241– 2244, August 2006.

123



### Impact of Depolarization Phenomena on Polarized MIMO Channel Performances

#### Nuttapol PRAYONGPUN, Kosai RAOOF

Laboratoire Grenoble Images Paroles Signal Automatique (GIPSA-LAB), UMR CNRS 5216 961, Rue de la Houille Blanche - BP 46 - 38402 Saint Martin d'Hères, France E-mail: {nuttapol.prayongpun, kosai.raoof}@gipsa-lab.inpg.fr

#### Abstract

The performance and capacity of multiple-input multiple-output (MIMO) wireless channels are limited by the spatial fading correlation between antenna elements. This limitation is due to the use of mono polarized antennas at receiver and transmitter sides. In this paper, in order to reduce the antenna correlation, the polarization diversity technique is employed. Although the spatial antenna correlation is attenuated for multipolarization configurations, the cross-polar components appear. This paper highlights the impact of depolarization effect on the MIMO channel capacity for a  $4\times4$  uniform linear antenna array. We assume that the channel is unknown at the transmitter and perfectly known at the receiver so that equal power is distributed to each of the transmit antennas. The numerical results illustrate that for low depolarization and spatial correlation, the capacity of single-polarization configurations.

**Keywords:** Multiple-input Multiple-output (MIMO), Channel Capacity, Spatial Fading Correlation, Multipolarized Antenna Arrays, Depolarization Effects.

#### 1. Introduction

For the next-generation of wireless communication systems, multiple antennas at both transmitter and receiver could be engaged to achieve higher capacity and reliability of wireless communication channels, under rich scattering environments, in comparisons with traditional single antennas. Due to the potential use of MIMO systems on a limited bandwidth and transmission power, the initial researches demonstrate that the uncorrelated channel capacity can be proportionally increased according to the number of antennas [4,5].

Unfortunately, in practice, the performances of MIMO communication channel are affected by spatial correlation and channel environments [6]. The spatial correlation depends on the array configuration such as radiation pattern, antenna spacing and array geometry. The channel environments are dependent on the environment characteristics such as number of channel paths, distribution and properties of scatterers, angle spread and cross-polarization discrimination [8–10].

Thus, the antenna arrays at transmitter and receiver should be properly designed to reduce the spatial correlation effects and to improve the communication performances [11].

However, it is possible to reduce this effect traditionally by increasing antenna array spacing, but it is not often suitable to apply in some wireless applications where the array size is limitted. Therefore, to eliminate spatial correlation effects with high transmission performances, there are essentially two diversity techniques; pattern and polarization diversity techniques [12,16]. For pattern diversity technique, the radiation of antennas should be generated in a manner to isolate the radiation pattern. For polarization diversity technique, the antennas are designed to radiate with orthogonal radiation polarizations to create uncorrelated channels. In general, there are more than two diversity techniques employed in MIMO wireless systems. However, there are also other techniques such as multimode diversity that exploits the difference of higher order modes to obtain low correlated channel [14]. Polarization diversity technique can be used with pattern

diversity technique in order to boost channel capacity.

Numerous MIMO channel models have already been proposed in literature. In this paper we focus on geometry-based stochastic channel models (GSCM) [15,16]. This calculates the channel response by taking into account the characteristics of wave propagation, Tx-Rx environments, and the scattering mechanisms. All parameters are statistically set to closely match the measured channel observation.

In this paper, we define a geometric scattering model based on a three-dimensional double bouncing model that takes into account the antenna configuration [17–19]. All antennas are provided as a uniform linear array with isotropic or dipole antennas at transmitter and receiver sides. However, all scatterers are uniformly distributed on scattering areas and take into account the cross-polar discrimination (XPD). This parameter indicates the ratio of the co-polarized average power to the cross-polarized average power. Therefore, scattering matrix is used to describe the depolarization of incident wave for each scatterer. Afterward, to simplify the simulated environment configuration, we assume that the angle of arrival and that of departure are uniformly distributed.



Figure 1. Geometries of MIMO channel

We present a simulation study of the spatial correlation and moreover the channel capacity of singleand dual-polarized antenna arrays applied to 4×4 MIMO system. All antenna elements are separated by a half wavelength even in the case of the dual polarization configuration. In addition, we examine the cross-polar discrimination effects on MIMO polarized channel capacity for different antenna configurations.

In Section 2, we provide electromagnetic patterns regarding different electric dipoles. These patterns are then used in Section 3 to create a channel model combining the effect of space separation, polarization antenna gains and depolarization mechanisms. In Section 4, we apply the information theory in order to examine the MIMO channel capacity. Finally, in Section 5, we analyze the numerical results of single- and dual-polarization configurations.

#### 2. Antennas

In practice, not only the propagation environment has an important role but the proper implementation of the

antennas plays also another dominant role for determining the multiple antenna transmission performances. The receiving signals on one element antenna can be correlated to that of another element antenna. Therefore, the systems, which can achieve the best performances, should properly reconfigure the transmitting or/and receiving antenna element arrays with the channel state information derived from the propagation channels.

Table 1. Patterns for different electric dipoles

	$G_x$	$G_y$	$G_{z}$
$G_{ heta}ig( heta, \phiig)$	$-\cos\theta\cos\phi$	$-\cos\theta\sin\phi$	$\sin \theta$
$G_{\phi}ig( heta,\phiig)$	sin Ø	$-\cos\phi$	0

Here we are interested in one array configuration. Orthogonally oriented antennas can offer orthogonal polarization, which corresponds to a complete separation between individual channels, although the antennas are co-located. Thus, using multiple polarization technique helps to guarantee an effective antenna deployment space. However, the receiving energy can be reduced due to the imbalance of depolarization mechanisms.

Three dipole antennas are concerned in this paper; x-, y- and z-oriented dipole antennas. Their patterns of electromagnetic radiations can be simplified by neglecting path loss and distance phase because the electromagnetic radiations are homogenously and identically diffused in the far field case. Their simple expressions of radiation patterns are given by [1]

$$G = G_{\theta} \left( \theta, \phi \right) \bar{\theta} + G_{\phi} \left( \theta, \phi \right) \bar{\phi}$$
(1)

where  $G_{\theta}(\theta, \phi)$  and  $G_{\phi}(\theta, \phi)$  are the antenna gains at elevation and azimuth directions. These gains also depend on the propagation direction. The radiation patterns of differently oriented dipoles are shown in Table 1.

In this paper, the propagation patterns of these antennas are normalized with the isotropic antenna which is specified as the reference antenna.

#### 3. Geometric Scattering Modelling

We focus on a useful model called "geometric scattering model" which is based on the assumption that scatterers around the transmitter and receiver organize the AOD and AOA respectively within transmit and receive scattering areas [15,16,18]. The scatterers are randomly located with according to a certain probability distribution. In particular, the scatterers are additionally used to represent the depolarization and attenuation mechanism of incident waves. To reduce the computational time, we consider that only one propagation path channel occurs when one of transmit and one of receive scatterers are randomly linked. Then the actual channel impulse response is established by a simplified ray-tracing route.

By using our simulated double bounce geometric scattering model as seen in Figure 1, we employ a uniform linear array at both transmitter and receiver. The height of transmitter and receiver has the same level. Moreover, transmit and receive scatterers are uniformly distributed within an angular region characterized by  $|\phi + \pi/2| \le \Delta \phi/2$  in elevation area and  $|\theta + \pi/2| \le \Delta \phi/2$  in azimuth area at transmitter and  $|\phi - \pi/2| \le \Delta \phi/2$  in elevation area and  $|\theta - \pi/2| \le \Delta \phi/2$  in elevation area at receiver.

Subsequently transmit and receive scatterers are randomly paired as previously mentioned. From one transmit scatterer to one receive scatterer, there is a double depolarization mechanism which is replaced by one scattering matrix. We also assume that the channel coherence bandwidth is larger than the transmitted bandwidth of the signal. This channel is usually called frequency non-selective or flat fading channel.

In the case of far field transmission without line-ofsight channel, the narrowband (flat fading) transmission channel between the antenna p at the transmitter and the antenna m at the receiver can be expressed as [20]

$$h_{mp}(t,f) = \frac{1}{\sqrt{N_S}} \sum_{i=1}^{N_S} a_m^{(i)} a_p^{(i)} \exp\left\{-j\vec{k}^{(i)} \cdot \vec{v}_{Rx}t - j\vec{k}^{\prime(i)} \cdot \vec{v}_{Tx}t + \varphi_{mp}\right\} \\ \left[G_{\theta}^m\left(\theta_i,\phi_i\right) \quad G_{\phi}^m\left(\theta_i,\phi_i\right)\right] \mathbf{S}_{mp}^{(i)} \begin{bmatrix}G_{\theta}^p\left(\theta_i,\phi_i\right)\\G_{\phi}^p\left(\theta_i,\phi_i\right)\end{bmatrix}$$
(2)

where  $N_s$  is the number of scatterers at the receiver and the transmitter;  $\vec{v}_{Tx}$  and  $\vec{v}_{Rx}$  are the velocity vector of the transmitter and the receiver;  $\vec{k}^{\prime(i)}$  and  $\vec{k}^{(i)}$  are the vectors of wave number in the direction of the *i*th transmit scatterer and the ith receive scatterer where  $\left|\vec{k}^{(i)}\right| = \left|\vec{k}^{\prime(i)}\right| = 2\pi/\lambda$ ;  $G^{p}_{\theta}\left(\theta_{i},\phi_{i}\right)$  and  $G^{p}_{\phi}\left(\theta_{i},\phi_{i}\right)$  are the gain in the direction of  $\vec{\theta}$  and  $\vec{\phi}$  of the *p*th transmit antenna in the direction of the *i*th transmit scatterer.  $G_{\theta}^{m}(\theta_{i},\phi_{i})$  and  $G_{\phi}^{m}(\theta_{i},\phi_{i})$  are the gain in the direction  $\vec{\theta}$  and  $\vec{\phi}$  of the *m*th receive antenna in the direction of the *i*th receive scatterer; *t* is time;  $a_m^{(i)}$  is the *m*th element of the local vector of the receive antenna, so that the local receive vector can be expressed  $\mathbf{a}_{\text{Rx}}^{(i)} = \begin{bmatrix} 1 & e^{-j\vec{k}^{(i)}\cdot\vec{r_1}} & \cdots & e^{-j\vec{k}^{(i)}\cdot\vec{r_{M-1}}} \end{bmatrix}$ ,  $a_p^{(i)}$  is the *p*th element of the local vector of the transmit antenna, where a local transmit vector is expressed as  $\mathbf{a}_{\text{Tx}}^{(i)} = \begin{bmatrix} 1 & e^{-j\vec{k}'^{(i)}\cdot\vec{r}'_{1}} & \cdots & e^{-j\vec{k}'^{(i)}\cdot\vec{r}'_{N-1}} \end{bmatrix} \quad ; \quad \mathbf{S}_{mp}^{(i)}$ are the

scattering matrix for the *p*th transmit scatterer and the *m*th receive scatterer and is also defined as

$$\mathbf{S}_{mp}^{(i)} = \begin{bmatrix} S_{\theta\theta}^{(i)} & S_{\phi\theta}^{(i)} \\ S_{\theta\phi}^{(i)} & S_{\phi\phi}^{(i)} \end{bmatrix}$$
(3)

The cross polarization discrimination (XPD) is defined as the average power ratio of the co-polarization component and the cross-polarization component.

$$\begin{aligned} XPD_{\theta} &= E\left\{ \left| S_{\theta\theta} \right|^{2} \right\} / E\left\{ \left| S_{\theta\phi} \right|^{2} \right\} \\ XPD_{\phi} &= E\left\{ \left| S_{\phi\phi} \right|^{2} \right\} / E\left\{ \left| S_{\phi\phi} \right|^{2} \right\} \end{aligned} \tag{4}$$

In some conditions such as the imbalance of depolarization and the use of different antenna patterns,  $XPD_{\theta} \neq XPD_{\phi}$ . We assume that the sum of the copolarized power and the cross-polarized power is constant. Therefore the scattering matrix can be written as

$$\mathbf{S}_{mp}^{(i)} = \begin{bmatrix} \sqrt{\frac{XPD_{\theta}}{1+XPD_{\theta}}} e^{\left\{j\varphi_{\theta\theta}^{(i)}\right\}} & \sqrt{\frac{1}{1+XPD_{\phi}}} e^{\left\{j\varphi_{\theta\theta}^{(i)}\right\}} \\ \sqrt{\frac{1}{1+XPD_{\theta}}} e^{\left\{j\varphi_{\theta\theta}^{(i)}\right\}} & \sqrt{\frac{XPD_{\phi}}{1+XPD_{\phi}}} e^{\left\{j\varphi_{\theta\theta}^{(i)}\right\}} \end{bmatrix}$$
(5)

where  $\varphi_{\theta\phi}^{(i)}$  denotes phase offset of *i*th incident wave which changes from  $\vec{\theta}$  direction to  $\vec{\phi}$  direction and superposing on *mp* channel.

#### 4. MIMO Capacity

In this section, we assume that the noise has a Gaussian distribution. Therefore, the optimal distribution of input signal is Gaussian for maximizing the mutual information (MI). The mutual information is given by [4,5]

$$I = \log_2 \det \left( \mathbf{I}_{N_R} + \mathbf{H} \Phi \mathbf{H}^{\dagger} \left( \mathbf{K}_{\mathbf{n}} \right)^{-1} \right)$$
 (6)

where  $\Phi = E(\mathbf{x}\mathbf{x}^{\dagger})$  is the spatial covariance matrix of the input vector x under the total transmitting power constraint  $tr(\Phi) = P_t$  and  $K_n$  is the covariance matrix of the noise vector n.  $(\cdot)^{\dagger}$  denotes the conjugate transpose operator,  $E(\cdot)$  is the expected value and  $tr(\cdot)$  is the trace operator.

When the MIMO channel state information (CSI) is known at the receiver but unknown to the transmitter and n is complex additive white Gaussian noise (AWGN) vector with zero mean, the covariance is equal to  $\mathbf{K}_{n} = \sigma_{n}^{2} \mathbf{I}_{N_{R}}$ . When CSI is not available at the transmitter, the transmitter splits equally the total power to each transmitting antenna. Then the input covariance matrix is a diagonal matrix  $\Phi = P_{t}/N_{T} \cdot \mathbf{I}_{N_{T}}$ .

Copyright © 2008 SciRes.



Figure 2. 4×4 MIMO channel capacity of isotropic antennas: (a) single-polarization system and (b) dualpolarization system

Therefore, the average MI, E(I), called the ergodic channel capacity, with equal-power allocation at transmitter can be written as

$$C_{noCSI} = E(I) = E_{\mathbf{H}} \left[ \log_2 \det \left( \mathbf{I}_{N_R} + \frac{P_I}{N_T \sigma_{\mathbf{n}}^2} \mathbf{H} \mathbf{H}^H \right) \right]$$
(7)

By applying an eigenvalue decomposition, (7) can be rewritten as

$$C_{noCSI} = E_{\mathbf{H}} \left[ \sum_{i=1}^{M} \log_2 \left( 1 + \frac{P_t}{N_T \sigma^2} \lambda_{\mathbf{H},i}^2 \right) \right]$$
(8)

where  $M = \min(N_T, N_R)$  that corresponds to the rank of channel matrix and  $\lambda_{H_i}^2$  is the *i*th eigenvalue of H.

#### 5. Simulation Results Based on Geometric Scattering Modelling

#### 5.1. Capacity Versus Angle Spread

The antenna correlation effect is an important indicator for transmission performance since lower correlation will tend to produce higher mean channel capacity for single polarization system as seen in Figure 2. Thus employing polarization and angular diversity techniques is an attractive way to improve MIMO systems. The 4×4 MIMO systems employ isotropic antennas for  $\lambda/2$ antenna spacing as shown in Figure 1. In order to estimate the channel capacity of different antenna configuration, the simulated environments must be identical. Hence the channel capacities are studied in terms of different antenna configurations. The radiation patterns of each antenna are normalized by the radiation pattern of an isotropic antenna.

As mentioned in previous section, the distribution of angles of departure is assumed to have a uniform elevation distribution  $|\phi + \pi/2| \le \Delta \phi/2$  and a uniform arrival azimuth distribution  $|\theta + \pi/2| \le \Delta \theta/2$  and the distribution of angles of arrival is assumed to have a uniform elevation distribution  $|\phi - \pi/2| \le \Delta \phi/2$  and a uniform arrival azimuth distribution  $|\theta - \pi/2| \le \Delta \phi/2$  and a uniform arrival azimuth distribution  $|\theta - \pi/2| \le \Delta \theta/2$  and a uniform arrival azimuth distribution  $|\theta - \pi/2| \le \Delta \theta/2$  where  $\Delta \phi = \Delta \theta = AS$  and  $XPD_{\theta}=XPD_{\theta}=XPD=0dB$  with 20 scatterers at both transmitter and receiver and 15 dB of SNR. The aim of this section is to study the effects of angle spreads and antenna radiation patterns in terms of ergodic capacity.

Figure 2 demonstrates 4×4 MIMO channel capacity of single and dual polarized configuration. For singlepolarization case, only azimuth isotropic antennas are employed and for dual-polarization case, we put successively azimuth and elevation isotropic antennas in order with  $\lambda/2$  antenna spacing. From Figure 2a, the MIMO channel capacity increases as the angle spread increases at transmitter and receiver for the same polarization antennas. In contrast, the dual polarization achieves better channel capacity due to the lower antenna correlation. It founds that the MIMO channel capacity is significantly dependent on the antenna correlation. The polarization diversity technique can diminish the spatial correlation effect and improve the system performances as shown in Figure 2b.



Figure 3. Difference between dual-polarized and singlepolarized channel capacity of 4×4 MIMO systems in functions of XPD and AS

Copyright © 2008 SciRes.



Figure 4. A 4×4 MIMO configuration of  $20^{\circ}$  transmit and  $180^{\circ}$  receive angle spreading: (a) Channel capacity and (b) Subchannel power

#### 5.2. Capacity Versus Depolarization Effects

If multi-polarized antenna array is employed, the spatial correlation effect can be reduced or eliminated due to low radiation pattern interference. Nevertheless, the cross-polarization discrimination (XPD) becomes the most important parameter because XPD represents the ratio of the co-polarized average received power to the cross-polarized average received power. Then, with high XPD value, less energy is coupled between the cross-polarized channels. Even if the capacity of multipolarized antenna arrays can remain high particularly at lower XPD and the higher K-factor values [17], single-polarized antenna array performance can effectively provide better than that of multi-polarized antenna array at higher XPD and lower spatial correlation value.

Figure 3 explains the difference between dualpolarized and single-polarized channel capacity of  $4\times 4$ MIMO systems  $(\Delta C = C_{\text{single-polar}} - C_{\text{dual-polar}})$  in functions of XPD and AS. We also consider that they have the same angle spreads (AS) at both transmitter and receiver sides. For a high XPD and a sufficiently wide angle spread, we note that the MIMO channel capacity of the single-polarized antenna is superior to that of the dual-polarized antenna because a product of the subchannel power is higher.

Figure 4 demonstrates the capacity variation in function of polarization decoupling and also subchannel power of channel matrix for isotropic and dipole antennas. We setup a  $4\times4$  MIMO system with  $20^{\circ}$  transmit and  $180^{\circ}$  receive angle spreading to achieve high transmit and low receive spatial correlation.

The channel capacity of the isotropic and dipole antenna configurations in Figure 4a is slightly different, because the transmission power is normalized with respect to the transmission power of the isotropic antennas. That is the reason why we have the same subchannel power for the isotropic and dipole antennas in Figure 4b. Although MIMO subchannel power of single polarization system is superior to that of dual polarization at high XPD, the single-polarized MIMO configuration cannot benefit of this high channel power due to the significant transmit correlation as shown in Figure 4a.

The subchannel power of channel matrix can be calculated by employing the Frobenius norm. The numerical results confirm that for high XPD case, the co-polarized channel components still have a significant value compared to the cross-polarized channels. The average transmission power of single-polarized isotropic antenna arrays is given by

$$\left\|\mathbf{H}\right\|_{F} = N_{R}^{x} N_{T}^{x} = 4 \times 4 = 16 \tag{9}$$

where  $(\cdot)^x$  represents the dipole orientation. In the case of low XPD, the average transmission power of singlepolarized antenna arrays tends to zero,  $\|\mathbf{H}\|_F \Rightarrow 0$ , because of the loss of co-polarized channel power as shown in Figure 4b. The channel power is directly proportional to the channel capacity as shown in Figure 4a and Figure 4b. In contrast, the average transmission power of dual-polarized antenna arrays is independent of the XPD value, and it approaches to

$$\|\mathbf{H}\|_{E} \approx N_{R}^{x} N_{T}^{x} + N_{R}^{z} N_{T}^{z} = 2 \times 2 + 2 \times 2 = 8$$
(10)

as illustrated in Figure 4b where  $(\cdot)^x$  and  $(\cdot)^z$  denote the dipole type shown in Table 1.

#### 6. Conclusions

The performance of MIMO communication systems is essentially affected by the spatial correlation and channel environments. The spatial correlation depends on the array configurations and the channel characteristics. Therefore to achieve the optimum performances with MIMO systems, the proper selection of array configuration is required. In this paper, we studied the MIMO wireless channel capacity of singleand multi-polarized antenna arrays applied to a uniform linear array with two isotropic antenna configurations.

The simulation results demonstrate that for the nonline-of-sight (NLOS) case, the use of multi-polarization antennas can provide capacity improvement over conventional single-polarization antennas for narrow angle spread. However, when the cross-polarization discrimination is superior than 0dB corresponding to high co-polarized channel power and low crosspolarized channel power, the subchannel power of single-polarization system can be higher by employing the same polarization as that of the co-polarized channel. Thus, with high XPD and low spatial correlation values, single-polarized antenna array performance can effectively provide better capacity than that of multipolarized antenna array. Finally, the cross-polarization discrimination should be also investigated before employing the polarization diversity technique.

#### 7. References

- C.A. Balanis, "Antenna Theory," Second Edition, John Wiley & Sons, New York, 1997.
- [2] L.M. Correia, Ed., Wireless Flexible Personalised Communication (COST 259 Final Report), Wiley, 2001
- [3] Mobile Broadband Multimedia Networks (COST 273 Final Report), Elsevier, 2006.
- [4] E. Telatar, "Capacity of multi-antenna Gaussian channels," European Trans. Telecommun., vol. 10, no. 6, pp. 585–595, November–December 1999.
- [5] G.J. Foschini and M.J. Gans, "On the limits of wireless communications in a fading environment when using multiple antennas," Wireless Personal communication, vol. 6, pp. 311–335, March 1998.
- [6] D.S. Shiu, G.J. Foschini, and M.J. Gans, "Fading Correlation and Its Effect on the Capacity of Multielement Antenna Systems," IEEE Trans. Comm.,vol. 48, no. 3, pp. 502–513, March 2000.
- [7] M. Ali Khalighi, K. Raoof, and G. Jourdain, "Capacity of Wireless Communication Systems Employing Antenna Arrays, a Tutorial Study," Kluwer Academic Publishers, vol. 23, no. 23, pp. 321–352, 2002.
- [8] P.Kyritsi et al., "Effect of antenna polarization on the capacity of a multiple element system in an in-door

environment," IEEE J. Sel. Areas Commun., vol. 20, pp. 1227–1239, August 2002.

- [9] "Correlation analysis based on MIMO channel measurements in an indoor environment," IEEE J. Sel. Areas Commun., vol. 21, pp. 713–720, June 2003.
- [10] J. Lempiainen, J. K. Laiho-Steffens, and A. F. Wacker, "Experimental results of cross polarization discrimination and signal correlation values for a polarization diversity scheme," in Proceeding VTC, pp. 1498–1502, May 1997.
- [11] M.A. Jensen and J.W. Wallace, "A Review of Antennas and Propagation for MIMO Wireless Communication," IEEE Trans. Antennas Propagat., vol.52 no. 11, November 2004.
- [12] M.R. Andrews, P.P. Mitra, and R. deCarvalho, "Tripling the capacity of wireless communication using electromagnetic polarization," Nature, vol. 409, pp. 316– 318, January 2001.
- [13] T. Svantesson, M.A. Jensen, and J.W. Wallace, "Analysis of electromagnetic field polarizations in multi antenna systems," IEEE on Wireless Commun., vol. 3, no.2, pp. 641–646, March 2004.
- [14] T. Svantesson., "On the potential of multimode antenna diversity," in Proceeding VTC, pp. 2368–2372, September 2000.
- [15] 3GPP–3GPP2 Spatial Channel Model Ad-hoc Group3GPP TR 25.996, "Spatial Channel Model for Multiple Input Multiple Output (MIMO) Simulations," v6.1.0 (2003–09).
- [16] C. Oestges, V. Erceg, and A.J. Paulraj, "Propagation Modeling of MIMO Multipolarized Fiwed Wireless Channels," IEEE Trans. Veh. Technol., vol. 53, no. 3, May 2004.
- [17] N. Prayongpun and K. Raoof, "MIMO Channel Capacities in Presence of Polarization Diversity with and without Line-of-Sight Path," Journal WSEAS Trans. on Commun., vol. 5, no. 9, pp. 1744–1750, September 2006.
- [18] K. Raoof and N. Prayongpun, "Channel Capacity Performance for MIMO polarized diversity systems," IEEE/WCNM2005, pp. 1–4, September 23–26, 2005.
- [19] N. Prayongpun and K. Raoof, "Impact of Depolarization Effects on Polarized MIMO Channel Performances," IEEE/WiCOM2007, pp. 1–4, September 21–23, 2007.
- [20] K.Raoof, A.Khalighi, and N.Prayongpun, "MIMO Systems: Principles, advanced polarization and iterative techniques," Adaptive Signal Processing inWireless Communications, Editor: M. Ibnkahla, CRC, pp.95–134, New York, 2008, ISBN–10: 1420046012.



## A Quadratic Constraint Total Least-squares Algorithm for Hyperbolic Location

Kai YANG<sup>1</sup>, Jianping AN, Zhan XU

Department of Electronic Engineering, School of Information Science and Technology, Beijing Institute of Technology, Beijing 100081, P.R. China E-mail: <sup>1</sup>yangkbit@gmail.com

#### Abstract

A novel algorithm for source location by utilizing the time difference of arrival (TDOA) measurements of a signal received at spatially separated sensors is proposed. The algorithm is based on quadratic constraint total least-squares (QC-TLS) method and gives an explicit solution. The total least-squares method is a generalized data fitting method that is appropriate for cases when the system model contains error or is not known exactly, and quadratic constraint, which could be realized via Lagrange multipliers technique, could constrain the solution to the location equations to improve location accuracy. Comparisons of performance with ordinary least-squares are made, and Monte Carlo simulations are performed. Simulation results indicate that the proposed algorithm has high location accuracy and achieves accuracy close to the Cramer-Rao lower bound (CRLB) near the small TDOA measurement error region.

Keywords: Location, Time Difference of Arrival, Total Least-squares

#### 1. Introduction

Determining the location of a source from its emissions is a critical requirement for the deployment of wireless sensor networks in a wide variety of applications [1]. Location finding based on time difference of arrival (TDOA), which does not require knowledge of the absolute transmission time, is the most popular method for accurate positioning systems [2]. The idea of TDOA is to determine the source location relative to the sensors by examining the difference in time at which the signal arrives at multiple measuring sensor units, rather than the absolute arrival time. Sensors at separate locations measuring the TDOA of the signal from a source can determine the location of the source as the intersection of hyperbolae for TDOA. For each TDOA measurement, the source lies on a hyperbola with a constant range difference between the two measuring sensors. However, finding the solution to the hyperbolic location equations is not easy as the equations are nonlinear. Furthermore, the nonlinear hyperbolic equations become inconsistent as errors occur in TDOA measurements and the hyperbolae no longer intersect at a single point.

In the past, the source location determination problem

has been mathematically formulated as a set of linear equations **Ax=b** which is in matrix form and ordinary least-squares (OLS) technique is utilized to find the maximum-likelihood solution by assuming the system matrix A is error-free and all errors are confined to the data vector **b** [3–6]. However, in the TDOA based location problem there are errors in both system matrix A and data vector **b**. Using OLS technique for this problem will result in biased solution and location accuracy will decrease due to the accumulation of the system matrix errors. To alleviate this problem, a generalization of the OLS solution, called total leastsquares (TLS) [7-9], is utilized to remove the noise in A and **b** using a perturbation on **A** and **b** of the smallest Frobenius norm which makes the system of equations consistent. It is shown that in the TDOA based location problem the unknown parameters in vector  $\mathbf{x}$  are quadratic constraint related, which could be realized via Lagrange multipliers technique to constrain the TLS solution.

For practical applications, the location estimation algorithm should be robust and easy to implement, and the source location estimation should minimize its deviation from the true location. In this paper, a novel TDOA based QC-TLS algorithm for source location

where

problem is presented. The rest of the paper is organized as follows. The formulation of TDOA based location estimation problem is described in Section 2. In Section 3, this problem is shown to be in the form of an overdetermined set of linear equations with errors in both **A** and **b** and the QC-TLS algorithm is applied to the TDOA measurement data for source location estimation. Computer simulations are used to verify the validity of the algorithm and simulation results are presented in Section 4, which is followed by conclusions.

#### 2. Problem Formulation

There are two basic ways to measure the TDOA in wireless sensor networks: the direct way and the indirect way. In the direct way, we could obtain the TDOA through the use of cross-correlation techniques, in which the received signal at one sensor is correlated with the received signal at another sensor. The timing requirement for this method is the synchronization of all the receivers participating in the TDOA measurements, which is more available in location applications. However, in multipath channels there is ambiguity in detecting the TDOA using cross-correlation techniques since the correlation peak to be detected may not be caused by the two direct line of sight (LOS) paths [10]. In the indirect way, we first obtain the arrival time of received signal transmitted from the source at spatially separated sensors and then subtract the measurements of arrival time at two sensors to produce a relative TDOA. TDOA also could be obtained by subtracting absolute measurements. but this requires TOA the synchronization of source and sensors. In this paper, because of these factors above we obtain the TDOA by subtracting the measurements of arrival time at two sensors.

When TDOA based location method is adopted to give source location estimation in wireless networks, according to TDOA measurements a set of hyperbolic equations is given by

$$r_{i1} = c\tau_{i1} = r_i - r_1, i = 2, 3, \cdots, M$$
  

$$r_i = \sqrt{\left(x_i - x_s\right)^2 + \left(y_i - y_s\right)^2}$$
(1)

where *c* is the speed of signal propagation,  $\tau_{i1}$  is the true value of TDOA measurement between sensor *i* and 1,  $r_i$  is the distance between the source and sensor *i*, *M* is the number of sensors, and  $(x_s, y_s)$  and  $(x_i, y_i)$  are the coordinates of the source and sensor *i* respectively.

Solving those nonlinear equations is difficult. Linearizing them and then solving is one possible way. From (1), we have

$$r_{i1} + r_1 = r_i \tag{2}$$

Substituting the expressions of  $r_i$  and  $r_1$  into (2) and then squaring both sides of (2) produces

$$(x_{i} - x_{1})(x_{s} - x_{1}) + (y_{i} - y_{1})(y_{s} - y_{1}) + r_{i1} \cdot r_{1}$$
  
=  $\frac{1}{2} ((x_{i} - x_{1})^{2} + (y_{i} - y_{1})^{2} - r_{i1}^{2})$  (3)

Formulation it in matrix form we have

$$\mathbf{A}\mathbf{x} = \mathbf{b} \tag{4}$$

$$\mathbf{A} = \begin{bmatrix} x_2 - x_1 & y_2 - y_1 & r_{21} \\ x_3 - x_1 & y_3 - y_1 & r_{31} \\ \vdots & & \\ x_M - x_1 & y_M - y_1 & r_{M1} \end{bmatrix}$$
$$\mathbf{x} = \begin{bmatrix} x_s - x_1 & y_s - y_1 & r_1 \end{bmatrix}^T$$
$$\mathbf{b} = \frac{1}{2} \begin{bmatrix} (x_2 - x_1)^2 + (y_2 - y_1)^2 - r_{21}^2 \\ (x_3 - x_1)^2 + (y_3 - y_1)^2 - r_{31}^2 \\ \vdots \\ (x_M - x_1)^2 + (y_M - y_1)^2 - r_{M1}^2 \end{bmatrix}$$

and the superscript T denotes the matrix transpose operation.

In the absence of noise and interference, hyperbolae from two or more sensors will intersect to determine a unique location which means that the overdetermined set of linear equations (4) is consistent. In the presence of noise, more than two hyperbolae will not intersect at a single point which means that (4) is inconsistent. In the following section, we will develop the solution to the overdetermined equations.

#### 3. Hyperbolic Location Solution

In this section we first analyze the OLS solution to (4) and then give the QC-TLS solution.

#### 3.1. OLS Approach

In our applications, the usual assumption for OLS approach is that the matrix **A** consists of M-1 observations on each of the 3 independent variables. The dependent variable is represented by vector **b**, in which we try keeping the correction term  $\Delta$ **b** as small as possible while simultaneously compensating for the noise present in **b** by forcing **Ax**=**b**+ $\Delta$ **b** [7].

Under these assumptions the least-squares estimator

$$\mathbf{x} = \left(\mathbf{A}^T \mathbf{A}\right)^{-1} \mathbf{A}^T \mathbf{b}$$
 (5)

Copyright © 2008 SciRes.

in which it is implicit that A is known exactly. Suppose now that the elements of the A matrix contain errors. The system matrix  $\mathbf{A}$  resulting from the calculation of (4) by using the measured values of TDOA can be given by  $A=A_0+E$ , where  $A_0$  and E represent the true value and error respectively. Under the assumption that the error E is reasonably small, retaining only the linear error terms, the OLS estimator is given by [11]

$$\mathbf{x}(\mathbf{E}) = \mathbf{x} + \left(\mathbf{A}_0^T \mathbf{A}_0\right)^{-1} \mathbf{E}^T \boldsymbol{\eta} - \left(\mathbf{A}_0^T \mathbf{A}_0\right)^{-1} \mathbf{A}_0^T \mathbf{E} \mathbf{x}$$
(6)

where  $\eta$  is the vector of residuals **b**-Ax which would be obtained if A were known accurately.

The sensitivity of each estimated parameter  $x_k$  to each element  $A_{ii}$  of system matrix **A** comes immediately from (6) by taking the derivative of  $x_k$  with respect to  $A_{ji}$  and this yields

$$\frac{\partial x_k}{\partial A_{ji}} = \left(\mathbf{A}_0^T \mathbf{A}_0\right)_{ki}^{-1} \eta_j - \left(\sum_{l=1}^3 \left(\mathbf{A}_0^T \mathbf{A}_0\right)_{kl}^{-1} A_{jl}\right) x_i$$
(7)

The sensitivities are different, both in value and ranking, from the sensitivities to change in the dependent variable. It would in particular give details of those observations most liable to cause estimation error.

#### **3.2. QC-TLS Approach**

In hyperbolic location problem, it can be argued that both the system matrix **A** and vector **b** are subject to error, which is out of accord with the assumption of OLS. To



Figure 1. Least-squares versus total least-squares

alleviate the effects of these errors, we propose to use the TLS solution for the source location problem. The total least-squares method [8,9], which is a natural extension of LS when errors occur in all data, is devised as a more global fitting technique than the ordinary least-squares technique for solving overdetermined sets of linear equations by trying to remove the data errors. The LS and TLS measures of goodness are shown in Figure 1. In the LS approach, it is the vertical distances that are important; whereas in the TLS problem, it is the perpendicular distances that are critical. So, from this

geometric interpretation, it's shown that the TLS method is better than the LS method with respect to the residual error in the curve fitting.

When errors exist in TDOA measurements, (4) can be represented by

$$(\mathbf{A}_0 + \mathbf{E})\mathbf{x} = \mathbf{b}_0 + \mathbf{w}$$
 (8)

where **E** and **w** are the perturbations of **A** and **b**, respectively. And (8) also can be put into the following form

 $\left(\left[-\mathbf{b}_{0} \stackrel{:}{:} \mathbf{A}_{0}\right] + \left[-\mathbf{w} \stackrel{:}{:} \mathbf{E}\right]\right) \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix} = \mathbf{0}$ 

or

(9b)  $(\mathbf{B}_0 + \mathbf{D})\boldsymbol{\alpha} = \mathbf{0}$ 

(9a)

Where

$$\mathbf{B}_{0} = \begin{bmatrix} -\mathbf{b}_{0} \\ \vdots \\ \mathbf{A}_{0} \end{bmatrix}, \ \mathbf{D} = \begin{bmatrix} -\mathbf{w} \\ \vdots \\ \mathbf{E} \end{bmatrix}, \ \boldsymbol{\alpha} = \begin{bmatrix} 1 \\ \mathbf{x}^{T} \end{bmatrix}^{T}$$

The rank of matrix  $\mathbf{B}_0$  is three, because it is equal to the rank of matrix  $A_0$ .

The TLS solution to this problem is to looks for the minimal (in the Frobenius norm sense) perturbation matrix  $\mathbf{E}$  and perturbation vector  $\mathbf{w}$  so that

$$\{\hat{\mathbf{x}}, \mathbf{E}, \mathbf{w}\} = \arg\min_{\mathbf{x}, \mathbf{E}, \mathbf{w}} \|\mathbf{D}\|_{F}^{2}$$
(10)

subject to

$$(\mathbf{A}_0 + \mathbf{E})\mathbf{x} = \mathbf{b}_0 + \mathbf{w}$$

where the subscript F denotes the Frobenius norm.

The TLS optimization involves to find the optimum  $\hat{\alpha}$  for  $\alpha$  that minimizes the cost function  $f(\hat{\alpha})$ ,

$$f(\hat{\boldsymbol{\alpha}}) = \|\mathbf{B}_0 \hat{\boldsymbol{\alpha}}\|_{F}^{2} = \hat{\boldsymbol{\alpha}}^{T} \mathbf{B}_0^{T} \mathbf{B}_0 \hat{\boldsymbol{\alpha}}$$
(11)

In practical applications, the value of matrix  $\mathbf{B}_0$  is unknown. We now make a singular value decomposition (SVD) of the matrix  $\mathbf{B} = [-\mathbf{b} \mathbf{A}]$ , that is,

$$\mathbf{B} = \mathbf{U}_{\mathrm{B}} \boldsymbol{\Sigma}_{\mathrm{B}} \mathbf{V}_{B}^{H}$$
(12)

where  $\Sigma_{\rm B}$  is composed of the singulars value of **B**. To solve (11), matrix  $\mathbf{B}_0$  can be replaced by a rank three optimum approximation of  $\mathbf{B}_0$ , that is,

$$\tilde{\mathbf{B}} = \mathbf{U}_{\mathrm{B}} \boldsymbol{\Sigma}_{3} \mathbf{V}_{B}^{H}$$
(13)

where  $\Sigma_3$  is composed of the maximum three singulars value of **B**. Now we have the cost function

-----

$$f(\hat{\boldsymbol{\alpha}}) = \hat{\boldsymbol{\alpha}}^T \tilde{\mathbf{B}}^T \tilde{\mathbf{B}} \hat{\boldsymbol{\alpha}}$$
(14)

The parameters in vector  $\hat{a}$  are subject to a quadratic constraint

#### A QUADRATIC CONSTRAINT TOTAL LEAST-SQUARES ALGORITHM FOR HYPERBOLIC LOCATION

$$1 + (x_1 - \hat{x}_s)^2 + (y_1 - \hat{y}_s)^2 - \hat{r}_1^2 - 1 = 0$$
 (15)

or, equivalently, as

$$\hat{\boldsymbol{\alpha}}^{\mathrm{T}}\boldsymbol{\Sigma}\hat{\boldsymbol{\alpha}}-1=0 \tag{16}$$

where  $\Sigma$ =diag(1, 1, 1, -1) is a diagonal and orthonormal matrix. The technique of Lagrange multipliers is used and the modified cost function is of the form

$$f(\hat{\boldsymbol{\alpha}}) = \hat{\boldsymbol{\alpha}}^T \tilde{\boldsymbol{B}}^T \tilde{\boldsymbol{B}} \hat{\boldsymbol{\alpha}} + \lambda \left\{ \hat{\boldsymbol{\alpha}}^T \boldsymbol{\Sigma} \hat{\boldsymbol{\alpha}} - 1 \right\}$$
  
$$= \hat{\boldsymbol{\alpha}}^T \left( \tilde{\boldsymbol{B}}^T \tilde{\boldsymbol{B}} + \lambda \boldsymbol{\Sigma} \right) \hat{\boldsymbol{\alpha}} - \lambda$$
 (17)

where  $\lambda$  is the Lagrange multiplier.

The required optimum parameter vector  $\hat{\alpha}$  is found by solving the following linear system of equations [12]

$$\left(\tilde{\mathbf{B}}^{T}\tilde{\mathbf{B}} + \lambda \boldsymbol{\Sigma}\right)\hat{\boldsymbol{\alpha}} = k\mathbf{e}_{1}$$
(18)

where,  $\mathbf{e}_1 = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}^T$ , and normalizing constant *k* is selected so that the first element of solution vector  $\hat{\boldsymbol{\alpha}}$  is equal to one. Solving (18) is very easy, and it is independent of normalization constant *k*, which is unknown. Calculating the inverse of matrix  $(\tilde{\mathbf{B}}^T \tilde{\mathbf{B}} + \lambda \boldsymbol{\Sigma})$  and then normalizing the first element of  $\hat{\boldsymbol{\alpha}}$ , we can obtain the solution vector. The QC-TLS solution vector to (9) is

$$\hat{\boldsymbol{\alpha}} = \frac{\left(\tilde{\boldsymbol{B}}^{T}\tilde{\boldsymbol{B}} + \lambda\boldsymbol{\Sigma}\right)^{-1}k\boldsymbol{e}_{1}}{\boldsymbol{e}_{1}^{T}\left(\tilde{\boldsymbol{B}}^{T}\tilde{\boldsymbol{B}} + \lambda\boldsymbol{\Sigma}\right)^{-1}k\boldsymbol{e}_{1}} = \frac{\left(\tilde{\boldsymbol{B}}^{T}\tilde{\boldsymbol{B}} + \lambda\boldsymbol{\Sigma}\right)^{-1}\boldsymbol{e}_{1}}{\boldsymbol{e}_{1}^{T}\left(\tilde{\boldsymbol{B}}^{T}\tilde{\boldsymbol{B}} + \lambda\boldsymbol{\Sigma}\right)^{-1}\boldsymbol{e}_{1}}$$
(19)

The estimated x and y coordinates of the source are

$$\hat{x}_s = \hat{\boldsymbol{a}}^T \mathbf{e}_2 + x_1$$

$$\hat{y}_s = \hat{\boldsymbol{a}}^T \mathbf{e}_3 + y_1$$
(20)

where

$$\mathbf{e}_{2} = \begin{bmatrix} 0 \ 1 \ 0 \ 0 \end{bmatrix}^{T}$$
$$\mathbf{e}_{3} = \begin{bmatrix} 0 \ 0 \ 1 \ 0 \end{bmatrix}^{T}$$

Lagrange Multiplier  $\lambda$  in the expression of solution vector  $\hat{a}$  is unknown. In order to find  $\lambda$ , we can impose the quadratic constraint directly by substituting (19) into (16), so that

$$\left[\frac{\left(\tilde{\mathbf{B}}^{T}\tilde{\mathbf{B}}+\lambda\Sigma\right)^{-1}\mathbf{e}_{1}}{\mathbf{e}_{1}^{T}\left(\tilde{\mathbf{B}}^{T}\tilde{\mathbf{B}}+\lambda\Sigma\right)^{-1}\mathbf{e}_{1}}\right]^{T}\Sigma\left[\frac{\left(\tilde{\mathbf{B}}^{T}\tilde{\mathbf{B}}+\lambda\Sigma\right)^{-1}\mathbf{e}_{1}}{\mathbf{e}_{1}^{T}\left(\tilde{\mathbf{B}}^{T}\tilde{\mathbf{B}}+\lambda\Sigma\right)^{-1}\mathbf{e}_{1}}\right]-1=0$$
(21)

By using an eigenvalue factorization, the matrix  $\tilde{B}^T \tilde{B} \Sigma$  can be diagonalized as

$$\tilde{\mathbf{B}}^T \tilde{\mathbf{B}} \boldsymbol{\Sigma} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^{-1} \tag{22}$$

Copyright © 2008 SciRes.

where  $\Lambda$ =diag( $\gamma_1$ ,  $\gamma_2$ ,  $\gamma_3$ ,  $\gamma_4$ ),  $\gamma_i$ , *i*=1, ..., 4 are the eigenvalues of the matrix  $\tilde{\mathbf{B}}^T \tilde{\mathbf{B}} \Sigma$ , and **U** is the corresponding eigenmatrix of  $\tilde{\mathbf{B}}^T \tilde{\mathbf{B}} \Sigma$ . Substituting (22) into (21) produces

$$\frac{\mathbf{e}_{1}^{T} \boldsymbol{\Sigma} \mathbf{U} \left( \boldsymbol{\Lambda} + \lambda \mathbf{I} \right)^{-2} \mathbf{U}^{-1} \mathbf{e}_{1}}{\left[ \mathbf{e}_{1}^{T} \boldsymbol{\Sigma} \mathbf{U} \left( \boldsymbol{\Lambda} + \lambda \mathbf{I} \right)^{-1} \mathbf{U}^{-1} \mathbf{e}_{1} \right]^{2}} - 1 = 0$$
(23)

For notational convenience, we define

$$\Sigma \mathbf{U} = \begin{bmatrix} \mathbf{p}_{1}^{T} & \mathbf{p}_{2}^{T} & \mathbf{p}_{3}^{T} & \mathbf{p}_{4}^{T} \end{bmatrix}^{T}$$

$$\mathbf{p}_{i} = \begin{bmatrix} p_{i1} & p_{i2} & p_{i3} & p_{i4} \end{bmatrix}$$

$$\mathbf{U}^{-1} = \begin{bmatrix} \mathbf{q}_{1} & \mathbf{q}_{2} & \mathbf{q}_{3} & \mathbf{q}_{4} \end{bmatrix}$$

$$\mathbf{q}_{i} = \begin{bmatrix} q_{1i} & q_{2i} & q_{3i} & q_{4i} \end{bmatrix}^{T}$$
(24)

Substituting (24) into (23), then the quadratic constraint in the form of  $\lambda$  becomes

$$f(\lambda) = \sum_{i=1}^{4} \frac{p_{1i}q_{i1}}{(\gamma_i + \lambda)^2} \cdot \prod_{i=1}^{4} (\gamma_i + \lambda)^2$$
$$-\left(\sum_{i=1}^{4} \frac{p_{1i}q_{i1}}{\gamma_i + \lambda} \cdot \prod_{i=1}^{4} (\gamma_i + \lambda)\right)^2 = 0$$
(25)

Lagrange Multiplier  $\lambda$  can be solved from (25) by using Newton's method with zero as the initial guess.

To sum up, a brief description of the proposed TDOA based location algorithm is as follows:

- 1) Calculate the optimum approximation of augmented matrix  $\mathbf{B}_0$  using (12) and (13).
- 2) Solve  $\lambda$  by finding the root of (25) using Newton's method.
- 3) Obtain the QC-TLS solution to (9) using (19) and determine x and y coordinates of the source using (20).

#### 4. Simulations

We have proposed a QC-TLS algorithm for TDOA based location problem. In this section, we evaluate its performance at practical measurement error levels using Monte-Carlo simulations.

Source and sensors were located in random positions in a square of area 100×100 m<sup>2</sup> as shown in Fig. 2. We assumed that the TDOA measurement errors were white random processes with zero mean and variance  $\sigma_{\text{TDOA}}^2$ , and the TDOA variance of all sensor inputs were identical. For simplicity, the TDOA variance was translated into corresponding distance variance  $\sigma^2 = (c\sigma_{\text{TDOA}})^2$ . Simulation results provided were averages of 10000 independent runs. We compare the proposed QC-TLS approach with OLS approach and CRLB.

I. J. Communications, Network and System Sciences, 2008, 2, 105-206

Tables 1 and 2 compare the mean absolute location errors (MALEs) of OLS and QC-TLS approach for various TDOA error variances, and MALE was defined as  $E\left[\sqrt{(x_s - \hat{x}_s)^2 + (y_s - \hat{y}_s)^2}\right]$ . The number of

sensors was set to 7 and 10 for Table 1 and Table 2, respectively. In Tables 1 and 2, CRLB, which is a fundamental lower bound on the variance of any unbiased estimator, is calculated using the functions in [4]. The distance unit is meters. We observe that increasing the number of sensors can increase the location accuracy, because increasing the number of sensors can provide more redundant information which could be helpful for solving the overdetermined linear equations.



Figure 2. Location in 2-D plane

Table 1. comparison of MALEs for QC-TLS and OLS method for various variances: seven sensors

$\sigma^2$	0.1	0.01	0.001	0.0001
OLS	0.9359	0.2962	0.0939	0.0293
QC-TLS	0.7735	0.2236	0.0709	0.0221
CRLB	0.6706	0.2134	0.0677	0.0214

Table 2. comparison of MALEs for QC-TLS and OLS method for various variances: ten sensors

$\sigma^2$	0.1	0.01	0.001	0.0001
OLS	0.5801	0.1851	0.0571	0.0182
QC-TLS	0.5410	0.1501	0.0472	0.0150
CRLB	0.4280	0.1361	0.0428	0.0136

The performance of QC-TLS approach is much better than OLS approach especially when the TDOA variance is large, and it is closer to CRLB. They verify that the proposed approach could inhibit the influence of TDOA measurement errors on the estimation results better than OLS. When TDOA variance is large, the noise in matrix **A**, which would be alleviated by QC-TLS approach, could significantly reduce the performance of OLS approach.

Mean absolute location error is one of the metrics to evaluate the performance of QC-TLS algorithm, and computational complexity is another important metric to evaluate the performance of algorithms. The price of the better performance for OC-TLS is that its computational complexity is larger than OLS. In the same simulation condition, the time of 10000 independent runs for QC-TLS and OLS are approximately 6153.6 ms and 1446.3 ms, respectively. It is obviously that QC-TLS algorithm requires matrix singular value decomposition once, eigenvalue decomposition once, Newton iteration once, inverse operation once and matrix multiply operation several times, whereas OLS algorithm only requires matrix pseudoinverse operation once and matrix multiply operation several times. It's noteworthy that QC-TLS algorithm performs better at a cost of larger computational complexity.

#### 5. Conclusions

A novel TDOA based quadratic constraint total leastsquares location algorithm is proposed in this paper. The proposed algorithm utilizes TLS to inhibit the influence of TDOA measurement errors. And the technique of Lagrange multipliers is utilized to exploit the quadratic constraint relation between the intermediate variables to constrain the solution to the location equations and improve the location accuracy. Simulation results indicate that the proposed QC-TLS algorithm gives better results than the OLS solution.

#### 6. References

- N. Patwari, J. N. Ash, S. Kyperountas, A. O. Hero III, R. L. Moses, and N. S. Correal, "Locating the Nodes: Cooperative Localization in Wireless Sensor Networks," IEEE Signal Processing Magazine, vol. 22, no. 4, pp. 54– 68, July 2005.
- [2] X. Jun, L. R. Ren, and J. D. Tan, "Research of TDOA based self-localization approach in wireless sensor network," in proceedings of IEEE International Conference on Intelligent Robots and Systems, Beijing, pp. 2035–2040, October 2006.
- [3] J. Smith and J. Abel, "Closed-form least-squares source location estimation from range-difference measurements," IEEE Transactions on Acoustics, Speech, and Signal Processing, vol. 35, pp. 1661–1669, December 1987.
- [4] Y. T. Chan and K. C. Ho, "A simple and efficient estimation for hyperbolic location," IEEE Transactions on Signal processing, vol. 42, no. 8, pp. 1905–1915, August 1994.
- [5] R. Schmidt, "Least squares range difference location," IEEE Transactions on Aerospace and Electronic Systems, vol. 32, no. 1, pp. 234–242, January 1996.
- [6] Y. Huang, J. Benesty, G. W. Elko, and R. M. Mersereau, "Real-time passive source localization: a practical linearcorrection least-squares approach," IEEE Transactions

On Speech And Audio Processing, vol. 9, no. 8, pp. 943–956, November 2001.

- [7] T. J. Abatzoglou, J. M. Mendel, and G. A. Harada, "The constrained total least squares technique and its applications to harmonic superresolution," IEEE Transactions on Signal Processing, vol. 39, no. 5, pp. 1070–1087, May 1991.
- [8] S. V. Huffel and J. Vandewalle, "The Total Least Squares Problem: Computational Aspects and Analysis," Philadelphia: SIAM, 1991.
- [9] I. Markovsky and S. V. Huffel, "Overview of total leastsquares methods," Signal Processing, vol. 87 no. 10, pp.

2283–2302, October 2007.

- [10] X. Li, "Super-Resolution TOA Estimation with Diversity Techniques for Indoor Geolocation Applications," Ph.D. Dissertation, Worcester Polytechnic Institute, Worcester, MA, 2003.
- [11] S. D. Hodges and P. G. Moore, "Data Uncertainties and Least Squares Regression," Applied Statistics, vol. 21, pp. 185–195, 1972.
- [12] J. A. Cadzow, "Spectral estimation: An overdetermined rational model equation approach," Proceedings of the IEEE, vol. 70, no. 9, pp.907–939, September 1982.



## Analysis of Lifetime of Large Wireless Sensor Networks Based on Multiple Battery Levels

Ruihua ZHANG<sup>1</sup>, Zhiping JIA<sup>1</sup>, Dongfeng YUAN<sup>2</sup>

<sup>1</sup> School of Computer Science and Technology, Shandong University, Jinan 250061, P.R.China
 <sup>2</sup> School of Information Science and Engineering, Shandong University, Jinan 250100, P.R.China E-mail: <sup>1</sup>ruihua\_zhang@sdu.edu.cn

#### Abstract

Due to the limited transmission range, data sensed by each sensor has to be forwarded in a multi-hop fashion before being delivered to the sink. The sensors closer to the sink have to forward comparatively more messages than sensors at the periphery of the network, and will deplete their batteries earlier. Besides the loss of the sensing capabilities of the nodes close to the sink, a more serious consequence of the death of the first tier of sensor nodes is the loss of connectivity between the nodes at the periphery of the network and the sink; it makes the wireless networks expire. To alleviate this undesired effect and maximize the useful lifetime of the network, we investigate the energy consumption of different tiers and the effect of multiple battery levels, and demonstrate an attractively simple scheme to redistribute the total energy budget in multiple battery levels by data traffic load. We show by theoretical analysis, as well as simulation, that this substantially improves the network lifetime.

Keywords: Wireless Sensor Networks, Energy Efficient, Network Lifetime, Battery Level

#### 1. Introduction

Sensor network applications have recently become of significant interest due to cheap single-chip transceivers and micro-controllers. Because sensor nodes are batterypowered, and their operational lifetime should be maximized, one of the most important design criteria for this type of network is energy efficiency.

Confirming the importance of the problem, many aspects of the problem have been extensively studied [1–4]. Medium access control (MAC) layer techniques [2] [3] aim to conserve battery energy by turning the receiver off whenever it is not needed. It is clear that the energy problem cannot be completely solved at any one single layer [4].

The motivation for our work stems from the observation that in a sensor network, the sensor nodes closer to the sink have to relay more packets than the ones at the periphery of the network. We assume that this increase in workload results in an increase in energy consumption, the nodes close to the sink will die first, leading to a premature loss of connectivity in the sensor network. To alleviate this undesirable effect, we study the energy consumption of different tiers, and demonstrate a scheme to redistribute the total energy budget for the sensor network, the lifetime of the network can be significantly improved over the case where all sensors have a uniform lifetime. The optimal solution is formulated theoretically and validated via simulations.

About this problem, in [5-9], non-uniform node deployment is exploited as an alternative manner to get over the effect of non-uniform energy depletion. The basic concept is that different node densities are assigned to different sub-regions trying to balance the communication load of each sub-region. Some works specifically consider the scenario for concentric subregions (rings); the node density of a ring is determined based on the number of hops to the sink, which is roughly approximated to be the same for all locations in a ring in [5–7]. In [8], the imbalanced energy utilization of a WSN was analyzed based on the hop counts, and the authors accordingly proposed a non-uniform node distribution depending on the hop-count to improve the long-term connectivity. In [9], the energy-hole problem was discussed in a more general form for a rectangular sensing area with multiple data sinks at different
locations. Another approach, heterogeneous deployment, overcomes the communication load unbalance problem by using two types of sensors, low energy sensors and high energy sensors, to construct a hierarchical network structure [10]. However, this scheme raises the deployment and implementation complexity and cannot be easily applied.

The remainder of the paper is organized as follows. First we present the network model and the energy consumption model in Section 2. Section 3 provides a mathematical formulation that attempts to estimate the average lifetime of a sensor network and analyzes the design problems. We provide the simulation results in section 4. We conclude the paper with a summary and discussion of future work in Section 5.

#### 2. Energy Consumption Model

Energy consumption models of the radio illustrated by Figure 1. In [11] a model for radio energy consumption is given for energy per bit single hop  $(e_b)$  as:

$$e_b = e_{tx} + e_{rx} \tag{1}$$

$$e_{tx} = e_{ta}d^{\alpha} + e_{te} \tag{1-1}$$



Figure 1. Radio energy consumption model

where  $e_{tx}$  and  $e_{rx}$  are the transmitter and receiver energy consumptions per bit, respectively,  $e_{te}$  is the energy per bit needed by the transmitter electronics,  $e_{ta}$  is the energy needed to successfully transmit one bit over one meter, dis the distance from transmitter to receiver and  $\alpha$  is a constant which depends on the attenuation the signal will suffer in that environment.

Consider a circular sensing field with the radius  $R_L$  meters, a total of N sensor nodes are uniformly distributed within the field. The sink node is in the center of the field. The transmission radius R of the sensor nodes assumed to be fixed. We divide the sensing field into T tiers, T=ceiling ( $R_L/R$ ). Figure 2 shows the "doughnut-like" distribution of nodes [9]. In such a circular field, consider the set of nodes close to the sink at the center that can communicate directly with it. We refer to these one-hop neighbors of the sink as the first tier of nodes. Similarly, the two-hop neighbors of the sink will be the second tier of nodes, etc. It is clear that the first tier nodes relay the largest amount of data

pockets and will die first. If the spatial distribution of nodes is close to uniform, then the traffic load is equally distributed spatially. Each first tier node will relay roughly the same amount of traffic, and all first tier nodes will die at times very close to each other, after the network is first put into operation. Once all of the first tier nodes are dead, no other node will be able to send data to the sink, and the lifetime of the network will be over.



Figure 2. Sensing field

We assume that sensed data is collected in a periodic manner; this period (interval P seconds) consists of the sensing of the data and the transmission of one packet containing the data sensed to the sink. We assume that each sensor has a constant amount of raw data pocket ( $b_s$  bits) to sense and send in interval P.

## 3. Analysis and Numerical Results

#### 3.1. Estimation of Network Lifetime

We define tier *i* as the set of the nodes that can reach the sink in *i* hops  $(1 \le i \le T)$ . We consider the energy budget *E* as the main cost function, *E<sub>i</sub>* is the sum of the energy available at the nodes of tier *i*, so

$$E = \sum_{i=1}^{T} E_i$$
 (2)

 $N_i$  is the number of nodes at tier *i*. With the assumption of a uniform density:

$$N_i = ((2i-1)N)/T^2$$
(3)

 $F_i$  is the total number of data packets, they are relayed to the sink by all nodes at tier *i*; each data packet is generated by one sensor node at tiers *i*+1...*T*.

$$F_{i} = N - \sum_{j=1}^{i} N_{j} = (N(T^{2} - i^{2}))/T^{2}$$
(4)

In tier *i*, nodes relay all data pockets generated by nodes at tiers  $j(i < j \le T)$ , and nodes at *i* tier generate and send all data pockets, let  $E_{relay}$  and  $E_{gen}$  denote their energy consumed, respectively;  $e_s$  is the energy spent sensing per bit;  $b_s$  is data pocket size. So we get the relation with energy consumption and time *t*.

$$E_i = E_{relav} + E_{gen} \tag{5-1}$$

$$E_{relay} = (t / p)F_i b_s e_b$$
(5-2)

$$E_{opn} = (t/p)N_{i}b_{s}(e_{s} + e_{tx})$$
 (5-3)

Thus

$$E_{i} = \frac{tNb_{s}(e_{b}(T^{2} - i^{2}) + (e_{s} + e_{ix})(2i - 1))}{T^{2}p}$$
(5)

 $L_i$  is the lifetime of the nodes at tier *i*, it can be determined by considering the energy consumption. Using (5), we can get:

$$L_{i} = \frac{E_{i}PT^{2}}{Nb_{s}((T^{2} - i^{2})e_{b} + (2i - 1)(e_{s} + e_{tx}))}$$
(6)

In accordance with all the above considerations, we define the lifetime of the network as the minimum of the lifetimes of its tiers:

$$L = \min_{i \in \mathcal{I}} L_i \tag{7}$$

In this case, each node will start with the same energy, namely E/N; then, the energy at each tier is

$$E_{i} = \frac{N_{i}E}{N} = \frac{(2i-1)E}{T^{2}}$$
(8)

Hence, (6) becomes

$$L_{i} = \frac{EP}{N\frac{T^{2} - i^{2}}{2i - 1}b_{s}e_{b} + Nb_{s}(e_{s} + e_{tx})}$$
(9)

Since the term  $(T^2 \cdot i^2)/(2i-1)$  is monotonously decreasing with *i* for  $1 \le i \le T$ ,  $L_i$  the lifetime of tier *i* is monotonously increasing with *i*. In other words, the lifetime of the network *L* is equal to the lifetime of the first tier:

$$L = L_1 = \frac{EP}{N(T^2 - 1)b_s e_b + Nb_s(e_s + e_{tx})}$$
 (10)

It is obvious that when the lifetime of the first tier has expired, the whole network has expired; the nodes in other tiers have still remaining energy. Using (8) and (9), the residual energy of other tiers  $E_{rest}$ :

$$E_{resi} = \frac{E}{T^2} [(2i-1) - \frac{(T^2 - i^2)e_b + (2i-1)(e_s + e_{tx})}{(T^2 - 1)e_b + e_s + e_{tx}}]$$
(11)

The above equations can also help us determine a

reasonable energy allocation for all nodes. One possible criterion is to let the all tiers have the same-targeted lifetime. Thus by balancing energy allocation, by using (2) and (6), we can maximize the network lifetime for a given fixed amount of energy E.

$$L_1 = L_2 = \dots = L_T \tag{12}$$

So we get the relation:

$$E_{i} = E_{1} \frac{(T^{2} - i^{2})e_{b} + (2i - 1)(e_{s} + e_{tx})}{(T^{2} - 1)e_{b} + e_{s} + e_{tx}}$$
(13)

Using (2) and (12), we get the relation:

$$E_{1} = \frac{((T^{2} - 1)e_{b} + e_{s} + e_{tx})E}{(4T^{3} - 3T^{2} - T)e_{b}/6 + T^{2}(e_{s} + e_{tx})}$$
(14)

If we allocate the energies at different tiers according as (13) and (14), we can maximize the network lifetime; all tiers' energies will be depleted at the same time.

#### 3.2. Design Problems

From the above, how to allocate the energies in different tiers is practical importance to the network designers. A reasonable alternative problem specification would provide, instead of the total budget E, the uniform battery level  $b_u$  for each node. In that case, the lifetime of the network would be independent of the total number of nodes N. The only topology relevant quantity would be the number of tiers T. It will be useful in nodes energy allocation.

According to (10), it is obvious that when the lifetime of the first tier and, therefore, the whole network has expired, the nodes in other tiers have still not expended their battery energy. The energy consumption for nodes at different tiers can be easily obtained from a consideration of (10) and (11). The hypothetical lifetime of tier *i* is  $L_i$ ; but, it is actually expending energy for only the lifetime *L* of the network. After that, communication stops, and there is no more expenditure of energy. Thus, the ratio of battery energy  $b_{e_i}$  actually consumed by a node in tier *i* to the uniform battery level  $b_u$  that it started with is:

$$C_{i} = \frac{b_{e_{i}}}{b_{u}} = \frac{L_{i}}{L_{i}} = \frac{\frac{T^{2} - i^{2}}{2i - 1}e_{b} + (e_{s} + e_{tx})}{(T^{2} - 1)e_{b} + e_{s} + e_{tx}}$$
(15)

We call the ratios  $C_i$  the ideal allocation ratios, because ideally each node would be allocated only the amount of battery it would actually consume during the lifetime. The farther out a node, the lower the fraction  $C_i$ of its battery that it has consumed. The unconsumed energy is wastage of the total energy budget. The broad design goal is to redistribute this energy budget nonuniformly so as to increase the lifetime of the whole

network.

#### 3.2.1. Problem 1

In many practice applications, a designer would have several known battery capacities to choose from e.g., AAA, AA, C and D. Then, the following problem can be formulated.

Given k available battery levels  $b_1 > b_2 > ... > b_k > 0$ , assign the battery level for each tier of nodes in a sensor network, such that the total battery budget E is minimized subject to maximizing the lifetime L of the network.

Two alternative flavors of the above problem can be articulated: (a) all nodes in any given tier are constrained to be assigned the same battery level, and (b) different nodes of the same tier can be assigned different battery levels. Thus, we are required to obtain one unique battery level  $b_i$  for tier *i* in Problem A; every node in this tier should be assigned this battery level. For problem B, we can provide more than one battery level for each tier accompanied by the proportions of the total number of nodes in that tier that should be assigned each battery level.

Below, we address problem 1A first.

The maximum network lifetime is obtained when all nodes have the maximum battery level  $b_1$ . Thus, the first tier of nodes should have batteries of capacity  $b_1$ , and the maximum lifetime of the network will be

$$L = L_1(b_1) \tag{16}$$

where, by  $L_i(b)$  we denote the lifetime of tier *i* (given by (6)) if each node in tier *i* initially has a battery capacity *b*.

The same maximum lifetime as in (16) may be obtained with a lower budget by assigning lower battery levels to nodes in the higher tiers (in tiers *i* with 1 < i < T); however, if the next smallest capacity size is too small, then the lifetime of the network would be reduced because nodes in higher tiers will deplete their batteries before the nodes in the first tier.

Given  $b_i$ , the ideal level of battery that each node of tier *i* should have is obviously  $C_ib_i$ , where  $C_i$  is the ideal allocation ratios given by (15), so that tier *i* will have exactly the same lifetime as tier 1. In short, for each tier *i*, in order to maximize the lifetime of the network and then conditionally minimize the total battery budget, we need to make sure that we assign the battery size  $b_j$  to tier *i* such that  $b_{j+1} < C_ib_1 < b_j$ . If the ideal level is less than the minimum level must be used in that tier instead; in this case some battery will indeed remain unconsumed at the end of the lifetime.

Now, if we consider problem 1B, we have the added flexibility of mixing the provided battery levels. First, we pick the highest battery level  $b_1$  for tier 1 as before: this maximizes the network lifetime. Now we predict the

ideal battery levels for each tier i as before using (15); but now, instead of picking the next highest available level from the available ones, we aim at attaining this ideal level exactly as the effective battery level by mixing the two provided battery levels in the appropriate ratio. In case the desired battery level is exactly equal to one of the provided battery levels, no mixing is needed.

How to mix the same tier nodes with the different battery levels? Usually, we would provide two battery levels to a tier. If we pick battery level  $b_1$  for tier 1, the ideal level of battery that each node of tier *i* should have been  $C_ib_1$ . We pick two battery levels  $b_i$  and  $b_{i-1}$  for tier *i* from available battery levels such that  $b_i < C_ib_1 < b_{i-1}$ . If we specified proportions  $f_1$  and  $f_2$  for the two levels with  $f_1+f_2=1$ , then the effective battery level of tier *i* would be given by  $b_if_1+b_{i-1}f_2$ , such that,

$$b_i f_1 + b_{i-1} f_2 = C_i b_1 \tag{17}$$

So we get the relation:

$$f_1 = \frac{b_{i-1} - C_i b_1}{b_{i-1} - b_i}$$
(18)

$$f_2 = \frac{C_i b_1 - b_i}{b_{i-1} - b_i}$$
(19)

## 3.2.2. Problem 2

In the above problem, minimizing the battery budget is a secondary goal of the optimal design. It is also possible that the battery budget may be strictly a constraint for the design. If the total cost of the batteries used for a particular network is upper bounded by a total battery budget E and the battery capacities are fixed and given, we can formulate the following problem:

Given k available battery levels  $b_1 > b_2 > ... > b_k > 0$ , and a total energy budget *E*, assign the battery level for each tier of nodes in a sensor network, such that the total lifetime *L* of the network is maximized.

As with Problem 1, we can conceive of two alternate problems, Problems 2A and 2B, in which the nodes in any tier are constrained to have the same battery level, or mixing is allowed, respectively.

We consider Problem 2A first.

As the solution to Problem 1A, we assign the battery levels at different tiers with the maximum lifetime unconstrained by total battery budget. Assume that  $L_i(b_j)$  is pre-computed lifetime for each tier *i* and each battery level  $b_j$ . Initialize "current network lifetime"  $L_c$  to  $L_1(b_1)$ . Whenever the total battery  $\sum_i N_i b_j$  becomes less than or

equal to the battery budget E, the algorithm is terminated.

If the current total battery exceeds E, repeatedly perform the following. Find the tier i such that the difference  $L_c-L_i(b_{j+1})$  is minimized. The tier i will be assigned the next lower battery level from the one it is currently assigned, that will result in the minimum reduction of the lifetime. Note that this minimum difference may well be zero at some iteration. Increment the current level for tier i by 1 and the update the network current lifetime to the new lifetime of tier i. Recalculate the total battery used. When the total battery first falls below the budget E, the algorithm will stop with an optimal solution.

Problem 2B turns out to be a modification of Problem 1B, where the total battery budget E is now a hard constraint. Accordingly, we take the following approach to solve it: first solve Problem 1B on the same parameters, ignoring the battery budget. If the total battery budget  $E_t$  of this solution does not exceed E, then this is the desired solution. If instead  $E_t > E$ , then obtain a new set of effective battery levels for each tier by scaling the battery levels by a factor of  $E/E_t$ : these are the new desired battery levels. A special case arises when some of these new desired battery level  $b_k$ . In this case the battery levels for the inner tiers have to be reduced to allow the outer tiers to have battery level  $b_k$ , further reducing the lifetime.

In considering these problems, any such algorithm would be executed not by the nodes themselves but offline during network design. In all realistic cases of deployment, most nodes will have positions in the appropriate annular regions, but some randomness will be introduced; lifetime will be then somewhat reduced from that achieved in the ideal case, but the ideal lifetime is a good indicator of actual lifetime in such cases.

#### 4. Simulation Results

To validate the results presented in the previous section, we decided to simulate the wireless sensor network in several scenarios. For simplicity reasons, we will assume that perfect scheduling is achieved at the MAC layer and routing layer. Any two nodes at a distance less than the transmission radius can communicate with no errors; any nodes at a distance larger than the transmission radius cannot communicate. We considered a network of N = 500 nodes, which corresponds to a five-hop route for the peripheral nodes (T=5), we assume that the total energy of the network is bounded by E=40000 joules in each case. Each node generates one data packet every minute; the size of pocket is 1024 bits.

Γ	abl	e 1.	Simu	lation	paran	neters
---	-----	------	------	--------	-------	--------

Parameter	Value
Transmitter circuitry, ete	2.34 µJ/bit
Receiver circuitry, erx	2.34 µJ/bit
Transmit one bit over one meter, eta	7.8 nJ/bit/m2
Sensing energy per bit, es	1.75µJ/bit
Bits sense per sensor, bs	1024bits

Table1 shows the values of the parameters in this sample circuit and the propagation environment [12].

#### 4.1. Base Case



Figure 3. The energy consumption at defferent tiers



Figure 4. The residual energies at different tiers



Figure 5. The allocating energies ratio at different tiers with maximizing the network lifetime

Using (5), we get Figure 3. It depicts the energy consumption at different tiers, because the first tiers relay the most data pockets, and consume the maximum energies. The energy consumption of the second tier, the third tier ... the fifth tier is ordinal decrease. The energy

consumption of the fifth tier is the minimum; nodes at the fifth tier don't relay other data pockets. If each node will start with the same energy (namely E/N), nodes at the first tier depleted their energies full out. Though nodes at other tiers have the residual energies, but their data pockets can't be relay to the sink, the network becomes disconnected; the lifetime of the network is the end. Using (11), we get Figure 4. The residual energies at different tiers aren't the same. The residual energies at fifth tier are the most; the residual energies at second tier are the least. For maximizing the network lifetime for a given fixed amount of energy E, we will balance energy allocation at different tiers. Using (13), we get Figure 5. It depicts the allocating energies ratio at different tiers. If we allocate the energies at different tiers according as the radio with Figure 5, we can maximize the network lifetime; all tiers' energies will be depleted at the same time. In the same parameters case, the network lifetime is 226.8427 seconds and 1494.8 seconds by the two energy allocation schemes, respectively. The lifetime of the network can be significantly improved.

## 4.2. The Effect of Multiple Battery Levels

On the assumption that we can provide 5 battery levels, they are 5178mAh, 2000 mAh, 1000 mAh, 500 mAh, 250 mAh, respectively.

For the base case, each node will be assigned with the same battery level 5178 mAh, the whole networks energy budget is 7767J.

For the problem 1A and maximizing network lifetime, the first tier of nodes should have batteries of capacity 5178mAh. According to the ideal allocation ratios  $C_i$ , the second tier, the third tier, the fourth tier, the fifth tier of nodes should have batteries of capacity 1656 mAh, 868.7 mAh, 472.1 mAh, 205.7 mAh, respectively, but as available battery levels limit, they have batteries of capacity 2000 mAh, 1000 mAh, 500 mAh, 250 mAh in application, respectively. Temporality, the whole networks energy budget is 1315.7J.

For the problem 1B and two battery levels, uniformity, the first tier of nodes should have batteries of capacity 5178mAh. Using formula (18) and (19), the optimum is achieved when 66% of the nodes have 2000 mAh batteries, and 34% of the nodes have batteries with 1000 mAh capacity in the second tier; 74% of the nodes have 1000 mAh batteries, and 26% of the nodes have batteries with 500 mAh capacity in the third tier; 89% of the nodes have 500 mAh batteries, and 11% of the nodes have batteries with 250 mAh capacity in the fourth tier; All nodes have 250 mAh batteries capacity in the fifth tier. Temporality, the whole networks energy budget is 1202.6J.

Using above battery deployment and the parameters in table 1, we can get figure 6 and Figure 7. They show the network lifetime of the different tiers. As can be seen, in three deployment schemes, the first tier lifetime is the shortest with 7.341 P intervals, so the whole networks lifetime is 7.341 P intervals. In base case deployment scheme, energy efficiency is 15.2%; there is much residual energy at different four tiers when the whole network has expired. In question 1A deployment scheme, energy efficiency is 89.6%. But available battery levels limit, there are some residual energies at different tiers when the whole networks has expired. In question 1B deployment scheme, energy efficiency is 98%. From the curve of question 1B, we can know the lifetime of nodes in the first, the second, the third, the forth is the same. Namely, the four tiers energy all has expired when the networks lifetime has expired.



Figure 6. The lifetime compare between the three schemes

There is some residual energy in fifth tier. Because available battery levels limit, we can't deploy the ideal appropriate battery level. If offering more battery levels, we can win the energy efficiency with 100%. All nodes in the different tiers have expired at the same time.

Just from the lifetime of networks, the problem 1B deployment scheme is ideality, but calculating the proportion of the mixing nodes is slightly complex. In the practice application, we can select one of the two schemes.



Figure 7. The lifetime compare between the two schemes

The network lifetime is further increased if more than two battery levels are considered, as seen in Figure 8. According to the scheme in problem 2, battery levels were used for the simulation. The same total energy budget 4000J was used for each simulation. Figure 8 show that the network lifetime will increase with the increase in the number of battery levels. In the same energy budget, the network lifetime with different battery levels will increase 188%, 356%, 487%, 559% relative to one with one battery level, respectively. This indicates that most of the increase in the lifetime of the network can be achieved with a relatively more number of battery levels.



Figure 8. The increase in network lifetime with the increase in the battery levels by the same energy budget at each battery level

## 4.3. Dependency on Number of Nodes

According to the scheme in problem 2, battery levels were used for the simulation. The same total energy budget 4000J was used for each simulation. Figure 9 shows the dependency of the network lifetime on the initial number of nodes. The density of the network was kept constant, so the area of the network was proportionally increased with the number of nodes.



Figure 9. Dependency of the network lifetime on the initial number of nodes for three battery levels (constant density)

The lifetime of the network decreases with the increase in the initial number of nodes. This is expected, as we already know a larger number of tiers results in more battery wastage.

#### 4.4. Dependency on Node Density

In this section, we study the effect of increasing the density of the network on the lifetime of the network, while keeping the network size constant. Figure 10 shows the dependency of the network lifetime on the node density. As can be seen in Figure 10, the node density has no influence on the network lifetime as long as it remains uniform. The reason is that the number of nodes in each tier will increase in the same proportion with the node density; hence, the number of load flows carried by each node does not change. The lifetime of the network remains constant even if the density of the network doubles or triples.



Figure 10. Dependency of the network lifetime on the network density for three battery levels (constant network area)

# 5. Conclusions and Future Work

In this paper, we have addressed a problem expected to occur in large, multi-hop wireless sensor networks. The nodes closer to the sink will die before the nodes at the periphery of the network. The main disadvantage of the expiration of the nodes close to the sink is that the network becomes disconnected while most of the nodes still have a considerable amount of energy left. To alleviate this undesirable effect, we proposed an energy allocation scheme, allocating different energy at different tiers by traffic. With this strategy, we have shown that the lifetime of the network can be significantly improved. Future work would explore similar issues, but MAC protocol will be considered.

# 6. References

[1] B.O. Priscilla Chen and E. Callaway, "Energy Efficient

System Design with Optimum Transmission Range for Wireless Ad-Hoc Networks," in Proceedings of ICC, vol. 2, pp. 945–952, 2002.

- [2] W. Ye, J. Heidemann, and D. Estrin, "Medium access control with coordinated adaptive sleeping for wireless sensor networks," IEEE/ACM Transactions on Networking, vol. 12, pp. 493–506, June 2004.
- [3] T. V. Dam and K. Langendoen, "An Adaptive Energy-Efficient MAC Protocol for Wireless Sensor Networks," Proceedings of first international conference an embedded networked sensor systems, pp. 171–180, 2003.
- [4] R. Min, M. Bhardwaj, N. Ickes, A. Wang, and A. Chandrakasan, "The hardware and the network: Totalsystem strategies for energy aware wireless micro sensors," in Proceedings of the IEEE CAS Workshop on Wireless Communications and Networking, (Pasadena, CA), September 2002.
- [5] S.C. Liu, "A Lifetime-Extending Deployment Strategy for Multi-Hop Wireless Sensor Networks," in Proceedings of IEEE Communication Networks and Services Research Conference, pp. 53–60, May 2006.
- [6] D. Wang, Y. Cheng, Y. Wang, and D.P. Agrawal, "Lifetime Enhancement of Wireless Sensor Networks by Differentiable Node Density Deployment," in Proceedings of IEEE International Conference on Mobile Ad hoc and Sensor Systems (MASS), pp. 546–549, October 2006.
- [7] X. Wu, G. Chen, and S.K. Das, "On the Energy Hole Problem of Non-uniform Node Distribution in Wireless

Sensor Networks," in Proceedings of IEEE International Conference on Mobile Ad-hoc and Sensor Systems (MASS), pp. 180–187, October 2006.

- [8] Y. Liu, H. Ngan, and L.M. Ni, "Power-Aware Node Deployment in Wireless Sensor Networks," in Proceedings of IEEE International Conference on Sensor Networks, Ubiquitous, and Trustworthy Computing (SUTC), pp.128–135, June 2006.
- [9] J. Lian, K. Naik, and G. Agnew, "Data Capacity Improvement of Wireless Sensor Networks Using Non-Uniform Sensor Distribution," International Journal of Distributed Sensor Networks, vol. 2, no. 2, pp. 121–145, April 2006.
- [10] J.J. Lee, B. Krishnamachari, and C.C.J. Kuo, "Impact of heterogeneous deployment on lifetime sensing coverage in sensor networks," in Proceedings of IEEE Conference on Sensor and Ad Hoc Communications and Networks (SECON), pp. 367–376, October 2004.
- [11] W. R. Heinzelman, A. Chandrakasan, and H. Balakrishnan, "Energy efficient communication protocol for wireless microsensor networks," in Proceedings of the Hawaii international Conference on System Sciences, January 2000.
- [12] R. Min, M. Bhardwaj, N. Ickes, A. Wang, and A. Chandrakasan, "The hardware and the network: Totalsystem strategies for energy aware wireless micro sensors," in Proceedings of the IEEE CAS Workshop on Wireless Communications and Networking, (Pasadena, CA), September 2002.



# Beacon-driven Leader Based Protocol over a GE Channel for MAC Layer Multicast Error Control

Zhao LI<sup>1</sup>, Thorsten HERFET<sup>2</sup>

<sup>1</sup> Student Member IEEE, USTC CS China & Saarland University, Germany <sup>2</sup> Senior Member IEEE, Saarland University, Germany E-mail: {li, herfet}@nt.uni-saarland.de

# Abstract

In wireless networks current standard MAC layer protocols don't provide any error correction scheme for broadcast/multicast. In this paper, we enhance a Leader Based Protocol (LBP) and propose a Beacon-driven Leader Based Protocol (BLBP) for the MAC layer multicast error control. To guarantee a very low Packet Loss Ratio (PLR) under strict delay constraints for video multicast over a Gilbert-Elliott (GE) channel, we analyze BLBP and compare it with LBP and different application layer multicast error control schemes via simulation experiments. Both the theoretical analysis and simulation results show that BLBP can correct nearly all the errors for all receivers in the MAC layer and is more efficient than LBP. BLBP is also more efficient than the application layer Automatic Repeat request (ARQ) scheme and the total multicast delay is much shorter. BLBP is very good for real-time multicast applications with strict delay constraints.

Keywords: BLBP, Multicast Error Control, Gilbert-elliott Channel

# 1. Introduction

With the rapid development of wireless networking technologies, it is becoming possible to supply wireless terminal users not only with data connections, but also with real-time communication services. The emerging real-time multicast applications in wireless networks include the local distribution of High Definition TV (HDTV) and Digital Video Broadcasting (DVB) [1], video-on demand, video conferencing, gaming, local distribution, VoIP. IPTV, Internet-Radio P2P broadcasting, etc. Most of these applications require a very low Packet Loss Ratio (PLR) under strict delay constraints.

The characteristics of wireless network can be summarized as a bandwidth variation and terminal heterogeneity plus a high degree of packet losses. It is known that the Gilbert-Elliot (GE) channel [2,4] with a 2-state Markov model is a good approximation for the packet loss model in wireless channels. However, the current standard MAC layer protocols don't provide any error correction scheme for broadcast/multicast. Hence the multicast error is controlled in the application layer. The existing application layer multicast error control schemes include automatic repeat request (ARQ), forward error correction (FEC) and hybrid error correction (HEC) [5–10]. Unfortunately the total multicast delays of the application layer schemes are always high and sometimes do not satisfy the strict application delay constraints, or these schemes are not efficient when the delay constraints are short.

Compared with application layer schemes, MAC layer multicast error control schemes recover the multicast loss locally and lead to much shorter delays. Currently, very few reliable MAC layer multicast schemes, such as Leader Based Protocol (LBP) [11], have been proposed for IEEE 802.11 based wireless networks. LBP elects one of the multicast group receivers as the leader. On erroneous reception of a data frame, the leader does not send an acknowledgement (ACK), prompting a retransmission. On erroneous reception of the data frame at the non-leader receivers, LBP allows negative acknowledgements (NACKs) from these receivers to collide with the ACK from the leader, thus destroying the ACK and prompting the sender to retransmit the data frame. We refer to this ACK/NACK iam as JACK.

However, LBP suffers from two main problems: First, when the entire data frame is lost, the non-leader

receivers can not reply NACKs because they don't know when or how to send them, as the destination is unknown especially when Request-To-Send (RTS) and Clear-To-Send (CTS) is not used for small data frames. As a result, LBP is not reliable for the non-leader receivers. Second, LBP has poor performance when the channel error rates are high. The non-leader receivers send NACKs whenever the received frame is in error, regardless of whether this erroneous frame has been received correctly before or not. This is because the receivers in LBP can not access the data frame sequence number before the frame is received, as there is no such field in the structure of RTS/CTS frames for multicast. So the sender has to retransmit until all receivers receive the data frame correctly at the same time. There are a lot of unnecessary transmissions; hence LBP is not efficient particularly for lossy channels.

In this paper, we enhance LBP and propose a Beacon-driven Leader Based Protocol (BLBP). The sender sends a beacon frame before the data frame to lead the non-leader receivers to set timers and to announce the sequence number of the following data frame. BLBP solves the problems of LBP well. Each of the non-leader receivers can send feedback when the timer times out. Both the leader receiver and non-leader receivers can send ACK and NACK respectively based on sequence check, hence avoids the unnecessary transmissions in LBP. To guarantee a very low PLR under strict delay constraints for video multicast over a GE channel, we analyze BLBP and compare it with LBP and two application layer multicast error control schemes via simulation experiments. One is an ARQ based scheme from [8], called HEC-PR, which combines a NACK based ARQ scheme with a packet repetition (PR) technique. The other one is a nearly optimal application layer multicast error control scheme called Hybrid ARQ (HARQ) Type I [9,10], which combines FEC and NACK based ARQ scheme together.

The remainder of this paper is organized as follows. Section 2 presents the related work. In section 3, we examine why it is necessary to correct multicast errors in the MAC layer. We describe BLBP in section 4 and analyze its performance over a GE channel in section 5. And in section 6, we evaluate the performance of BLBP and compare it with LBP, HEC-PR and HARQ Type I via simulation experiments. We conclude in section 7.

# 2. Related Work

For the application layer multicast error control, many authors [5,8] studied the ARQ based schemes and concluded that when combined with feedback suppression and other accessorial techniques, ARQ is effective to repair multicast packet losses for small groups with low error rates. However, the application layer ARQ always take a long time and they are not efficient at high error rates and with large numbers of receivers due to feedback implosion and the limitation to scale.

Another technique commonly used to handle losses for multicast in the application layer is FEC, whereby redundant information in the data stream enables the receiver to correct losses without contacting the sender. Rizzo [6] studied the feasibility of software encoding/decoding for packet-level FEC. A (n, k) block erasure code converts k source packets into a group of ncoded packets, such that any k of the encoded packets can be used to reconstruct the k source packets. Usually, the first k packets in each group are identical to the original k data packets; the remaining n-k packets are referred to as parity packets. The advantage of using block erasure codes for wireless multicasting is that a single parity packet can be used to correct independent single-packet losses among different receivers.

The integrated FEC/ARQ schemes or any other kinds of combination of more than one error control schemes are referred to as HEC schemes in this paper. Previous works [7-10] indicate that HEC schemes are much more efficient for recovering data packets than the schemes with either FEC or ARQ alone. We consider the HEC-PR scheme from [8], which combines a NACK based ARQ scheme and a packet repetition technique. The number of feedback/retransmission rounds and the number of packet repetitions in each round are adapted to the network condition. HEC-PR is actually an ARQ based scheme without FEC coding. We also consider HARQ Type I from [9,10], in which the sender sends a certain amount of parity packets using FEC following the original k data transmissions. If the loss rate obtained after reconstruction at the receiver is still too high, ARQ is used to retransmit more parity packets. Tan [10] developed formulas to optimize the performance of HARQ Type I while guaranteeing the required PLR under strict delay constraints. HARQ Type I is a nearly optimal application layer multicast error control scheme.

However, these application layer multicast schemes always take long multicast delays or they are not efficient when the delay constraints are short. So we study MAC layer multicast error control schemes to support real-time multicast with strict delay constraints. For the reliable MAC layer multicast, besides LBP [11], Tourrihes [12] proposed a robust broadcast using a collision detector to inform the sender whether the broadcast packet is successful or not. However, this scheme can not guarantee the reliability of multicast transmissions because the feedbacks are only from the detector instead of all receivers themselves. Gupta et al. [13] proposed a tone-based solution for multicast in both infrastructure and ad-hoc 802.11 networks. They use dual busy tones to simulate NACKs or Negative CTS (NCTS). Although this scheme is good to detect and correct the multicast errors, it requires an additional

channel for the tone, which is not always feasible in practice.

Our BLBP enhances LBP with a beacon frame to lead the non-leader receivers to set timers and to announce the sequence number of the following data frame. BLBP avoids the problems of LBP and is more efficient.

# 3. Motivation

The emerging real-time multicast applications in wireless networks (such as wireless HDTV, DVB, game, video conference) require strict delay constraints. Error recovery based on application layer ARQ is suboptimal because the end-to-end application layer feedback and retransmission take a too long time due to application layer protocol waiting, MAC layer queuing, hardware handling, etc. Moreover, the application layer ARQ based schemes are not efficient when the error rates are high and the numbers of receivers are large due to feedback implosion and the limitation to scale. The FEC coding based application layer schemes can satisfy the strict delay constraints but they are not efficient when the delay constraints are very strict particularly for small multicast groups. And the FEC based ones are not adaptive to the heterogeneity of receivers because the code has to be set based on the receiver with the worst channel condition.

Compared with application layer schemes, MAC layer multicast error control schemes take a much shorter time due to the faster feedback and retransmission. Due to the JACK scheme, BLBP and LBP even achieve complete feedback suppression. So the MAC layer multicast error control schemes are very good for real-time multicast applications. For non-real-time multicast applications such as file dissemination and shared whiteboards, reliable MAC layer multicast saves time as well as both network and end-system resources.

Moreover, for multi-hop multicast with wired network and wireless LAN as the last hop, the need for additional transmissions due to errors in the wireless LANs puts unnecessary processing burden on the original remote sender. These additional transmissions go over the entire wired multicast tree and also the wireless links, taking a long time, wasting bandwidth and also leading to processing of unwanted redundant retransmissions at those receivers which might have already received the packet. The similar thing also happens in multi-hop wireless networks such as wireless mesh networks and wireless ad hoc networks. If the access points (AP) (or base stations) were to take the responsibility of supplying retransmissions rather than the original sender, then the load of supplying retransmission gets distributed across access points and takes a shorter time. The total error correction cost will be much shorter and it is easier to guarantee the final PLR under strict delay constraints in the application layer.

#### 4. Main Scheme of BLBP

The MAC layer reliable multicast BLBP requires a slight modification to the IEEE 802.11 MAC layer protocols. As mentioned earlier, 802.11 DCF (Distributed Coordination Function) unicast – assumed RTS/CTS is switched on to solve the hidden terminal problem – is more reliable than broadcast/multicast, because unicast uses RTS/CTS signaling and ACK/retransmission scheme in the MAC layer and broadcast/multicast does not.

BLBP enhances LBP with a MAC control frame called beacon shown in Figure 1. Besides the same fields in RTS/CTS frames, such as frame control header, transmission duration, receiver address (RA), transmitter address (TA) and frame check sequence (FCS), the beacon frame also includes the sequence number of the following data frame. The use of the beacon frame is to lead the non-leader receivers to set timers and to announce the sequence number of the following data frame.

2	2	6	6	2	4
Frame Control	Duration	RA	ТА	Sequence Control	FCS

Figure 1. The format of the beacon frame

The main scheme of BLBP is shown in Figure 2. A receiver is selected as the leader for the multicast group. The AP first sends a RTS frame to all receivers, and only the leader receiver transmits a CTS frame in reply to the AP. The AP is then assured that the channel is granted and starts the transmission of the beacon frame with the sequence number of the following data frame. On receipt of the beacon frame, each of the non-leader receivers sets a timer according to the beacon frame. The AP then transmits the data frame following the beacon frame. The leader receiver replies an ACK frame if the data is correct or it has already got the data based on sequence check, or does nothing otherwise. When the timer times out, each non-leader receiver replies a NACK if the data is error and it has not received it correctly yet based on sequence check, or does nothing otherwise. Then if the AP receives an ACK, this transmission is done. Otherwise, the AP repeats the whole procedure and retransmits again until the number of times is beyond the retransmission limit. For example, in the retransmission phase in Figure 2, although this time the data frame is lost, the leader receiver still replies an ACK because it knows this data frame has been received correctly already in the first transmission, thanks to the beacon frame.

For the determination of a leader, we use a scheme

from LBP [11]: The first receiver that joins the multicast group acts as the leader. The group is cancelled if there is no leader. The other group members can rebuild/rejoin the multicast group if necessary when time out. Please note that it is possible to reduce the amount of control traffic flow for leader election purposes when a higher layer group management protocol like IGMP (Internet Group Management Protocol) is running above the link layer [11]. The leader does not affect the performance because BLBP supplies fair service for all receivers.



Figure 2. Main scheme of BLBP (Ri denotes receiver i)

BLBP solves the problems of LBP well. All the nonleader receivers can send feedbacks when the timers time out. Both the leader receiver and non-leader receivers send ACK and NACK respectively based on sequence check thanks to the beacon frame, hence it avoids the unnecessary transmissions in LBP. Clearly, BLBP can correct all the errors for all receivers due to the ACK/NACK feedback and retransmission in the MAC layer. BLBP is even more efficient than the application layer ARQ schemes because BLBP suppresses the multiple feedbacks into just a JACK. However, the loss of beacon frames will decrease the performance (reception rate) of the non-leader receivers. Fortunately, the beacon frames are much more reliable (nearly error free) than data frames because they are much smaller and are transmitted using the lowest data rate, like other control frames in 802.11 (RTS, CTS, ACK). Moreover, due to RTS/CTS signaling, the beacon frames also avoid collision loss.

Please also note that BLBP can run without RTS/CTS exchanges for small data frames just like 802.11 DCF unicast. Although our discussion is in the context of 802.11 DCF, BLBP is actually applicable to all ACK/retransmission based MAC protocols, such as 802.11 PCF (Point Coordination Function) etc.

## 5. Performance Analysis

In this section we first introduce the GE channel model and then analyze the theoretical performance of BLBP over the GE channel model with both temporal error correlation and spatial error correlation.

#### **5.1. GE Channel Model**

The GE channel model is a two-state Markov chain shown in Figure 3. In the Good state (G) errors occur with (low) probability  $P_G$  while in the Bad state (B) they occur with (high) probability  $P_B$ .

The errors occur in clusters or bursts with relatively long error-free intervals (gaps) between them. The state transition is summarized by its transition probability matrix in formula (1).



Figure 3. GE channel model

$$P = \begin{bmatrix} \beta & 1 - \beta \\ 1 - \alpha & \alpha \end{bmatrix}$$
(1)

To reflect most reasonable choices for real scenarios, it is assumed that  $P_G=0$  and  $P_B=1$ . This model is always referred to as the simplified GE model. Our analysis and simulation experiments in the following sections are based on the simplified GE model.

The occupancy times for state B and G are both geometrically distributed with respective means  $(1-\alpha)^{-1}$  and  $(1-\beta)^{-1}$ , which are also called as the expected error burst length and the expected error free length respectively. The steady state probabilities of being in states G and B are  $\pi_G = (1-\alpha)/(2-\alpha-\beta)$  and  $\pi_B = (1-\beta)/(2-\alpha-\beta)$  respectively. So the average packet loss rate produced by the GE channel model is

$$p = P_G \pi_G + P_B \pi_B = \frac{P_G (1 - \alpha) + P_B (1 - \beta)}{(1 - \alpha + 1 - \beta)}$$
(2)

For the simplified GE channel model, the PLR will be

$$p = \frac{(1-\beta)}{(1-\alpha+1-\beta)}$$
(3)

Following [14], the variance of the error symbol (or packet) *X* is  $\sigma^2 = E(X - \overline{p})^2 = p(1-p)$ . So we get the correlation of two consecutive error symbols *X*<sub>1</sub> and *X*<sub>2</sub>:

$$\tau = \frac{E((X_1 - p)(X_2 - p))}{\sigma^2} = \alpha + \beta - 1$$
 (4)

which is also referred to as the temporal error correlation.

Finally, the two parameters of the simplified model ( $\alpha$  and  $\beta$ ) can be expressed in the terms of the more meaningful quantities p and  $\tau$  by solving formulas (3) and (4). These yields

$$\alpha = p + \tau(1 - p) \tag{5}$$

$$\beta = (1-p) + \tau p \tag{6}$$

The transition probability matrix then becomes

$$P = \begin{bmatrix} 1 - p(1 - \tau) & p(1 - \tau) \\ (1 - p)(1 - \tau) & 1 - (1 - p)(1 - \tau) \end{bmatrix}$$
(7)

And the *I*-step transition probability matrix is:

$$P^{I} = \begin{bmatrix} 1 - p(1 - \tau^{I}) & p(1 - \tau^{I}) \\ (1 - p)(1 - \tau^{I}) & 1 - (1 - p)(1 - \tau^{I}) \end{bmatrix}$$
(8)

Now we compute P[a,b], the probability of *a* errors in a sequence of *b* symbols following [14]. Let  $P_G[a,b]$ be the probability of *a* errors in *b* transmissions with the channel ending in state G. Similarly, let  $P_B[a,b]$  be the probability of *a* errors in *b* transmissions with the channel ending in state B. Then

$$P[a,b] = P_{G}[a,b] + P_{B}[a,b]$$
(9)

For  $b = 1, 2, 3 \dots$  and  $a = 0, 1, 2 \dots b$ , assuming the simplified GE channel, then

$$P_{G}[a,b] = P_{G}[a,b-1]\beta + P_{B}[a,b-1](1-\alpha)$$
(10)

$$P_{B}[a,b] = P_{B}[a-1,b-1]\alpha + P_{G}[a-1,b-1](1-\beta)$$
 (11)

The initial conditions for the recursion are

$$P_G[0,0] = (1-\alpha)/(2-\alpha-\beta)$$
 (12)

$$P_{B}[0,0] = (1-\beta)/(2-\alpha-\beta)$$
 (13)

and  $P_G[a,0] = P_B[a,0] = 0$  for  $a \neq 0$ . Note that with these initial conditions, all numerical values computed will be steady state results.

# **5.2. BLBP over the GE Channel Model with** Temporal Error Correlation

Now we analyze BLBP over the GE channel model. As in most referenced papers, it is assumed that the MAC layer control frames (RTS/CTS/Beacon/ACK) are error free from the error model. This also makes sense in practice because the control frames are very small and are sent in the lowest data rate, and hence they are more reliable than data frames. We first consider only the temporal error correlation. In other words, it is assumed that the error events at different receivers are independent. To be clear, we show the used symbols here.

- *P*: The original packet error rate for all receivers;
- *R*: The number of receivers;
- *m*: The retransmission limit;

- *N*: The total number of transmissions required to transmit a multicast packet correctly to all the *R* receivers;
- *N<sub>r</sub>*: The number of transmissions required for receiver *r* to receive a packet correctly;
- *PLR<sub>target\_mac</sub>*: The PLR target in the MAC layer;
- *D<sub>target\_mac</sub>*: The Delay target in the MAC layer.

First we consider the final PLR for receiver r, shown in formula (14).

$$PLR(r) = P[m+1,m+1] = p\alpha^{m}$$
 (14)

About the determination of the retransmission limit m, there are two constraints, the PLR constraints and the delay constraints which are shown as follows:

$$p\alpha^m \le PLR_{target mac} \tag{15}$$

$$m * T_{BLBP} \le D_{t \arg et_mac} \tag{16}$$

where  $T_{BLBP}=T_{CC} + T_{RTS} + T_{CTS} + T_{BEACON} + T_{DATA} + T_{ACK} + DIFS + 4SIFS$  is the time of one transmission in BLBP.  $T_{CC}$  denotes the channel contention time which can be calculated theoretically following [15] or by measurements in practice.  $T_{RTS}$ ,  $T_{CTS}$ ,  $T_{BEACON}$ ,  $T_{DATA}$ , and  $T_{ACK}$  are the transmission times of frames RTS, CTS, BEACON, DATA and ACK respectively. *DIFS* denotes the Distributed Inter Frame Space while *SIFS* is the Short Inter Frame Space. Note that the PLR target and the delay target may not be satisfied at the same time sometimes, especially when the delay constraint is too strict. We will explore this further in simulation experiments.

Now we consider the expected number of transmissions for one multicast data packet. The probabilities that  $N_r \le n$ ,  $N_r = n$  and  $N_r > n$  are shown in formulas (17), (18) and (19) respectively.

$$P[N_r \le n] = 1 - P[n, n]$$
  
= 
$$\begin{cases} 1 - p\alpha^{n-1}, & n = 1, 2...m + 1 \\ 0, & n = 0 \end{cases}$$
 (17)

$$P[N_r = n] = P[N_r \le n] - P[N_r \le n-1]$$
  
= 
$$\begin{cases} p\alpha^{n-2} - p\alpha^{n-1}, & n = 2, 3...m+1\\ 1 - p, & n = 1 \end{cases}$$
 (18)

$$P[N_r > n] = P[n,n] = \begin{cases} p\alpha^{n-1}, & n = 1, 2...m \\ 1, & n = 0 \end{cases}$$
(19)

So we get the expected number of transmissions for one multicast data packet required for receiver *r*:

$$E(N_r) = \sum_{n=0}^{m} P[N_r > n] = 1 + \frac{p(1 - \alpha^m)}{(1 - \alpha)}$$
(20)

Next the probabilities that  $N \le n$ , N = n and N > n are shown in formulas (21), (22) and (23) respectively.

$$P[N \le n] = \prod_{r=1}^{R} P[N_r \le n]$$

$$\left[ (1 - nq^{n-1})^R \quad n = 1, 2, m+1 \right]$$

$$=\begin{cases} (1-p\alpha - 1)^{R}, & n=1,2...,m+1\\ 0, & n=0\\ P[N=n] = P[N \le n] - P[N \le n-1] \end{cases}$$
(21)

$$=\begin{cases} \left(1-p\alpha^{n-1}\right)^{n} - \left(1-p\alpha^{n-2}\right)^{n}, & n=2,3...m+1\\ \left(1-p\right)^{R}, & n=1 \end{cases}$$
(22)

$$P[N > n] = 1 - P[N \le n]$$
  
= 
$$\begin{cases} 1 - (1 - p\alpha^{n-1})^{R}, & n = 1, 2...m \\ 1, & n = 0 \end{cases}$$
 (23)

Finally we get the expected number of transmissions for one multicast data packet for all receivers, shown in formula (24).

$$E[N] = \sum_{n=0}^{m} P[N > n]$$
  
= 1 +  $\sum_{n=1}^{m} \left( 1 - \left( 1 - p \alpha^{n-1} \right)^{R} \right)$  (24)

Similarly, we get the expected number of transmissions for one multicast data packet for all receivers in LBP, shown in formula (25).

$$E_{LBP}[N] = \sum_{n=0}^{m} \left(1 - \left(1 - P[1,1]\right)^{R}\right)^{n}$$
$$= \sum_{n=0}^{m} \left(1 - \left(1 - p\right)^{R}\right)^{n}$$
(25)

And the redundant information (RI) of BLBP is:

$$RI = E[N] - 1 \tag{26}$$

Similar as in [16], among pure ARQ based multicast error control schemes, BLBP needs the minimum number of transmissions (shown in formula (24)) to let all receivers receive the packet correctly. There are no unnecessary transmissions due to the sequence check scheme and the complete feedback suppression based on the JACK scheme. We will explore this further in simulation experiments by comparing BLBP with LBP and an application layer pure ARQ scheme.

# 5.3. BLBP over the GE Channel Model With Both Temporal and Spatial Error Correlation

Now we consider BLBP over a GE channel model with both the temporal error correlation and the spatial error correlation. The error events at different receivers are in a certain correlation. We assume two kinds of error events, the error event at the sender which leads to the correlated packet loss among all receivers and the error events at different receivers which lead to independent packet losses. Some new symbols are shown as follows.

- *P<sub>m</sub>*: The error rate caused at receivers (called input error), independent of each other;
- $P_{out}$ : The error rate caused at the sender (called output error), which leads to correlated loss covering all receivers;  $p = p_{out} + (1 p_{out})p_{in}$ ;
- $P_{in} = \begin{bmatrix} \beta_{in} & 1 \beta_{in} \\ 1 \alpha_{in} & \alpha_{in} \end{bmatrix}$ : The transition probability matrix

of the input errors at all receivers;

•  $P_{out} = \begin{bmatrix} \beta_{out} & 1 - \beta_{out} \\ 1 - \alpha_{out} & \alpha_{out} \end{bmatrix}$ : The transition probability

matrix of the output error at the sender;

•  $\lambda$ : The spatial error correlation among different receivers,  $\lambda = p_{out}/p$ .

Combining the output error and the input error, we can compute the total error model, which is a 4-states Markov chain, shown as follows.

$$P = \begin{bmatrix} \beta_{out} \beta_{in} & \beta_{out} (1 - \beta_{in}) & (1 - \beta_{out}) \beta_{in} & (1 - \beta_{out}) (1 - \beta_{in}) \\ \beta_{out} (1 - \alpha_{in}) & \beta_{out} \alpha_{in} & (1 - \beta_{out}) (1 - \alpha_{in}) & (1 - \beta_{out}) \alpha_{in} \\ (1 - \alpha_{out}) \beta_{in} & (1 - \alpha_{out}) (1 - \beta_{in}) & \alpha_{out} \beta_{in} & \alpha_{out} (1 - \beta_{in}) \\ (1 - \alpha_{out}) (1 - \alpha_{in}) & (1 - \alpha_{out}) \alpha_{in} & \alpha_{out} (1 - \alpha_{in}) & \alpha_{out} \alpha_{in} \end{bmatrix}$$

$$(27)$$

For the convenience of analyzing, we use a GE channel model to approximate the total error model. This is also confirmed by the simulation experiments. The total GE channel model is

$$P = \begin{bmatrix} \beta & 1 - \beta \\ 1 - \alpha & \alpha \end{bmatrix}$$
(28)

where  $\alpha = 1 - (1 - \beta)(1 - p) / p$  and  $\beta = \beta_{out} \beta_{in}$ .

So we get the final PLR for all receivers, shown in formula (29).

$$PLR(r) = P[m+1, m+1] = p\alpha^{m}$$
 (29)

Now we consider the expected number of transmissions for one multicast data packet for all receivers. The probability that N > n is shown as follows.

$$P[N > 0] = 1 \tag{30}$$

$$P[N > 1] = p_{out} + (1 - p_{out}) \left( 1 - (1 - p_{in})^{R} \right)$$
(31)

$$P[N > n] \approx P_{out}[n, n] + (1 - p_{out}) \left( 1 - (1 - (p_{in} / p)P[n, n])^{R} \right) \\ + \sum_{i=1}^{n-1} \left( P_{out}[n - i, n - i](1 - \alpha_{out}) \left( 1 - (1 - (p_{in} / p)P[i, i])^{R} \right) \right) \\ = p_{out} \left( \alpha_{out} \right)^{n-1} + (1 - p_{out}) \left( 1 - (1 - p_{in}\alpha^{n-1})^{R} \right) \\ + \sum_{i=1}^{n-1} \left( p_{out} \left( \alpha_{out} \right)^{n-i-1} (1 - \alpha_{out}) \left( 1 - (1 - p_{in}\alpha^{i-1})^{R} \right) \right)$$
(32)

Finally the expected number of transmissions for one multicast data packet for all receivers can be calculated as in formula (33).

$$E[N] = \sum_{n=0}^{m} P[N > n]$$

$$\approx 1 + \sum_{n=1}^{m} \left( p_{out} (\alpha_{out})^{n-1} + (1 - p_{out}) (1 - (1 - p_{in} \alpha^{n-1})^{R}) + \sum_{i=1}^{n-1} (p_{out} (\alpha_{out})^{n-i-1} (1 - \alpha_{out}) (1 - (1 - p_{in} \alpha^{i-1})^{R})) \right)$$
(33)

Note that it is direct and simple to calculate the average channel holding time of BLBP from the expected number of transmissions in practice. We will evaluate the performance of BLBP and compare BLBP with LBP and different application layer multicast error control schemes by simulation experiments in the following section.

# 6. Performance Evaluation

In this section, we first evaluate the performance of BLBP and confirm the theoretical results by simulation experiments. Then we compare BLBP with the application layer multicast error control schemes HEC-PR and HARQ Type I. The metrics used for evaluation include the average number of transmissions, the maximum multicast delay and the total RI. Consistent with many references [8,10,14], we also consider the redundant transmission only as RI.

Table 1. Application targets and parameters

PLR Requirement	1e-6
Delay Constraint	20-100ms
RTP Payload Length	1316Bytes
Multicast load interval	2.5ms
RTT	≈ 3.5ms
Original Error Rate	$\leq 10\%$
Packet sent	40-100e6

 Table 2. The retransmission limit and the temporal error correlation (PLR constraint 1e-6)

Error	The temporal error correlations					
Rates	0.0	0.1	0.2	0.3	0.4	0.5
0.05	4	6	8	10	13	17
0.10	6	7	10	12	15	20

We conduct our simulation study using ns-2 and implement BLBP based on the 802.11e simulation model from [17]. All client nodes are one hop to the AP and at most two hops to each other. We use IEEE 802.11a parameters to model the physical layer. The data rate we choose is 24Mbps. The first receiver that joins the multicast group acts as the leader. The application targets and parameters are presented in Table 1. Note that the PLR target (1e-6) is very strict. We use this final PLR target as the MAC layer PLR target. The total payload length in the MAC layer is 1356 bytes, and there is no fragmentation in the MAC layer or the network layer. The retransmission limit of BLBP is determined following only the PLR constraint shown in formula (15). Some examples are shown in Table 2. We guarantee the PLR target first and explore the total multicast delays of BLBP in different network scenarios and channel conditions.

The application layer multicast error control schemes HEC-PR and HARQ Type I are implemented based on the real-time transport protocol (RTP) [18]. We use unicast for the feedback in HEC-PR and HARQ Type I instead of broadcast because unicast is more reliable. The simplified GE channel model is implemented in the physical layer, but it is used only for data frames. The MAC control frames (RTS/CTS/Beacon/ACK) are error free from the error model. (The control frames also may be lost because they might collide with the background traffic.) This also makes sense in practice because the control frames are very small and are sent in the lowest data rate, and hence they are more reliable than data frames.

First we compare the average numbers of transmissions for BLBP and LBP in different channel conditions. The results are shown in Figure 4. The simulation result and the theoretical result of BLBP match very well. As expected, BLBP is more efficient than LBP particularly when the error rates are high and the numbers of receivers are large. This is because BLBP allows all receivers to send feedback based on sequence check thanks to the beacon frame and LBP can not.

Then we explore the effect of the temporal error correlation. The average numbers of transmissions and the maximum multicast delays of BLBP with different temporal error correlations are shown in Figure 5 and Figure 6 respectively. About the expected number of transmissions, we can see that the theoretical analysis and the simulation results match very well. The temporal error correlation affects the average number of transmissions very much. However, the multicast delays are still very low even when the temporal error correlations are high, thanks to the fast ACK and retransmission in the MAC layer. So BLBP is very good for real-time multicast applications with strict delay constraints.

Figure 7 shows the average numbers of transmissions for BLBP with different spatial error correlations. The theoretical analysis and the simulation result match very well. BLBP can take full advantage of the spatial error correlation because of the complete feedback suppression thanks to the JACK scheme.

Then we compare BLBP with the application layer multicast error control schemes. The temporal error correlation and the spatial error correlation are set to 0.10 and 0.20 respectively to simulate a near realistic channel condition. First we compare BLBP with HEC-PR in different channel conditions. The total RIs and maximum multicast delays are shown in Figure 8 and Figure 9 respectively. When the error rates are high, we can see BLBP is much more efficient than HEC-PR which is a pure ARQ based scheme (without FEC coding). This is because BLBP suppresses the multiple feedbacks into just a JACK and it is more efficient than the application layer feedback and retransmission suppression in HEC-PR. The results also show that the multicast delays of BLBP are much shorter than the delays in HEC-PR. This is due to the fast MAC layer ACK/retransmission in the MAC layer. Moreover, because of the much longer delays among each transmission or retransmissions, HEC-PR isn't much affected by the temporal error correlation.

Figure 10 shows the total RIs of BLBP, HEC-PR and HARQ Type I under different delay constraints. As shown in Figure 9, HEC-PR takes a long time and hence only satisfies long delay constraints. Both BLBP and HARQ Type I can satisfy all the delay constraints from 20ms to 100ms. We can see that BLBP is more efficient than HARQ Type I when the delay constraints are short and the numbers of receivers are small. This is because that HARQ Type I has to switch to pure FEC scheme (no ARQ) when the delay constraint are very short, hence it is not efficient.

Finally we compare the total RIs of BLBP, HEC-PR and HARQ Type I with the heterogeneity of receivers. The error rate for receiver 1 is variable and all the other receivers have a fixed error rate 0.01. Figure 11 shows the result. We can see that both BLBP and HEC-PR are much more efficient than HARQ Type I because HARQ Type I has to set FEC code and other parameters according to the receiver with the worst channel condition but BLBP and HEC-PR are more adaptive. Moreover, due to the effect of the temporal error correlation (shown in Figure 5), here BLBP is a little less efficient than HEC-PR.



Figure 4. The expected number of transmissions with different error rates ( $\lambda = 0$   $\tau = 0$  for all receivers)

#### 7. Conclusions

In this work, we enhance LBP and propose BLBP for the MAC layer multicast error control in wireless networks. The use of the beacon frame is to lead the non-leader receivers to set timers and to announce the data frame sequence. On erroneous reception of a data frame that has not been correctly received before, the leader does not send an ACK, prompting a retransmission. On erroneous reception (The timer times out) of the data frame that has not been correctly received before, the non-leader receivers send NACKs to collide with the potential ACK from the leader, thus prompting the AP to retransmit the packet. To guarantee a very low PLR under strict delay constraints for video multicast over a GE channel with both the temporal error correlation and the spatial error correlation, we analyze BLBP and evaluate its performance via simulation experiments.

Both the theoretical analysis and simulation results show that BLBP can correct nearly all the errors for all receivers in the MAC layer. BLBP needs the minimum number of redundancy transmissions among all pure ARQ based schemes. BLBP is more efficient than the application layer ARQ schemes and the total multicast delay is much shorter. BLBP is even more efficient than the best application layer multicast error control scheme when the delay constraints are short or with the heterogeneity of receivers. BLBP is very good for the real-time multicast applications with strict delay constraints, especially for small groups.

In the future, we plan to extend BLBP with FEC coding and further improve the QoS of multicast cross the application layer and MAC layer in wireless networks.



Figure 5. The expected number of transmissions with different temporal error correlations (error rate 0.10  $\lambda = 0$  for all receivers)



Figure 6. The maximum multicast delay with different temporal error correlations (error rate 0.10  $\lambda$  =0 for all receivers)



Figure 7. The expected number of transmissions with different spatial error correlations (error rate 0.10  $\tau$  =0.10 for all receivers)



Figure 8. The total RI with different error rates (  $\tau$  =0.10  $\lambda$  =0.20 for all receivers)



Figure 9. The maximum multicast delay with different error rates ( $\tau$  =0.10  $\lambda$  =0.20 for all receivers)



Figure 10. The total RI with different delay constraints (error rate 0.10  $\tau = 0.10 \lambda = 0.20$  for all receivers)



Figure 11. The total RI with a bad receiver 1 (error rate 0.01  $\tau$  =0.10  $\lambda$  =0.20 for other receivers, delay constraint 100ms)

# 8. References

- U.H. Reimers, "DVB-The Family of International standards for Digital Video Broadcasting," in Proceedings of IEEE, vol.94, no.1, pp. 173–182, January 2006.
- [2] E.N. Gilbert, "Capacity of a burst-noise channel," Bell Syst. Tech. J., vol.39, pp. 1253–1265, September 1960.
- [3] E.O. Elliott, "Estimates of error rate for codes on burstnoise channels," Bell Syst. Tech. J., vol.42, pp. 1977– 1997, September 1963.
- [4] M. Mushkin and I. Bar-David, "Capacity and coding for the Gilbert-Elliott channels," IEEE Trans. on Inform. Theory, vol. 35, pp. 1277–1290, November 1989.
- [5] J. Nonnenmacher and E.W. Biersack, "Optimal multicast feedback," in IEEE Infocom, (San Francisco, California), pp. 964, March/April 1998.
- [6] L. Rizzo, "Effective erasure codes for reliable computer communication protocols," ACM Computer Communication Review, April 1997.
- [7] D. Qiao and K.G. Shin, "A two-step adaptive error recovery scheme for video transmission over wireless networks," in Proceedings of IEEE Infocom, 2000, March 2000.
- [8] G.P. Tan and T. Herfet, "The Optimization of an RTP Level Hybrid Error Correction Scheme for DVB Systems in Wireless Home Networks under Strict Delay Constraint," IEEE Trans. on Broadcasting, November 2006.
- [9] G. Carle and E.W. Biersack, "Survey of error recovery techniques for IP-based audio-visual multicast applications," IEEE Network, vol.11, no 6, pp. 24–36, November –December 1997.

- [10] G. Tan and Th. Herfet, "Application Layer Hybrid Error Correction with Reed-Solomon Code for DVB Services over Wireless LANs," the 3rd IEEE International Conference on Wireless Communications, Networking and Mobile Computing (WiCOM), Shanghai, China, September 2007.
- [11] J. Kuri, S.K. Kasera. "Reliable multicast in multi-access wireless LANs," in Proceedings of Infocom '99.
- [12] J.Tourrilhes, "Robust broadcast: improving the reliability of broadcast transmission on CSMA/CA," in Proceedings of the Ninth IEEE International Symposium on. Personal, Indoor and Mobile Radio Communications, 1998.
- [13] S.K.S. Gupta, V.Shankar, and S. Lalwani, "Reliable Multicast MAC Protocol for Wireless LANs," IEEE International Conference on Communications, vol. 1, pp. 93–97, 2003.
- [14] J.R. Yee and E.J. Weldon, Jr., "Evaluation of the performance of error-correcting codes on a Gilbert channel", IEEE Tran. on Communications, vol. 43, no. 8, August 1995.
- [15] G. Bianchi, "Performance analysis of the IEEE 802.11 distributed coordination function," IEEE Journal on Selected Areas in Communications, vol. 18, no. 3, March 2000.
- [16] D.F. Towsley, J. Kurose, and S. Pingali, "A Comparison of Sender-Initiated and Receiver-Initiated Reliable Multicast Protocols," IEEE Journal on Selected Areas in Communications, vol. 15, no.3, April 1997.
- [17] S. Wiethölter and C. Hoene, An IEEE 802.11e EDCA and CFB Simulation Model for ns-2, http://www.tkn.tuberlin.de/research/802.11e\_ns2/.
- [18] J. Ott, S. Wenger, N. Sato, C. Burmeister, and J. Rey, "Extended RTP profile for RTCP-based feedback," draftietf-avt-rtcp-feedback-11.txt, August 2004.



# Routing and Wavelength Assignment in GMPLS-based 10 Gb/s Ethernet Long Haul Optical Networks with and without Linear Dispersion Constraints

# Le Nguyen BINH<sup>1,2</sup>

<sup>1</sup> Department of Electrical and Computer Systems Engineering, Monash University, Australia <sup>2</sup> Lehrstuhl für Nachrichen- und Übertragungstechnik, Technische Fakultaet der Christian Albretchs Universitaet zu Kiel, Germany E-mail: le.nguyen.binh@eng.monash.edu.au

# Abstract

Given a set of lightpath connection requests in an all-10 Gb/s optical dense wavelength division multiplexed (DWDM) Ethernet network, lightpaths are designed. In addition the wavelength channels are assigned subject to minimization of the channel blocking and provisional requests satisfying the limits due to accumulative linear dispersion effects over the hops.

This paper proposes a routing and wavelength assignment scheme for DWDM long-haul optical networks that includes routing, assignment and reservation of different wavelength channels operating under the Generalized Multiprotocol Label Switching (GMPLS) environment. The GMPLS framework can offer an approach to implement IP over DWDM with variable weighting assignments of routes based on the limitations due to residual dispersion accumulated on the lightwave path.

The modeling is implemented under the framework of an object-oriented modeling platform OMNeT++. Network performance tests are evaluated based mainly on a long-haul terrestrial fiber mesh network composed of as well as three topologies structured as chain, ring, and mesh configurations. Blocking probability of lightpath connection requests are examined with the average link utilization in the network employing variable number of wavelength channels in association with the limits of route distance due to linear chromatic and polarization mode dispersion effects.

Keywords: DWDM Optical Networks, Optical Transmission Systems, GMPLS, Routing and Wavelength Assignment (RWA), Wavelength Routers.

## 1. Introduction

In recent years, since the invention of optical amplifiers in late 1980s advances in dense wavelength division multiplexing (DWDM) technology has enabled the implementation of multi-Tera/bits/sec. capacity in intelligent optical networks (ION) [1]. In particular 10 Gb/s Ethernet optical networks over DWDM have been standardized and emerges as the most likely ultra-fast information networks in the near future. Undoubtedly systems and networks need to be interconnected in the optical wavelength domain or layer via optical routers and photonic switches and managed under appropriate protocols. The DWDM technology offers the capability of building very large wide area network (WAN) consisting of thousands of nodes with per-node throughputs of the order of tens of Gb/s.

However, at 10 Gb/s and possibly higher with the possibility of upgrading certain number of channels to 40 Gb/s or even 100 Gb/s, the detrimental effects of linear dispersion effects such as the chromatic dispersion (CD) and polarization mode dispersion (PMD) and are critical for error free transmission and interconnect between the IP routers. Furthermore, the number of wavelength channels are high with the average total power may be over the nonlinear threshold for the transmission lightpaths that would lead to the reduction of the eye opening of the receiver signals. We must note that these effects are not very critical for 2.5 Gb/s transmission bit rate which is now commonly used in

present optical networks.

In this paper we include only the residual linear dispersion and the PMD effects into the investigation of traffic engineering in IP over WDM optical networks, especially the routing and wavelength assignment (RWA) of wavelength channels over fiber paths of the networks. We understand that these transmission issues have not ever been taken into account in published research works. Due to the linear dispersion tolerance and nonlinear effects currently several advanced modulation formats have been studied [1] in order to combat these impairments. Usually these distortion effects are dependent on the length of the transmission that is commonly faced in all optical networks, especially in backbone WDM networks.

Routing and wavelength assignment (RWA), i.e. a cross-connecting of several wavelength allocated channels between DWDM optical networks, is very important [3]. One unique feature of DWDM networks is the tight coupling between routing and wavelength selection. Routing and wavelength assignment (RWA) problem is a major and complex problem associated with optical networks. Thus, given that (i) a set of dynamic and randomly chosen lightpaths that need to be established; (ii) constraints of the total number of wavelengths propagating in an optical fiber; and (iii) wavelength continuity constraint (WCC), which means a lightpath must use the same wavelength on all the links along its path from source to destination edge node, The controller for routing and assignment of wavelength channels must determine the routes over which these lightpaths should be set up and determine the wavelengths to be assigned to these lightpaths so that an optimum number of lightpaths may be established for the routing. Each wavelength channel is associated with a dispersion factor. Due to the dispersion slope of the transmission fiber, the dispersion compensation using dispersion compensating fiber is normally fully compensated at the center wavelength of the wavelength band. Thus there exists residual chromatic dispersion (CD) at other channels. At 10 Gb/s these residual dispersion effects are very critical. Furthermore the dispersion due to different traveling velocity of polarized modes of the single mode fiber, the polarization mode dispersion (PMD) is also important and must be taken into account.

Lightpaths that cannot be set up due to constraints on routes or wavelengths or the linear CD and PMD at which the eye opening penalty (EOP) at the receive end would suffer 3 dB, are said to be blocked, so the corresponding network optimization problem is to maximize the probability of set-up for a current connection request while minimizing the blocking probability of future connection requests. Furthermore in long haul terrestrial and or intercontinental optical networks, a generic management of the whole network status is preferred.

The concept of GMPLS offers an approach to implement IP over DWDM [4–10]. Indeed, GMPLS becomes evident that it is the best control plane solution for next-generation optical networking. In GMPLS, the MPLS label is generalized so that a label can also be encoded as a time slot, a wavelength, or a spatial identifier. The embrace by GMPLS is extremely large; the most critical part of it is to solve the general problem of dynamic lightpath establishment.

Therefore, in this paper various schemes for RWA algorithms are studied, compared and selectively implemented on a long haul optical terrestrial mesh network based on a common control plane GMPLS when residual chromatic framework. especially dispersion and polarization mode dispersion effects are taken into account. The nonlinear phase noises and impairments in network routing and wavelength assignment will be reported in a future article. Other network topologies also are investigated in this work such as a 4-node/8-node chain, a 5-node ring, and a terrestrial mesh network. We studied two kinds of wavelength assignment methods in all-optical networks: First-fit (Fixed-order) and random wavelength assignment, and two kinds of link metrics: simple total and available wavelength (TAW) and enhanced TAW. We use OMNeT++ [12] to simulate GMPLS-based alloptical networks and compare the blocking probabilities of these RWA methods on different topologies and traffic models. The simulation results show that first-fit wavelength assignment has a lower blocking probability than random wavelength assignment, and enhanced TAW performs better than simple TAW when link utilization becomes high. We also found that under some circumstances the difference in blocking probability may be widened and worsened due to additional limitations due the total accumulated dispersion along a routed lightpath.

The paper is organized as follows: Section 2 outlines the physical constraints on the routings of wavelength channels over long haul optical networks whose parameters are then declared. Also an ultra-fast automatic wavelength router is proposed. The specific characteristics of GMPLS are given. Section 3 and 4 propose the RWA algorithm and structures of simulation platform. The issues involved RWA problems in such linear dispersion (CD and PMD) effects are given. Our simulation model and simulation results are given in Section 5. Traffic performances of the proposed RWA under different network scenario under with and without dispersion effects are obtained and described. Finally, conclusion and further works are found in Section 6.

# 2. Network Physical Constraints

Routing and wavelength selection are coupled tightly in

WDM all-optical networks. When establishing a connection in optical networks, both routing (path selection) and wavelength assignment (allocating a wavelength along the selected path) must be considered. This is referred to as routing and wavelength assignment (RWA) problem. This problem is more complex than routing problem in electronic networks. Optical networks operating at lower bit rates less than 2.5 Gb/s can be routed under constraints of [3]: (i) Wavelength continuity constraint (WCC): a lightpath must use the same wavelength on all the links along its path from source to destination edge node; (ii) Distinct wavelength constraint: all lightpaths using the same link (fiber) must be allocated distinct wavelengths; (iii) Limited number of wavelength per fiber link is restricted within the Cband and by the nonlinear threshold, mainly due to the self phase modulation (SPM) effects. Furthermore for higher bit rates 10 Gb/s and above additional constraints are (iv) Total fiber routing path is limited due to the total residual dispersion allowable on the total routing distance that suffer an eye opening penalty of 3 dB; and (v) Total accumulated noise and hence EOP due to amplified stimulated amplification within optical amplifiers placed in each fiber spans that is important for backbone long haul networks. These noise contributions are ignored in this work.

It is noted here that the dispersion compensation can normally be made for a specific wavelength of the Cband and there are always residual dispersion mismatched at other wavelength channels due to dispersion slopes. Thus the effects of residual dispersion are very critical in the design of the DWDM network routing and wavelength assignment as these dispersion effects are length dependent.

# 2.1. Wavelength Continuity Constraint (WCC)

The WCC is a unique constraint of WDM all-optical networks. In WDM networks, if a connection is accepted, it should be assigned a path through the network and a wavelength. This wavelength must be the same wavelength on all links along the path. Consequently, it is likely that any algorithm engaging WCC would suffer higher blocking probability. If another lightpath has to be set up between node 1 and node 3, it will be blocked, because there is no available wavelength along its lightpath (i.e. no common available wavelength on both links). Due to high costs of all-optical wavelengthconversion photonic components, it is assumed that none of the optical cross connectors (OXCs) has wavelength conversion capability.

The algorithm must establish dynamically an end-toend path between any ingress nodes and any egress nodes in the all-optical network. Our main goal is to minimize the blocking of lightpaths in the network. Different metrics for dynamic routing algorithms are examined. For those parameters which are not wavelength-based then constraints are not taken into account, unlike in other published works physical constraints are integrated them in the routing decision process [1].

## 2.2. Distinct Wavelength Constraint

Different lightpaths on the same physical link (fiber) must use different wavelengths. Each fiber in WDM networks may contain several WDM wavelengths subject to power below specific power threshold due to nonlinear effects. However these nonlinear effects are ignored in this report. It is assumed that the operating wavelengths are placed in a 50, 100 or 200 GHz grid (0.4nm wavelength spacing) with a nominal center frequency of 193.1THz (1552.52 nm) in the middle of the 1.55 $\mu$ m fiber and EDFA pass-band [4]. More details are specified in ITU G.692.

#### 2.3. Number of Wavelengths per Fiber Link

Consider one link between two adjacent nodes. Each link cable may be composed of several fibers in which several multiplexed wavelength channels are transmitted. Figure 1 shows a typical automatic wavelength optical cross connect (OXC) including wavelength converters, wavelength switching matrix whose control signals come from the central management systems complying the state of the designed algorithm and global information obtained from the nodes of the networks. The wavelength converter is a semiconductor optical amplifier biased to operate for wavelength conversion. This conversion is associated with the wavelength multiplexer so that once the conversion is completed to a wavelength output port the channel is routed immediately to the output. Thus this structure allows an ultra-fast routing. The limit of the operation of this OXC depends on the conversion speed of the SOA which is, under present technology, in order of picoseconds for GHz operation.

Let  $\lambda_{i, j}$  be the wavelength number i in the fiber j. Assuming the operating wavelength range is the C-band. The wavelengths range in the C-band is 1,530nm to 1,565 nm. Each wavelength is assigned a specific emission wavelength that is specified by the ITU-T. Consequently, it is assumed that there is a maximum of 44 wavelengths per fiber with 100 GHz spacing between channels. The wavelengths are ordered as follows:  $\lambda_{1,j} =$ 1,530.33nm and  $\lambda_{44, j} =$  1,564.68nm, with j being the number of fibers used. In a nutshell, the physical-based wavelength constraints are:

$$\begin{cases} \lambda_{i,j} = wavelength \ i \ in \ j^{th} fiber....i \in 1 \div 44 \\ \lambda_{1,j} = 1,530.33nm; \ \lambda_{44,j} = 1,564.68nm \ (C-band) \end{cases}$$
(1)

wavelength channels propagating in the same fiber path.



Figure 1. Schematic of an optical cross connect (OXC) incorporating wavelength conversion and routing by an SOA and demultiplexer and switching matrix for routing a specific wavelength to a designated output port for further routing.

Still, it can be shown that a multi-fibre RWA problem is algorithmically equivalent to a single fibre based optical network. Thus, in the following, an ordered list of wavelengths with only one fibre is considered. Further it is very unlikely that the 44 wavelengths transmitted or available in the C band will be all practically employed. Usually only a few wavelength channels are lit. Thus it is reasonable to assume that the population of the WDM slots is limited to a multiple of 8 wavelengths, up to 32 or 40 with 100 GHz spacing between wavelength channels. As the number of fibers and the optical power launched in each wavelength, there are possibilities that the total average power propagating through a fiber path would reach above the nonlinear threshold (normally 3 dBm) and hence nonlinear phase distortion. In this work we set the limit at 10 dBm taking into account the randomness of the chance of simultaneous transmission of all channels.

## 2.4. Residual Linear Dispersion Constraints

For high capacity long-haul transmission employing external modulation, the dispersion limit can be estimated in the following equation [38,39].

$$L_{D} = \frac{10^{5}}{D.B_{R}^{2}}$$
(2)

where D is the dispersion factor and  $B_R$  is the bit rate,  $L_D$ 

is the dispersion length.

#### 2.5. Polarization Mode Dispersion Limit

The length limit by the PMD effect can be defined as follows. Constraint on the dispersion length due to PMD must be included and an estimate of the transmission length limit due to PMD effect is given as [39]

$$L_{PMD} = \frac{0.02}{\left\langle \Delta \tau \right\rangle^2 \cdot B_R^2} \tag{3}$$

with  $L_{PMD}$  is the maximum length due to the PMD effect that generates a 1dB penalty on the eye opening.

## 2.6. GMPLS Framework

Generalized Multiprotocol Label Switching (GMPLS) extends the label switching architecture proposed in MPLS to other types of non-packet based networks, such as SONET/SDH based networks and WDM networks. It has recently become evident that GMPLS is the unified control plane solution for ultra-high capacity, ultra-fast optical networking. In GMPLS, the MPLS label is generalized so that a label can also be encoded as a time slot, a wavelength, or a spatial identifier. The basis of the GMPLS framework is defined in [5,9], including two main parts of routing and signaling.

GMPLS routing proposes extensions to interior gateway protocols (IGP) of routing being the open shortest path first extended for traffic engineering (OSPF-TE) [7] and intersystem to intersystem extended for traffic engineering (ISIS-TE) [8]. The protocols provide the distribution of link state information with traffic engineering information and attributes such as network topology, resource availability, and administrative constraints.

GMPLS signaling protocols include reservation protocol extended for traffic engineering (RSVP-TE) [10] and the constrained label distribution protocol (CR-LDP) of the MPLS-TE framework [11]. The GMPLS signaling protocols extend certain base functions of the RSVP-TE and CR-LDP signaling, which originally configure and control the distributed label switched paths (LSP). The extensions impact basically on LSP properties on how labels are requested and communicated, unidirectional nature of LSPs, how errors are propagated, and information provided for synchronizing the ingress and egress nodes.

## 3. Routing and Wavelength Assignment

RWA algorithm must be able to identify the maximum of explicit routes in IONs in order to satisfy dynamic routing requests. The RWA algorithm separates routing, wavelength assignment (WA) and reservation

mechanisms. The search of an explicit route in the network that would meet the physical constraints of the ION is the principal aspect. Such a separation between routing and reservation protocols eases further works.

RWA algorithm is designed so that it gives an output an explicit route in the ION that determines each node to be traversed by the lightpath. Besides, a WA algorithm allows the determination of a wavelength channel satisfying the WCC that is reserved on the lightpath. Explicit routes and their associate wavelengths are then treated as the inputs of a reservation protocol in order to establish the optical circuit in the network. Present work develops a RWA scheme model of a GMPLS-based ION based on an adaptive link-state routing using global network knowledge including the rise time and dispersion budget of proposed routing link/distance. This choice is due to the followings. Table 1 shows a summary of the WA algorithms on a RWA problem.

Table 1. WA (wavelength assignment) algorithms in the RWA problem

WA heuristics	VA euristics Characteristics		Drawback
First-fit	The wavelengths are indexed, and a lightpath will attempt to select the wavelength with the lowest index before attempting to select a wavelength with a higher index	If global information, outperforms random heuristics.	Possible blocking if simultaneous lightpath connections.
Random	Select one of the wavelengths at random	If local information, outperforms first-fit heuristics. Very simple.	Does not provide optimal wavelength assignment.
Least-used	The wavelength which is the most used in the rest of the network is selected.	Spreads the load evenly across all wavelengths	Global information needed.
Most-used	The wavelength which is the most used in the rest of the network is selected.	Provides maximum wavelength reuse in the network	Global information needed.

# 3.1. Routing

Several types of RWA algorithms have been studied. Three main types of routing schemes are fixed, fixedalternate and dynamic routing. Fixed and fixed-alternate routing principles are very simple. However, they have poor capabilities when dealing with self-healing capabilities and can yield a very high blocking probability under dynamic network operations. Thus, dynamic routing is considered.

The main disadvantage of dynamic routing is possibility for high overhead due to route advertisements. Nevertheless, the choice of an algorithm is always a trade-off between network performance and traffic overhead. Considering earlier research studies presented in [3,12,13], it has been thought that a dynamic routing is one of the best ways to solve the RWA problem.

In the first instance, both adaptive routing schemes, that is link-state routing and distance-vector routing are considered. Contrary to link-state based routing algorithms that flood packets onto the network distancevector algorithms send their route advertisements only to their neighbors. In link-state algorithms, the link state update is flooded onto the whole network or area of OSPF.

#### 3.1.1. Distance-vector Routing Algorithm

A distance-vector routing approach needs to use a distributed algorithm such as the popular Bellman-Ford algorithm. It is possible that a distance-vector algorithm can minimize the traffic overhead in the network while still giving relatively low blocking probabilities under high loads. A distance-vector algorithm is also interesting when considering constraint routing. Indeed, it is a property of the Bellman-Ford algorithm that, at its *h*-th iteration, it identifies the optimal (in our context: maximal number of wavelengths) path between the source and each destination, among paths of at most h hops. This is recommended for QoS routing implementations [14]. Furthermore the distance vector is also tagged with the dispersion and distortion estimated for the routing distance.

Bellman-Ford algorithm progresses by increasing hop count, it essentially provides for free the hop count of a path as a second optimization criterion. This property is very interesting when applied to all-optical networks. However, even if the shortest path is sacrificed for a longer – in hops – route, the lightpath should not be too long as this could lead to a poor optical power and dispersion budgets.

Routing Information Protocol (RIP) is the most popular implementation of a distance-vector based protocol because it is simple and it is well suited to small networks. However, RIP has several flaws that make it unsuitable for ION. Particularly, RIP is unsuitable for large configurations and the convergence of the algorithm can also be lengthy; it suffers also from the count-to-infinity problem, of which the best remedy is to implement link-state algorithms.

#### 3.1.2. Link-state Based Routing Algorithm

In link-state based routing, information is only sent when changes occur. A node builds up first a description of the topology of the network. Then it may use any routing algorithm to determine the route.

In the context of the ION, in order to compute an explicit route, it is also much easier to use a link-state based routing algorithm. Indeed, the lightpath to be established is more optimal if each node has a global knowledge of the network. It is also a property of the

Dijkstra algorithm that a complete route from a source to any other node in the network can be easily found by recursive iteration on the graph.

Link-state based algorithms can use any routing algorithm. Different RWA algorithms can be implemented based on different optimization schemes (partial or total wavelength knowledge).

The open shortest path first (OSPF) protocol, known mostly for its second version [23], is the most widely known link-state based routing protocol and is employed for the simulation platform of this work. It has been increasingly popular over RIP, because it is most suitable for large networks. OSPF is an open source algorithm that can be found in different languages, including C++. This is a particular advantage for simulation program based on OMNeT++, a C++ based simulator.

#### 3.1.3. Route Advertisements

The route advertisement messages are built according to the OSPF extensions for GMPLS. This work proposes to extend the sub type length value (TLV) field relative to the interface switching capacity descriptor (ISCD) field of the GMPLS OSPF extensions.

In order to dynamically monitor the state of the network, each GMPLS router keeps track of all the wavelength capabilities of the whole network in a links database. This database is constantly updated each time a "routing" message is received. The "routing" message contains the description of the wavelength capabilities of a certain link that have changed very recently.

Those "routing" messages are actually flooded on the whole network by the node which one of its links' capabilities changed containing the address of the extremities of the link of which the wavelength capability has changed and the state of the changed wavelength.

The links database is composed of records, where each record describes one link of the network. Each record contains the node addresses of the link's extremities, a wavelength capability field and a metric field. The wavelength capability field describes explicitly the wavelength resources of the considered link. The metric field is the cost to use this link when performing the shortest path calculation. In this work, two different metrics have been implemented. They are both function of the total number of wavelengths and the number of available wavelength(s) on the link as described in the next part.

### 3.1.4. Link Metrics

The link metrics represent the costs to use certain links in the network. Intuitively, the link cost is a linear function of the total number of wavelength and the number of available wavelengths and the residual dispersion if taken into account.

The routing scheme is tested under different link metrics: simple total and available wavelength (TAW) or enhanced TAW. Let  $\lambda_{i,j}^a$  be the number of available (unused) wavelengths on the link (i,j) and the total number of possible wavelengths on that link. The simple TAW metric represents the load assigned to a link and is defined by:

$$w_{i,j} = 1 - \frac{\lambda_{i,j}^a}{\lambda_{i,j}^T}, \forall (i, j) \in E$$
(4)

The enhanced TAW metric is to test the metric proposed by [24] that equivalently minimizes the probability of blocking on an explicit route in which the weights are assigned following a log scale as:

$$w_{i,j} = -\log\left[1 - \left(1 - \frac{\lambda_{i,j}^a}{\lambda_{i,j}^T}\right)^{\lambda_{i,j}^a}\right], \forall (i,j) \in E$$
(5)

If residual dispersion is accounted then a constant or a linear or quadratic function of dispersion factor can be added into (4) or (5).

#### 3.1.5. Path Calculation

The *links database* that has been freshly updated by the GMPLS router serves as the basis of the path discovery calculation based on the Dijkstra algorithm. This path computation is performed by each node in the network when it receives a "*routing*" message, which is an update of a certain link capability in the network.

The routing algorithm actually allows each node to build its own photonic database that contains N records, where N is the number of nodes of the network. Each record is based on the following structure: (i) a node destination address — this identifies the record; it is the node to reach; (ii) a total cost to this destination node — this is the total cost when taking the shortest path route to the destination node; (iii) an address of the last-but-one node on the shortest path route to the destination node; (iv) an end-to-end available wavelength capability, which determinates explicitly the possible wavelengths to be assigned, if any, on the shortest path route; and (v) an explicit route holding an ordered list of all the addresses of the nodes on the shortest path route.

Given the links database, the first three fields are actually the direct output of the Dijkstra shortest path algorithm [14]. The last two fields are the result of a sub-routine of the Dijkstra algorithm. Indeed, it is a property of the Dijkstra algorithm that it is possible to find the list of the nodes of each shortest path by a simple recursive call.

We have also set constraints on the dispersion and attenuation of the fiber paths, in particular the polarization mode dispersion factor (PMD) as PMD is unlike the chromatic dispersion dependent on the fiber length. For chromatic dispersion we assume that the (standard single mode fiber) SSMF is compensated with dispersion compensating fiber and only residual dispersion is used. At 10 Gb/s these dispersion effects are very critical. Additional nonlinear phase and hence distortion effects are not considered in this work.

## 3.1.6. Lightpath Establishment

There are two categories of lightpath establishment, one is the static lightpath establishment (SLE), and the other is dynamic lightpath establishment (DLE).

#### 3.1.7. Static Lightpath Establishment

Static lightpath establishment (SLE) is used in the design and capacity planning phase of architecting an optical network. It can be logically decomposed into four subproblems. Assuming no wavelength conversion, the subproblems are listed as follows [33]: (i) Topology subproblem: determine the logical topology to be imposed on the physical topology, that is, determine the lightpaths in terms of their source and destination edge nodes; (ii) Lightpath routing sub-problem: determine the physical links which each lightpath consists of, that is, route the lightpaths over the physical topology; (iii) Wavelength assignment sub-problem: Determine the wavelength each lightpath uses, that is, assign a wavelength to each lightpath in the logical topology so that wavelength restrictions are obeyed for each physical link; and (iv) Traffic routing sub-problem: Route packet traffic between source and destination edge nodes over the logical topology obtained.

The SLE is also referred as static RWA problem, can be formulated as an integer linear program (ILP) where the objective is to maximize the number of connections that are successfully routed. A number of studies have investigated the static RWA problem.[12,15,17,29]. The ILP formulations are NP-complete and therefore may only be solved for very small systems. For large systems, heuristic methods must be used. A number of heuristics were proposed for the problem, they can be divided in four classes [17]: (i) heuristic solutions of the mixed integer linear programming problem; (ii) maximizations single-hop traffic flows; (iii) heuristic of the maximizations of the single-hop and multi-hop traffic flows; (iv) algorithms based on the adoption of a preestablished regular logical topology and on the optimization of the nodes placement according to the traffic pattern.

#### 3.1.8. Routing Sub-problems

Fixed Routing — The most straightforward approach to routing a connection is to always choose the same fixed route for a given source-destination pair. This is called fixed routing. Fixed routing approach is easy to implement. However, it is very limited in terms of routing options. If resources (wavelengths) along the path are tied up, it can potentially lead to high blocking probabilities in the dynamic case, or may result in a large number of wavelengths being used in the static case. Moreover, fixed routing may be unable to handle fault situations in which one or more links in the network fail.

Fixed-Alternate Routing — An approach to routing that considers multiple routes is fixed-alternate routing. It increases the likelihood of establishing a connection by taking into account network state information. In fixed-alternate routing, each node in the network is required to maintain a routing table that contains an ordered list of a number of fixed routes to each destination node. Fixed-alternate routing provides simplicity of control for setting up and tearing down lightpaths and it may also be used to provide some degree of fault tolerance upon link failures. It can significantly reduce the connection blocking probability compared to fixed routing.

Adaptive Routing — In adaptive routing, route form a source node to a destination node is chosen dynamically, depending on the network state. The network state is determined by the status of all connections that are currently in progress. In order to choose an optimal route, a cost is assigned to each link in the network based on current network state information, such as wavelength availability on links, the total residual dispersion amount. A least-cost algorithm is then executed to find the minimum-cost route. Whenever a connection is established or taken down, the network state information is then updated.

## 3.2. Wavelength Assignment

In general, if there are multiple feasible wavelengths between a source node and a destination node, then a wavelength assignment algorithm is required to select a wavelength for a given lightpath. The wavelength selection may be performed either after a route has been determined, or in parallel with finding a route.

Since the same wavelength must be used on all links in a lightpath (noted that wavelength conversion is out of scope of this paper), it is important that wavelengths are chosen in a way which attempts to reduce blocking for subsequent connections. For the case that there are multiple feasible wavelengths between a source node and a destination node, heuristic methods must be used to assign wavelengths to the lightpath. A review of wavelength-assignment approaches can be found in [13] as follows: (i) Random: this scheme first searches the space of wavelengths to determine the set of all wavelengths that are available on the required route, randomly choose one wavelength (usually with uniform probability). (ii) First-Fit or fixed-order, in which all wavelengths are numbered. After searching for the available wavelengths, a wavelength with lowest index is chosen. The idea behind this scheme is to pack all of the in-use wavelengths toward the lower end of the wavelength space so that more wavelengths with higher indexes can be used for longer-hop connections. (iii) Least-Used selects the available wavelength that is the least used in the network, so the load is balanced on all the wavelengths. (iv) Most-Used is the opposite scheme of Least-Used, it always choose the available wavelength that is most used in the network. It packs connections into fewer wavelengths. (v) Min-Product is used in multi-fiber networks. In single-fiber networks, it is the same as First-Fit. It packs wavelengths into fibers, minimizes the number of fibers in the networks. (vi) Least-Loaded is also used in multi-fiber networks. It chooses the wavelength that has the largest residual capacity on the most-loaded link along route p. (vii) Max-Sum is proposed for multi-fiber networks but can also be used in single-fiber case. Generally speaking, the scheme considers all possible paths (lightpaths with their pre-selected routers) in the network and attempts to maximize the remaining path capacities after lightpath establishment. (viii) Relative Capacity Loss (RCL) is based on Max-Sum. RCL calculates a Relative Capacity Loss for each path on each available wavelength and then chooses the wavelength that minimizes the sum of the relative capacity loss on all the paths [36]. (ix) Wavelength Reservation reserves some wavelengths for multi-hop connections. This scheme reduces the blocking for multi-hop traffic but sometimes may overpunish single-hop traffic; and (x) Protecting Threshold, in which a single-hop connection is assigned a wavelength only if the number of idle wavelengths on the link is at or above a give threshold.

In this work the first-fit and random schemes are also selected for modeling the lightpath setup and RWA schemes in GMPLS-based IP-over-DWDM optical network. First-fit chooses the first wavelength available in an ordered list of wavelength in contrast to the random scheme which chooses a wavelength randomly between the different available wavelengths. Those two schemes are used directly at the output of the path calculation algorithm in order to determine the wavelength to be reserved on the path to each destination in the network from a source node.

## 3.3. Signaling and Resource Reservation

In order to set up a lightpath, a signaling protocol is required to exchange control information among nodes and to reserve resources along the path. In most cases, the signaling protocol is closely integrated with the RWA algorithms.

Signaling and resource reservation protocols may be categorized based on whether the resources are reserved on each link in parallel, reserved on a hop-by-hop basis along the forward path, or reverse path (backward). Algorithms will also differ depending on whether global information including the dispersion factors of each hop is available or not.

# 4. Model Implementation

# 4.1. OMNeT++ Simulation Model

OMNeT++ is a free discrete event network simulation tool, which was primarily designed for the simulation of communication networks, multi-processors and other distributed systems [5]. The OMNeT++ based model is built to test two link metrics mentioned previously: the simple TAW and the enhanced TAW; and two wavelength assignment heuristics: first fit and random assignments. Different topologies are included for simulation such as single link, chain, ring, and mesh topology.

The model is composed of simple modules referred to as GMPLS router and compound modules referred to as EndSystem. The GMPLS router module contains all the RWA and reservation algorithms, while the EndSystem module is responsible of generating and analyzing the response to connection request messages.

The EndSystem module composes of Generator and Sink simple modules.

The Generator module generates a certain number of connection requests to randomly chosen destination nodes in the ION.

The Sink module receives the responses to each connection request generated by its counterpart Generator module. The principal task of the Sink module is to analyze the responses and to generate a statistical table.

Note that the Endsystem compound module actually models an end system in an ION. Effectively, the EndSystem module represents one or several end systems such as switches or routers, e.g. such end systems may be an ATM switch or IP routers using MPLS-based protocols. It is also essential to standardize the physical interface between the clients (end systems) and the transport network (ION) which can be designated as the UNI interface. An implementation agreement for the UNI has been proposed by the OIF [30].

Each GMPLS router module is associated to only one Endsystem module. Those modules have the same address. The network model is based on the USA Abilene network that composes of 12 GMPLS router modules scattered at different GigaPoPs; each GMPLS router module is associated to its own Endsystem module. An object-oriented framework has been developed with an UML class diagram (implemented with Rational Rose) of the RWA model which describes the relation and contents of different of classes of the simulation framework which are Generator, Sink, GmplsRouter, Node, LinkInfo, NodeInfo, ExplicitRoute and LambdaCap. The details of these classes can b explained on request of readers.

#### 4.2. Reservation Protocol

For simplicity, a reservation protocol based on parallel reservation has been implemented in order to test the different routing and wavelength assignment schemes implemented. The reservation scheme and process of messages exchanges are shown in Figure 2 and described as follows:

The Generator sends a "request" message for a new connection to another randomly chosen node. When GMPLS Router receives a "request" message, it calculates an explicit route to the destination requested and assigns a wavelength for this connection and then estimation of the total residual and PMD dispersion, thence broadening time.

If both route and wavelength are available, the GMPLS router reserves in parallel the route and the wavelength. It sends in parallel a "reserve" message to each node of the explicit route, excepted itself. When a node receives the "reserve" message, it checks if the requested wavelength is available on the link to its predecessor in the explicit route. Then it sends a "response" message back to inform the source whether the reservation is successful.



Figure 2. Parallel reservation model



Figure 3. Illustration of the use of TAKEDOWN messages

If the requested wavelength is not available in one or several node on the explicit route, the source will send a TAKEDOWN message to all the nodes, which have already reserved the wavelength on the explicit route to reset it as available (see Figure 3).

#### 4.3. System Parameters

The network traffic is generated in terms of connection request, which is a request to establish a lightpath from a source node to a randomly chosen destination node. The connection request arriving at each router is assumed to follow an exponential distribution with mean  $\lambda_{poisson}$  per unit of time. The system variable parameter is  $\lambda_{i,j}^{T}$ , the total number of wavelengths on each link and the connection requests arrivals.  $\lambda_{i,j}^{T}$  changes between 8 and 24 by an increment of 8. Performance parameters considered are blocking probability *P* and link utilization *U*.

The blocking probability P is the probability that a connection request is blocked due to no available wavelength for the selected path. The average link utilization U is defined as the percentage of time that all wavelengths of each link in the network are fully utilized, as

$$U = \frac{\sum_{i,j}^{i < j, i \neq j} \left( 1 - \frac{\lambda_{i,j}^{a}}{\lambda_{i,j}^{T}} \right)}{\sum_{i,j}^{i < j, i \neq j} \lambda_{i,j}^{T}}$$
(6)

Under the residual dispersion, the constraint for routing is

$$\sum_{k}^{N_{l}} L_{k,l} \leq L_{D}, L_{PMD}$$
(7)

with the subscript k indicating fiber path and  $N_l$  is the total number of the fiber paths that a routing path l is selected. The length limit is set by the dispersion length  $L_D$  and the PMD effect,  $L_{PMD}$ , the maximum length due to the PMD effect that enforces a 1dB penalty on the eye opening. It is noted that in calculating the total length of the routing route both the linear chromatic and polarization mode dispersion impairments are summed up as they would be superimposed on each other. PMD is a random process and it is modeled with the total average group velocity delay and a Maxwellian distribution. Since the lightpath travels the whole distance from one node to the other, it is reasonable that the average and standard deviation is sufficient for our estimation of the distance and distortion effects.

In addition to the CD and PMD, nonlinear phase distortion is also critical that would be investigated in a future article. The nonlinear effects [38] include the self phase modulation (SPM), the cross phase modulation (XPM), the four waves mixing (FWM), stimulated

## ROUTING AND WAVELENGTH ASSIGNMENT IN GMPLS-BASED 10 GB/S ETHERNET LONG 163 HAUL OPTICAL NETWORKS WITH AND WITHOUT LINEAR DISPERSION CONSTRAINTS

Raman scattering (SRS) and stimulated Brillouin scattering (SBS). Of these effects SPM and XPM are generated with the change of the phase of the carrier due to the total intensity of the light paths exerted over the core of the fiber and length dependence. The XPM would create a low frequency noise component in the frequency spectra of other lightpaths from lightwave signals of a lightpath. When the utilization is increased the nonlinear effects SPM and XPM are enhanced and thus distortion would affect the quality of the transmitted signals in the lightpaths. Similarly for SRS but the distortion is generated to the lightpaths whose spectra are about 100 nm away from the signal bands. This would be significant when both C- and L- bands are employed. These nonlinear effects are not included in this work but will be reported in a future article.

## 5. Simulation Results and Discussions

One of our goals is to estimate the performances of different RWA algorithms by simulations under without and with residual dispersion effects. The blocking probability P versus the average link utilization U is used as the performance criterion. Under the condition that the wavelength conversion is not used, the wavelength continuity constraint would lead to a high level of blocking when the load increases when no dispersion constraints are imposed. We then investigate the performance under the dispersion constraints. Firstly, wavelength assignment schemes are compared. The blocking in the network as a function of the average link utilization is described, analyzed and discussed. Then, various RWA algorithms are tested based on some specific network topologies.

# 5.1. Performance under No Constraints of Dispersion Effects

The first-fit and random wavelength assignment schemes are simulated and presented in this sub-section under non constraints of linear dispersion and nonlinear effects. The blocking probability *P* is plotted as a function of the link utilization *U* for different values of  $\lambda_{i,j}^T$ , the total number of wavelengths in a fiber.  $\lambda_{i,j}^T$  is varied between 8 and 24 by increment of 8 for different schemes, including simple TAW metric, enhanced TAW metric, first-fit and random wavelength assignment are tested. Figure 4 shows the difference of performance between the two wavelength assignments. They represent blocking probability *P* versus the normalized link utilization for: Figure4(a) and (b)  $\lambda_{i,j}^T = 8/\text{simple TAW}$ ,  $\lambda_{i,j}^T = 8/\text{enhanced TAW}$ , (c) and (d)  $\lambda_{i,j}^T = 16/\text{simple}$ TAW,  $\lambda_{i,j}^T = 16/\text{enhanced TAW}$ , and (e) and (f)  $\lambda_{i,j}^T = 24/\text{simple TAW}$  and  $\lambda_{i,j}^T = 24/\text{enhanced TAW}$ , as indicated. It is found that the first-fit wavelength assignment performs superior, with much less blocking problem, than the random scheme when the average link utilization is between 40% and 60%, and when the link utilization becomes relatively high, the difference between these two wavelength assignments is very small. The results are independent with either the TAW type or the number of total wavelengths.



Figure 4. Blocking probability versus link utilization with total number of wavelength for First-fit and random wavelength assignments — no linear dispersion and nonlinear effects.

Shown in Figure 5 is the difference of performance between two metrics used simple and enhanced TAW. When the link utilization is relatively low, the difference is not much, but with link utilization becomes high, enhanced TAW performs better. In general, it has been observed that the enhanced TAW metric performs better than the simple TAW metric at higher utilization. One reason is that the enhanced metric tries to minimize the weight of the explicit route while still trying to minimize the probability of blocking for further requests. Conversely, the simple TAW metric does not take into account any probabilities. It just tries to minimize the total link utilization on a possible explicit route from the



ingress node to the egress node.

Figure 5. Blocking probability versus link utilization with total number of wavelength for simple and enhanced TAW.

## 5.2. Performance under Constraints of Dispersion Effects for Specific Networks

The purpose of this section is to test the affection of network topologies on different RWA algorithms under the constraints of linear dispersion and nonlinear effects. Specific networks of 5 node ring and a specific long haul backbone mesh network (see Figure 6) are chosen so that the lengths of the fiber physical paths can be determined.

Now let

$$A = \{A_1 \div A_4\} = \{4 \text{-node chain, } 8 \text{-node chain, } 5 \text{-node ring, Specific Net}\}$$
(8)

be the set of four topologies (Networks with assigned codes of A2, A3 and A4) to be investigated including a 4-nodes (A1) and 8-nodes chains (A2), a 5-nodes ring (A3), and the Specific Net (Figure 6 as A4), respectively. A total wavelengths number of 24 which is chosen with 50 GHz channel spacing for the study of the traffic performance o the networks under no dispersion effects and with linear CD and PMD effects. Note that the dispersion due to nonlinear self phase effect is not considered as a constraint [40]. The network routing is subjected to the linear dispersion constraints given in eq.(6) and (3) and (4). For the same reason as mentioned above, we compare two wavelength assignment heuristics — the first-fit and random wavelength assignments.

The simulation results are shown in Fiugre 8 and we can make the following observations. Firstly, under the case of the simplest topology - single link, suppose there are n total wavelengths on the link, then when the number of requests is smaller than n, the block will never happen. This is because there must be unused wavelength(s) on the link, and no matter which wavelength assignment heuristic is adopted, the source chooses the wavelength from those unused. Secondly, when the link utilization is low, there is no difference between First-fit and random wavelength assignments, since the connection requests can always be satisfied. When link utilization becomes higher, the block happens. In the same topology, the First-fit WA has a lower probability than random blocking wavelength assignment. But when link utilization is very high, the difference between these two wavelength assignments becomes narrower.

It also can be seen that the difference between first-fit and random wavelength assignment heuristics of the 4nodes chain is quite small, but that of 8-nodes chain becomes larger. Naturally the reason is destinations of connections are chosen randomly, so with 8-nodes chain there should be more multi-hop connection requests than that with 4-nodes chain. The first-fit heuristic tends to pack all of the in-use wavelengths toward the lower end of the wavelength space so that more wavelengths with higher indexes can be used for longer-hop connections. Thus on 8-nodes chain the advantage of First-fit heuristic becomes more obvious.



Figure 6. Specific Net with Cities and Distance (in number of spans). A typical long haul back bone optically amplified

dispersion managed network — the number indicates the number of spans in multiples of 100 km spans with dispersion compensating module and optical amplifiers. Capital letters indicate the sites for IP routers and OXCs. Residual dispersion is taken as 2% of the link distance.

The ring and 4-node chain topologies give much lower blocking probability as compared with 8-node chain and mesh topologies (The Specific Net of Figure 6).

We now demonstrate the two link metrics: simple TAW and enhanced TAW. Link metrics are used in path selection. In the chain topologies, there is only one path between the source and the destination, so there is no difference between simple TAW and enhanced TAW. The blocking probability for a 5-nodes ring (A3) and Specific Net (A4) topologies with a total of 24 wavelength channels. A3 network is set under the condition of equal number of spans per link of 5 between nodes with 100 km completely dispersion compensated per span - that is a total of 500 km of dispersion compensated fibers and optical amplifiers. A residual dispersion is set at 2% of the dispersion of the uncompensated dispersion of the lowest dispersion value (e.g. 2% of 17 ps/nm/km of standard single mode optical fiber). An average PMD first order value of 0.1  $ps/(km)^{1/2}$  is also included. As can be seen from Figure 7, the enhanced TAW performs better than simple TAW, because the enhanced TAW takes into consideration the further requests and tries to minimize the probability of blocking for future, while simple TAW only choose the link with a light load. The effects of CD and PMD have shown clearly an increase in the blocking probability when the utilization of the link paths is increased. This is expected due to the limitation of the dispersion budget given in (2) and (3). We note here that no statistical property of the PMD is taken into account and that we may expect some sudden fluctuation of the blocking probability in practice.



Figure 7. Blocking probability versus link utilization with

total number of wavelength is 24 for 2 kinds of topology under with and without constraints of linear chromatic dispersion (CD bold circle dotted line) and CD + PMD (random) effects (square dotted line). Legend: dotted line indicates "random".



Figure 8. Simple TAW vs. Enhanced TAW on node-ring and sample network (Figure 6) (a) without (b) with dispersion constraints using random and first-fit strategy.

Also, numerous simulations have been conducted to test these RWA heuristics on different traffic schemes and total available wavelengths. The traffic schemes include uniform, exponential and Poisson distribution, and the number of total available wavelengths is various among 8, 16 and 24. There is no obvious impact on the blocking probability due to different traffic schemes and the number of total wavelengths. The first fit WA performs better than random WA, and enhanced TAW performs better than simple TAW. Figure 8(a) and (b) shows a contrast of the blocking probability of the simple TAW and enhanced TAW under without and with dispersion constraints with different algorithms.

Concerning with TAKEDOWN message of reservation scheme, in our previous simulations, the

TAKEDOWN message never appears. This is because when a wavelength is reserved on some link(s), the UPDATE messages will be sent to every node in the network. When a source calculates the explicit routes, it requires the knowledge of whether a wavelength is unused or not, so it would not seek other nodes to reserve those unavailable ones. There is one case that the TAKEDOWN message would be used: after a wavelength is reserved on some link, and before the UPDATE message reaches a source, this source receives a connection request and estimates some explicit route using that wavelength on that link. Then the source seeks for the downstream node of that link to reserve the wavelength, which is already reserved for some other connection. Clearly, the node responds negatively, and TAKEDOWN messages are to be sent from the source to other nodes on that explicit route.

# 6. Concluding Remarks

In this paper we have presented a simulation model on GMPLS-based all-optical networks using OMNeT++ to study the traffic performance of RWA of channels within the C-band of optical fibers under the influences of the dispersion effects due to chromatic and polarization mode dispersion. The blocking probabilities of two types of wavelength assignment heuristics and two types of link metrics are investigated.

It has been found that when the link utilization is low, there is no difference in blocking probability between first fit and random wavelength assignments. Almost all the connection requests can be satisfied. When link utilization becomes high, first fit performs better than random wavelength assignment. However, when link utilization is close to 100%, the difference between the two wavelength assignment heuristics becomes small again. At higher link utilizations, blocking increases exponentially when link utilization increases.

For the link metrics, we find, in a practical network such as a ring or a mesh network, the enhanced TAW link metric has a lower blocking probability than simple TAW. Further the enhanced TAW metric performs better at high link utilizations than the simple TAW metric particularly when the number of wavelengths per fiber is low. Conversely, at low link utilization, both metrics have very similar results whatever the number of wavelengths are used. Moreover, we also design a contention case to illustrate the use of TAKEDOWN message. Under this situation, the random wavelength assignment performs better than first fit wavelength assignment by spreading the chosen wavelengths.

In the simulation we only examine two simplest wavelength assignment heuristics: first-fit and random wavelength assignment schemes. They yield similar performance when associated with any routing and reservation schemes implemented in this work. We believe more heuristics should be tested and compared. Also, the wavelength convention (both full-convention and part-convention) can be considered. For the contention case, we wonder some improvements can be done on the first fit wavelength assignment heuristic, to make its blocking probability not worse than random heuristic.

The effects of chromatic dispersion and PMD are also included as constraints of the routing and wavelength assignment. These effects are significant for 10 Gb/s optical networks. The random scheme may not offer the best due to the effects of the residual dispersion. So the first-fit would offer better performance under these dispersion constraints. Nonlinear dispersion effects in RWA in optical networks such as self phase modulation will be included in future works. These effects are critical when the number of wavelength channels are increased, thus the total average power and intensity imposed on the fiber guiding region. This is very critical for network management in practice as the symbol rates are increased. The 40 Gb/s is emerging as the rate of next generation networks and the linear dispersion and nonlinear phase distortion are even much more critical with much lower tolerance for distortion and must be taken into account in the RAW management strategy. This will be reported in a future article. The randomness and statistical property of the PMD is included only with its mean value. This clearly allows the expected traffic performance and not to the details of the blocking that may happen higher or lower from time to time.

We have also assumed that the lightwaves of the lightpaths are amplitude modulated with either NRZ or RZ formats. It is more likely that the RZ format is used as it is the preferred format. Other modulation formats such as different phase shift keying [39], continuous phase shift keying or frequency shift keying have also been intensively investigated. They would be taken as the critical feature to combat impairments in long haul optical DWDM networks and RWA management strategy.

# 7. Acknowledgement

The author acknowledges the OMNeT++ programming assistance from C. Cieutat.

## 8. References

- [1] IEEE Workshop on Advanced Modulation Formats, San Francisco, USA, 2004.
- [2] N.Y. Ken-ichi Sato, et al., "GMPLS-Based Photonic Multilayer Router (Hikari Router) Architecture: An Overview of Traffic Engineering and Signaling Technology," IEEE Communications, vol. 40, no. 3, pp. 96–101, March 2002.
- [3] G.N. Rouskas and H.G. Perros, "A Tutorial on Optical Networks," Networking 2002 Tutorials, Lecture Notes in Computer Science, vol. 2497, 2002, pp. 155–193.

## ROUTING AND WAVELENGTH ASSIGNMENT IN GMPLS-BASED 10 GB/S ETHERNET LONG 167 HAUL OPTICAL NETWORKS WITH AND WITHOUT LINEAR DISPERSION CONSTRAINTS

- [4] R. Ramaswami and K.N. Sivarajan, "Optical Networks: A Practical Perspective," 2nd Edition, Morgan Kaufmann, 2002.
- [5] P.A. Smith, et al., "Generalized Multi-Protocol Label Switching (GMPLS) Architecture," draft-ietf-ccampgmpls-architecture-02.txt, 2002.
- [6] K. Kompella, et al., "Routing Extensions in Support of Generalized MPLS," draft-ietf-ccamp-gmpls-routing-04.txt, 2002.
- [7] D. Katz, D. Yeung, and K. Kompella, "Traffic engineering extensions to OSPF," draft-katz-yeung-ospf-traffic-10.txt, December 2003.
- [8] H. Smith and T. Li, "IS-IS Extensions for Traffic Engineering," draft-ietf-isis-traffic-05.txt, August 2003.
- [9] L. Berger, et al., "Generalized MPLS Signaling Functional Description," draft-ietf-mpls-generalizedsignaling-08.txt, 2002.
- [10] L. Berger, et al., "Generalized MPLS Signaling RSVP-TE Extensions," draft-ietf-mpls-generalized-rsvpte-07.txt, 2002.
- [11] B. Jamoussi, ed., et al., "Constraint-Based LSP Setup using LDP," RFC 3212, IETF Standards Track.
- [12] J. Varga, OMNeT++, 2003. http://whale.hit.bme.hu/omnetpp/.
- [13] D. Banerjee and B. Mukherjee, "A Practical Approach for Routing and Wavelength Assignment in Large Wavelength-Routed Optical Networks," IEEE Journal Selected Areas in Communications, vol. 14, no. 5, pp. 903–908, June 1996.
- [14] H. Zang, J.P. Jue, and B. Mukherjee, "A Review of Routing and Wavelength Assignment Approaches for Wavelength-Routed Optical WDM Networks," Optical Networks, vol. 1, no. 1, pp. 47–60, January 2000.
- [15] Apostolopoulos, et al., "QoS Routing Mechanisms and OSPF Extensions," RFC 2676, 1999.
- [16] R. Dutta and G.N. Rouskas, "A Survey of Virtual Topology Design Algorithms for Wavelength Routed Optical Networks," Optical Networks, vol. 1, no. 1, pp. 73–89, January 2000.
- [17] G. Xiao and Y. Leung, "Algorithms for Allocating Wavelength Converters in All-Optical Networks," IEEE/ACM Transactions on Networking, vol. 7, no. 4, pp. 545–557, August 1999.
- [18] E. Leonardi, M. Mellia, and M.A. Marsan, "Algorithms for the Logical Topology Design in WDM All-Optical Networks," Optical Networks, vol. 1, no. 1, pp.35–46, January 2000.
- [19] Y. Zhang, et al., "An Efficient Heuristic for Routing and Wavelength Assignment in Optical WDM Networks," IEEE ICC 2002, vol. 5, pp. 2734–2739, May 2002.
- [20] A. Birman, "Computing Approximate Blocking Probabilities for a Class of All-Optical Networks," IEEE Journal Selected Areas in Communications, vol. 14, no. 5, pp. 852–857, June 1996.
- [21] R. Ramaswami and K.N. Sivarajan, "Design of Logical Topologies for Wavelength Routed Optical Networks," IEEE Journal Selected Areas in Communications, vol. 14, no. 5, pp. 840–851, June 1996.

- [22] O. Gerstel and S. Kutten, "Dynamic Wavelength Allocation in All-Optical Ring Networks," IEEE ICC'97, vol. 1, pp. 432–436, June 1997.
- [23] J. Moy, OSPF Version 2, RFC 2328, 1998.
- [24] T.F. Asztalos, N.M. Bhide, and K.M. Sivalingam, "Adaptive Weight Functions for Shortest Path Routing Algorithms for Multi-Wavelength Optical WDM Networks," IEEE ICC'2000, New Orleans, LA, pp. 1330–1334, June 2000.
- [25] O. Gerstel, "On The Future of Wavelength Routing Networks," IEEE Network, vol. 10, no. 6, pp. 14–20, November–December 1996.
- [26] P.H. Ho and H.T. Mouftah, "Path Selection with Tunnel Allocation in the Optical Internet Based on Generalized MPLS Architecture," IEEE ICC 2002, vol. 5, pp. 2697– 2701, May 2002.
- [27] J. Zheng and H.T. Mouftah, "Routing and Wavelength Assignment for Advance Reservation in Wavelength-Routed WDM Optical Networks," IEEE ICC 2002, vol. 5, pp. 2722–2726, May 2002.
- [28] R.S. Ramaswami, "A. Distributed network control for wavelength routed optical networks," IEEE Infocom, San Francisco, CA, USA, pp. 138–147, March 1996.
- [29] R. Ramaswami and K.N. Sivarajan, "Routing and Wavelength Assignment in All-Optical Networks," IEEE/ACM Transactions on Networking, vol. 3, no. 5, pp. 489–500, October 1995.
- [30] Working Group: Architecture, O.P., PLL, & Signaling Working Groups, User Network Interface (UNI) 1.0.
- [31] Signaling Specification, 2001, Optical Internetworking Forum (OIF).
- [32] A. Birman and A. Kershenbaum, "Routing and Wavelength Assignment Methods in Single-Hop All-Optical Networks with Blocking," IEEE Infocom'95, vol. 2, pp. 431–438, April 1995.
- [33] B. Mukherjee, et al., "Some Principles for Designing a Wide-Area WDM Optical Network," IEEE/ACM Transactions on Networking, vol. 4, no. 5, pp. 684–695, October 1996.
- [34] I. Widjaja and A.I. Elwalid, "Study of GMPLS Lightpath Setup over Lambda-Router Networks," IEEE ICC 2002, vol. 5, pp. 2707–2711, May 2002.
- [35] D. Banerjee and B. Mukherjee, "Wavelength-Routed Optical Networks: Linear Formulation, Resource Budgeting Tradeoffs, and Reconfiguration Study," IEEE/ACM Transactions on Networking, vol. 8, no. 5, pp. 598–607, October 2000.
- [36] X. Zhang and C. Qiao, "Wavelength Assignment for Dynamic Traffic in Multi-fiber WDM Networks," IEEE ICC'98, pp. 479–485, October 1998.
- [37] G.P. Agrawal, "Fiber-optic communications systems," 3rd Edition, John Wiley & Sons, 2001.
- [38] I. Kaminov and T. Koch, "Optical Fiber Communications," vol. 3.1, Academic Press.
- [39] C. Wree, "Differential phase shift keying for long haul fiber optic transmission based on direct detection," Doctoral Thesis Dissertation, CA University of Kiel, 2002.



Jun LIU<sup>1</sup>, Dafang ZHANG<sup>2</sup>, Junhang JIN<sup>2</sup>

<sup>1</sup> School of Computer and communication, Hunan University, Changsha, P.R.China <sup>2</sup> School of software, Hunan University, Changsha, P.R.China *E-mail:* <sup>1</sup>a-7@163.com

# Abstract

Packet-pair sampling, also called probe gap model (PGM) is proposed as a lightweight and fast available bandwidth measurement method. But measurement tools based on PGM gives results with great uncertainty in some cases. PGM's statistical robustness has not been proved. In this paper we propose a more precise statistical model based on PGM. We present the new approach by using probability distribution and statistical parameters. We also investigate the use of a PGM bandwidth evaluation method considering a non-fluid cross traffic and present the alternative approach where the bursty nature of the probed traffic could be taken into account. Based on the model, measurement variance and sample size can be calculated to improve the measurement accuracy. We evaluated the model in a controlled and reproducible environment using NS simulations.

Keywords: Active Probing, Probe Gap Model, Available Bandwidth, Variance, Sample Size, Traffic Burst

# 1. Introduction

For many applications, such as congestion control [1], service level agreement verification [2], rate-based streaming applications [3], Grid applications [4,5], efficient and reliable available bandwidth measurement remains a very important goal. Researchers have been trying to create end-to-end measurement algorithms for available bandwidth over the last 15 years. The objective was to measure available bandwidth accurately, quickly, and without affecting the traffic in the path. However, the diversity of network conditions makes it a very challenging task. We also worked at bandwidth measurement in recent 5 years [6–10]. Existing measurement techniques fall into two broad categories [11]:

The first class of schemes is based on statistical cross-traffic model, known as the probe gap model (PGM), also called direct probing. Measurement tools such as Delphi [12] and Spruce [13] are based on PGM. The main component of PGM is the mathematical relation between the input and output rates of a probing packet pair, under a fluid traffic model. For PGM, the probing packets are sent to the Internet from the probe

host with a known separation. By measuring the change of the output gap, the utilization of the bottleneck link can be calculated.

The second class of schemes is the probe rate model (PRM), also known as iterative probing. PRM is based on the concept of self-induced congestion. Pathload [14], PathChirp [11], PTR [16], and TOPP [17] use the PRM model.



Figure 1. PGM model

Our analysis in this paper focuses on PGM. PGM based tools sample the arrival rate at the bottleneck by sending pairs of packets spaced so that the second probe packet arrives at a bottle-neck queue before the first packet departs the queue. These tools then calculate the number of bytes that arrived at the queue between the two probes from the inter-probe spacing at the receiver. A PGM tool computes the available bandwidth as the difference between the path capacity and the arrival rate



at the bottleneck.

Figure 1 illustrates a typical PGM probing tool: A probe pair is sent with time gap  $\mathbf{g}_i$ . Specifically, if the size of the probing packets is  $\mathbf{I}_c$  and the packet pair is sent to the path at the rate of the bottleneck capacity  $\mathbf{b}_0$  (back-to-back), the input gap of the packet pair is set to  $\mathbf{g}_i = \mathbf{I}_c / \mathbf{b}_0$ . If the queue does not become empty between the departure of the first probe packet in the pair and the arrival of the second probe packet, the output gap  $\mathbf{g}_0$  would be the time taken by the bottleneck to transmit the second probe packet of the pair and the cross traffic arrived during  $\mathbf{g}_i$ . Thus, the time to transmit the cross traffic is  $\mathbf{g}_0 - \mathbf{g}_i$ , and the rate of the cross-traffic is :

$$b_c = \frac{g_o - g_i}{g_i} b_o \tag{1}$$

The available bandwidth is:

$$A = b_o - b_c = \left(1 - \frac{g_o - g_i}{g_i}\right) b_o$$
<sup>(2)</sup>

The PGM approaches assume [13]:

- 1) FIFO queuing at all routers along the path;
- 2) Cross traffic follows a fluid model (Non-probe packets have an infinitely small packet size);
- 3) Average rates of cross traffic change slowly and is constant for the duration of a single measurement.
- 4) The single bottleneck is both the narrow and tight link for that path.
- 5) The cross traffic follows the same path with the measurement traffic [18].

These assumptions are necessary for analysis but the model might still work even when some of the assumptions do not hold [16]. In this paper, we are questioning the assumption 2:

Published paper only illustrates PGM in its simplest form with stationary and fluid cross traffic. In real tests [19], the time gap  $\mathbf{g}_0$  between probing packets will grow discretely because long or short cross traffic packets are inserted between the probing packets. Several works in this field have been published [13,16,20,21] and all show very similar features. Paper [20] published a detailed analysis of paths characteristics and its influence on  $\mathbf{g}_0$  through lots of experiments. But no solid theory can be found to support these results.

To cope with the burst of cross traffic, both the PGM and PRM tools use a train of probe packets to generate a single measurement. They use statistical methods to estimate the cross traffic, computes the available bandwidth from the average of several sample measurements. Spruce [13] averages individual samples using a sliding window over 100 packets. Abing [20] sends 20 probes for a single test. Why the statistical averages can reflect the real cross traffic is still not clarified in published research. We put all related symbols in Table1. If  $E(g_0)$  is the excepted value of output gap,  $b_c$  is the cross traffic throughput,  $b'_c$  is probe result, we get:

$$b_c = \frac{E(g_o) - g_i}{g_i} b_o$$
(3)

$$b_c' = \frac{\overline{g}_o - g_i}{g_i} b_o \tag{4}$$

Equation (4) and ideal PGM equation (1) are different. After above analysis, we have got several questions:

First, PGM should be statistical robust even the cross traffic follows do NOT follows fluid model. Why the statistical averages can reflect the real cross traffic load?

Second, PGM still has great uncertainty in accuracy. Can we get the variance of PGM test?

Third, what should be the reasonable sample size of PGM probing on a specific network environment?

Fourth, how much is the PGM's accuracy affected by traffic burst?

This paper is organized as follows: In Section 2, we present a new approach for analyzing the PGM model by using probability distribution and statistical parameters on CBR traffic. Section 3 is aimed to derive the coefficient of variation (CV) of the probe result for the crossing traffic. The theoretical study is used to evaluate the impact on CV due to different characteristics of probing and crossing traffic. We present background traffic burst measurement methodology in Section 4. Section 5 is the NS-2 simulation. The paper concludes in Section 6.

Table 1. Units for PGM

Symbol	Quantity	Conversion
$\mathbf{g}_{o}$	The output gap	
l <sub>c</sub>	Probing packet length	
$\mathbf{g}_{\mathbf{i}}$	The initial gap	$\mathbf{g}_{i} = \frac{\mathbf{l}_{c}}{\mathbf{b}_{o}}$
b <sub>o</sub>	Bottleneck bandwidth	
l <sub>p</sub>	Traffic packet length	
b <sub>c</sub>	Cross traffic throughput	
$\mathbf{g}_{\mathrm{on}}$	Time gap of ON period	$\mathbf{g}_{\mathrm{on}} = \frac{\mathbf{l}_{\mathrm{p}}}{\mathbf{b}_{\mathrm{in}}}$
$\mathbf{g}_{\mathrm{off}}$	Time gap of OFF period	$\mathbf{g}_{\mathrm{off}} = \frac{\mathbf{l}_{\mathrm{p}}}{\mathbf{b}_{\mathrm{c}}} - \frac{\mathbf{l}_{\mathrm{p}}}{\mathbf{b}_{\mathrm{in}}}$
Т	The statistic period	$T=g_{on}+g_{off}=\frac{l_{p}}{b_{c}}$
Δt	Offset of first probing packet in <b>T</b>	· ·
g <sub>c</sub>	Time for bottleneck to process one traffic packet	$\mathbf{g}_{c} = \frac{\mathbf{l}_{p}}{\mathbf{b}_{o}}$

# 2. Modeling

## 2.1. Modeling of PGM Probing



Figure 2. PGM structural process model

Common study of the network layer always has focused on service models: the routing algorithms and the protocols. To analysis the interaction between the probing packets and the competing traffic, we have to consider the switching function of a router, find out the details of the actual queuing and transfer of the packets from incoming links to the outgoing links.

Generally, FIFO queuing is supported on the router's output port. The packet from different input queue is forwarded to the switching fabric and put on the output port follows the FIFO scheduling discipline. Figure 2 shows a scene where four cross traffic packets (black) and two probing packets (shaded black) at the front of two different input queues of a router are destined for the same output port.

Simple traffic models [22] such as ON/OFF sources have been very popular for describing Internet traffic flows. Informally these models assume that the traffic alternates between active states (ON periods) and idle states (OFF periods). During ON periods packets are sent at a constant rate, during OFF periods no packets are transmitted.

We use similar technique to analysis process model of PGM. The time segment for receiving one cross traffic packet is defined as  $g_{on}$ . The time segment between two adjacent cross traffic packets is defined as  $g_{off}$ . A queuing period T is defined to be the time segment from the first bit of cross traffic packet 1 received in queue 1 to the first bit of cross traffic packet 2 received in the same queue.

The switch fabric always chooses to transfer the firstcome packet from input queues to the output queues. On this case (Figure 2), the black packet in the up-left queue must wait the first probing packets (shaded black) to go first for it comes first. So all the cross traffic packet arrived during time gap  $g_i$  are "captured" by the output gap. We put an icon of eye  $\triangleleft$  to describe the fact.

Because how much the time gap is for the first probing packet being ahead of the first cross traffic packet is complete random, the probability of  $\Delta t$  fall across any point in a queuing period T is equally distributed. So  $\Delta t$  is a random variable with continuous uniform distribution. There are FOUR possible scenarios according to the relationship between  $g_{on}$  and  $g_i$ . We now investigate them one by one:

**Scenario1.**  $g_i < g_{on}$ , for  $\Delta t$  is a continuous random variable in [0,T]. We have to identify the three more specific conditions (Figure 3(a)):

- a) When  $\Delta t \in [0, T g_{on}]$ , no cross traffic packets are inserted between two probing packets, so  $g_{o} = g_{in}$ .
- b) When  $\Delta t \in [T g_{on}, T g_{on} + g_i]$ , if and only if one cross traffic packet is inserted between two probing packets,  $g_0 = g_i + g_c$ .
- c) When  $\Delta t \in [T g_{on} + g_i, T]$ , no cross traffic packets are inserted between two probing packets,  $g_0 = g_{in}$ .

Conditions a) – c) are summarized as the probability distribution Figure 3(a). We can find that  $\mathbf{g}_0$  is a continuous random variable. Then the expected value of output gap  $\mathbf{E}(\mathbf{g}_0)$  is:

$$E(g_{o}) = \int_{0}^{T-g_{on}} g_{o}dt + \int_{T-g_{on}+g_{i}}^{T-g_{on}+g_{i}} g_{o}dt + \int_{T-g_{on}+g_{i}}^{T} g_{o}dt$$
$$= \frac{(T-g_{i})g_{i} + g_{i}(g_{i}+g_{c})}{T}$$
$$= \frac{l_{c}}{b} \left(1 + \frac{b_{c}}{b}\right)$$



Figure 3. PGM statistical model and probability distribution curve

**Scenario2.**  $T > g_i > g_{on}$ , (Figure 3(b)):

- a) When  $\Delta t \in [0, \mathbf{g}_i \mathbf{g}_{on}]$ , If and only if one cross traffic packet is inserted between two probing packets,  $\mathbf{g}_0 = \mathbf{g}_i + \mathbf{g}_c$ .
- b) When  $\Delta t \in [\mathbf{g}_i \mathbf{g}_{on}, \mathbf{T} \mathbf{g}_{on}]$ , no cross traffic packets are inserted between two probing packets,  $\mathbf{g}_0 = \mathbf{g}_i$ .
- c) When  $\Delta t \in [T g_{on}, T]$ , if and only if one cross traffic packet is inserted between two probing packets,  $g_o = g_i + g_c$ .

$$E(g_o) = \int_0^{g_i - g_{on}} g_o dt + \int_{g_i - g_{on}}^{T - g_{on}} g_o dt + \int_{T - g_{on}}^T g_o dt$$
$$= \frac{(T - g_i)g_b + g_i(g_b + g_c)}{T}$$
$$= \frac{l_c}{b_o} \left(1 + \frac{b_c}{b_o}\right)$$

**Scenario3.**  $(n+1)T > g_i > g_{on} + nT.$  (Figure 3(c)) :

- a) When  $\Delta t \in [0, g_i nT g_{on}]$ , n+1 cross traffic packets are inserted between two probing packets.  $g_o = g_i + (n+1)g_c$ .
- b) When  $\Delta t \in [\mathbf{g}_i \mathbf{n}T \mathbf{g}_{on}, T \mathbf{g}_{on}]$ , n cross traffic packets are inserted between two probing packets.  $\mathbf{g}_o = \mathbf{g}_i + \mathbf{n}\mathbf{g}_c$ .
- c) When  $\Delta t \in [T g_{on}, T]$ , n+1 cross traffic packets are inserted between two probing packets.  $g_o = g_i + (n+1)g_c$ .

$$E(g_o) = \int_0^{g_i - nT - g_{on}} g_o dt + \int_{g_i - nT - g_{on}}^{T - g_{on}} g_o dt + \int_{T - g_{on}}^T g_o dt$$
$$= \frac{((n+1)T - g_i)(g_b + ng_c) + (g_i - nT)(g_b + (n+1)g_c)}{T}$$
$$= \frac{l_c}{b_o} \left(1 + \frac{b_c}{b_o}\right)$$

Scenario 4.  $nT < g_i < g_0 + nT$ .(Figure 3(d))

- a) When  $\Delta t \in [0, T g_{on}]$ , n cross traffic packets are inserted between two probing packets.  $g_0 = g_i + ng_c$ .
- b) When  $\Delta t \in [T g_{on}, g_i (n 1)T g_{on}]$ , n+1 cross traffic packets are inserted between two probing packets.  $g_0 = g_i + (n+1)g_c$ .
- c) When  $\Delta t \in [\mathbf{g}_i (n-1) \mathbf{T} \mathbf{g}_{on}, \mathbf{T}]$ , n cross traffic packets are inserted between two probing packets.  $\mathbf{g}_o = \mathbf{g}_i + n\mathbf{g}_c$ .

$$\begin{split} E(g_o) &= \int_0^{T-g_{on}} g_o \, dt + \int_{T-g_{on}}^{g_i - (n-1)T-g_{on}} g_o \, dt + \int_{g_i - (n-1)T-g_{on}}^T g_o \, dt \\ &= \frac{\left( (n+1)T - g_i \right) (g_b + ng_c) + (g_i - nT) (g_b + (n+1)g_c)}{T} \\ &= \frac{l_c}{b_o} \left( 1 + \frac{b_c}{b_o} \right) \end{split}$$

We can summarize from scenarios 1 to 4 that under any circumstance,  $\mathbf{E}(\mathbf{g}_0)$  is constant. The equation (3) is proved. It can also tell us the PGM model is statistical robust.

## 2.2. PGM Probing Variance

This section is aimed to derive the variance of the probe result. Variance can be calculated as:  $D(b_c)=E(b_c^2)$  –

 $(\mathbf{E}(\mathbf{b}_{c}))^{2}$ . From equation (3), we get:  $\mathbf{b}_{c} = \mathbf{E}(\mathbf{g}_{o}) \frac{\mathbf{b}_{o}^{2}}{\mathbf{l}_{c}} - \mathbf{b}_{o}$  So:

$$E(b_{c}^{2}) = E\left(\frac{(E(g_{o}))^{2}b_{o}^{4}}{l_{c}^{2}}\right) - E\left(\frac{2E(g_{o})b_{o}^{2}}{l_{c}}\right) + b_{o}^{2}$$
(5)

To get  $\mathbf{E}(\mathbf{b_c}^2)$ , we divide equation (5) in to 3 parts, then calculate them one by one. Consider below two scenarios:

**Scenario1.**  $g_i < T$ , according to probability distribution of  $g_0$  (Figure 3(a), Figure 3(b)), we get:

$$E\left(\frac{\left(E(g_{o})\right)^{2}b_{o}^{4}}{l_{c}^{2}}\right) = b_{o}^{4}E\left(\frac{E^{2}(g_{o})}{l_{c}^{2}}\right)$$
$$= b_{o}^{4}\left(\int_{0}^{T-g_{on}}\left(\frac{g_{o}}{l_{c}}\right)^{2}dt + \int_{T-g_{on}}^{T-g_{on}+g_{l}}\left(\frac{g_{o}}{l_{c}}\right)^{2}dt + \int_{T-g_{on}+g_{l}}^{T}\left(\frac{g_{o}}{l_{c}}\right)^{2}dt\right)$$
$$= \frac{l_{p}b_{c}b_{c}}{l_{c}} + 2b_{o}b_{c} + b_{o}^{2}$$

(6)

For the same probability distribution curve, we can get:

$$E\left(\frac{2E(g_{o})b_{o}^{2}}{l_{c}}\right) = 2b_{o}b_{c} + 2b_{o}^{2}$$
(7)

Then from equation (5)–(7), we get:

$$D(b_{c}) = E(b_{c}^{2}) - (E(b_{c}))^{2} = \frac{l_{p}b_{o}b_{c}}{l_{c}} - b_{c}^{2}$$
(8)

 $\label{eq:scenario2.nT} \begin{array}{l} \text{Scenario2. nT} < g_i < g_{on} + \text{ nT } \text{ or } (n{+}1)\text{T} > g_i > g_{on} + \\ n\text{T} \ (n{=}1,2{\cdots}): \end{array}$ 

According to  $\mathbf{g}_{0}$  probability distribution (Figure 3(c), Figure 3(d)):

$$E\left(\frac{\left(E(g_{o})\right)^{2}b_{o}^{4}}{l_{c}^{2}}\right)$$

$$=-\frac{n^{2}b_{o}^{2}l_{p}^{2}}{l_{c}^{2}}-\frac{nb_{o}^{2}l_{p}^{2}}{l_{c}^{2}}+\frac{2nl_{p}b_{o}b_{c}}{l_{c}}+\frac{l_{p}b_{o}b_{c}}{l_{c}}+2b_{o}b_{c}+b_{o}^{2}$$

$$E\left(\frac{2E(g_{o})b_{o}^{2}}{l_{c}}\right)=2b_{o}b_{c}+2b_{o}^{2}$$
(10)

From equation (5), (9), (10) we get:

$$D(b_{c}) = E(b_{c}^{2}) - (E(b_{c}))^{2}$$
  
=  $-\frac{n^{2}l_{p}^{2}b_{o}^{2}}{l_{c}^{2}} - \frac{nl_{p}^{2}b_{o}^{2}}{l_{c}^{2}} + \frac{2nl_{p}b_{o}b_{c}}{l_{c}} + \frac{l_{p}b_{o}b_{c}}{l_{c}} - b_{c}^{2}$  (11)

The probed packet number n can be calculated as:  $n=floor(\frac{l_c b_c}{l_n b_o})$ , from (11) we get:

$$D(b_{c}) = -\frac{floor\left(\frac{l_{c}b_{c}}{l_{p}b_{o}}\right)^{2} l_{p}^{2} b_{o}^{2}}{l_{c}^{2}} - \frac{floor\left(\frac{l_{c}b_{c}}{l_{p}b_{o}}\right) l_{p}^{2} b_{o}^{2}}{l_{c}^{2}} + \frac{2 floor\left(\frac{l_{c}b_{c}}{l_{p}b_{o}}\right) l_{p} b_{o} b_{c}}{l_{c}} + \frac{l_{p}b_{o}b_{c}}{l_{c}} - b_{c}^{2}$$
(12)

When n=0, the equation (12) becomes the same to equation (8). So equation (12) can be treated as the uniform equation for variance calculation under all conditions.

#### 2.3. Coefficient of Variation

Only specifying the standard deviation is more or less useless without the additional specification of the mean value. It makes a big difference if  $D(b_c)=5$  with a mean of  $E(b_c)=100M/s$  or  $E(b_c)=3M/s$ . Relating the standard deviation to the mean resolves this problem. The coefficient of variation represents the ratio of the standard deviation to the mean, and it is a useful statistic for comparing the degree of variation from one data series to another, even if the means are drastically different from each other. The coefficient CV is defined

by: 
$$CV = \frac{\sqrt{Variance}}{Mean}$$
, so for PGM probing:

$$CV_{bc} = \frac{\sqrt{D(b_c)}}{E(b_c)}$$

$$= \frac{1}{b_c} \sqrt{\frac{floor\left(\frac{l_c b_c}{l_p b_o}\right)^2 l_p^2 b_o^2}{l_c^2} - \frac{floor\left(\frac{l_c b_c}{l_p b_o}\right) l_p^2 b_o^2}{l_c^2} + \frac{2floor\left(\frac{l_c b_c}{l_p b_o}\right) l_p b_c}{l_c} + \frac{l_p b_o b_c}{l_c} - b_c^2}$$
(13)

# 2.4. CV Curve

We use graphical methods to analysis CV according to equation (13). We put four CV curves on Figure 4 to explore how network parameters can affect the measurement accuracy.

1) Figure 4(a): CV(b<sub>c</sub>) and traffic packet length

Figure 4(a) clearly shows that the  $CV(b_c)$  shows an upward trend as traffic packet length increases . There is some variability about this trend, with some CV increasing over and some decreasing.

2) Figure 4(b): CV(b<sub>c</sub>) and cross traffic

In fact cross traffic  $\mathbf{b}_c$  is our goal of probing. Analysis here is only for showing the relationship between  $\mathbf{b}_c$  and  $\mathbf{CV}(\mathbf{b}_c)$ .Note that if the  $\mathbf{b}_c$  is very small,  $\mathbf{CV}(\mathbf{b}_c)$  becomes very big .That means: Probing with accuracy for small cross traffic is very difficult.

3) Figure 4(c): CV(b<sub>c</sub>) and probing packet length

Figure 4(c) shows that  $CV(b_c)$  a downward trend. If the probing packet length  $I_c$  is very small,  $CV(b_c)$ becomes very big. That means: Small packet is not suitable for probing.

4) Figure 4(d): CV(b<sub>c</sub>) and bottleneck capacity

Figure 4(d) shows that the  $CV(b_c)$  exhibit an upward trend as bottleneck capacity increases.



Figure 4. CV curve

## 3. Sample Size

Sample size is the number of observations in a sample. Lindebergh-Levy Theorem [15] describes the large sample behavior of random variables that involve sums of variables. The theorem is often written as  $\sqrt{N(\bar{X}_{N},\mu)} \rightarrow N(0,\sigma^2)$  which points out that  $\bar{X}_N$  converges to its mean  $\mu$  at exactly the rate as  $\sqrt{N}$  increases, so that the product eventually "balances out" to yield a random variable with normal distribution. The proof of this result is a bit more involved, requiring manipulation of characteristic functions, and will not be presented here.

From above analysis, obviously we can treat distribution of  $\overline{g_o}$  as a random variable with normal distribution. From equation (4), we can know  $\overline{b_e}$  and  $\overline{g_o}$  share the similar distribution. Thus  $\overline{b_e}$  can be treated as a random variable with normal distribution. We can use below equation to calculate the sample sizes for variables of normal distribution:

$$N = D\left(b_c'\right) \frac{Z^2}{e^2} \tag{14}$$

where N is the sample size, Z is the confidence level, e
is the desired level of precision (in the same unit of measure as the variance), and  $D(b'_c)$  is the probe variance. We also can calculate the sample sizes from CV:

$$N = CV \left(b_c\right)^2 \frac{Z^2}{m^2}$$
(15)

where **m** is the desired relate level of precision, it is often expressed in percentage points relative to mean  $\overline{\mathbf{b}_{e}}$ .

# 4. Measurement over Bursty Network Traffic

#### 4.1. PGM over Internet

PGM is working perfect over the CBR traffic flows. But for Internet probing, how to estimate accuracy and precision of PGM is still an open question. Recent empirical studies have provided ample evidence that actual network traffic is self similar [23] or fractal in bursty nature over a wide range of time scales. These observations are useful to make PGM probing over Internet more significant to cope with bursty traffic. We suggest a solution to predict the variance of PGM. An observer can easily predicted traffic patterns from actual measured traffic traces to help PGM probing evaluation.

#### 4.2. Traffic Burst Coefficient

ON/OFF models assume that traffic alternates among an active state or ON period and idle state or OFF period. We define the traffic burst coefficient  $\mathbf{B}_u$  as: variance of OFF period  $\mathbf{D}(\mathbf{g}_{off})$ . We capture n+1 packets (n can be set to more than 5000 in our experiments) from the link running PGM probing (The sample is often got from MIB database of router of the link). The timestamp of first bit and last bit of packets arrived in router is recorded as  $\{\mathbf{t}_1, \mathbf{t}_2, \mathbf{t}_3, \mathbf{t}_4, \dots, \mathbf{t}_{2n-2}, \mathbf{t}_{2n-1}\}$ , and then we have:

$$B_{u} = D\left(g_{off}\right) = \sum_{j=1}^{n} \left(g_{off,j} - E\left(g_{off}\right)\right)^{2}$$
  
=  $\sum_{j=1}^{n} \left(t_{2j} - t_{2j-1}\right) - \frac{1}{n} \sum_{k=1}^{n} \left(t_{2k} - t_{2k-1}\right)^{2}$  (16)

#### 4.3. PGM CV Estimation

Clearly, the coefficient of variation depends on whether or not the network traffic being bursty. We can solve this problem by the following reasoning:

First we have to estimate the maximum and minimal value of  $\mathbf{B}_{\mathbf{u}}$ :

Minimal  $B_u$ : minimal  $B_u$  is zero, when traffic is CBR source. All packets are equality distributed.

Maximum  $\mathbf{B}_{\mathbf{u}}$ : Maximum burst means all n+1 packets are connected. The traffic self-similarity gives raise to structural models with the distinctive feature that their ON and or OFF periods are heavy tailed with infinite variance. So the Maximum  $\mathbf{B}_{\mathbf{u}}$  is:

$$Bu_{max} = D(g_{off})_{max} = n(E(g_{off}))^2 = \frac{1}{n} \sum_{k=1}^n (t_{2k} - t_{2k-1})^2$$
(17)

We define  $I'_p$  as the average length of the packet trains. It can be proved that  $I'_p$  is in direct proportion to with  $B_u$ . So:

$$l'_{p} \approx n E \left( l_{p} \right) \left( \frac{Bu \left( g_{off} \right)}{Bu \left( g_{off} \right)_{max}} \right)$$
(18)

 $I'_p$  can be use in place of  $I_p$  in equations to calculate the  $CV(b_c)$ . Thus  $CV(b_c)$  can be calculated as:

$$\begin{cases} When: 0 < l'_{p} < \frac{l_{c}b_{c}}{b_{o}}, CV'(b_{c}) = CV(b_{c})_{max} = 2\frac{b_{o}l'_{p}}{b_{c}} \\ When: l'_{p} > \frac{l_{c}b_{c}}{b_{o}}, CV'(b_{c}) = CV(b_{c})_{max} = \frac{1}{b_{c}}\sqrt{\frac{b_{c}b_{o}l'_{p}}{l_{c}} - b_{c}^{2}} \end{cases}$$
(19)

#### 5. Evaluation Methodology

#### 5.1. Test Bed Illustration

Our evaluation is based on NS-2 [24] simulations, since we need to carefully control the competing traffic load in the network. The topology is shown in Figure 5. In this topology, there are one probing source and one receiver. **Cbr1, Cbr2** and **Cbr3** are used to generate competing traffic. **R1** and **R2** are FIFO-based routers. **Probe sender** sends out UDP packets pairs, served as a PGM probing pair. The time gap  $g_i$  is calculated by packet length and bottleneck bandwidth to ensure the packets are back-to-back. The cross traffic is generated using CBR stream.



Figure 5. Testbed topology

Timer driven stratified random sampling [25] is used in all our tests, for there is often positive correlation between competing traffic packets within the sample space. We use a timer to trigger the sending of probing packets. When the timer expires, we send out a pair of probing packet.

The timer is set to a random time with an average of 50ms. Our test result indicates that the accuracy performance of stratified sampling technique is wellabove other sampling methods in eliminate the positive correlation, and improves measured precision greatly. Further discussion is beyond the scope of this paper.

Our evaluation includes two parts: continues probing tests and grouped probing tests.

First, we perform continues probing using 4000 packet pairs, focusing on the measurement accuracy and the convergence time. We inject cross traffic at a rate of **b**<sub>c</sub>=3Mb/s. The bottleneck bandwidth is set to 10Mb/s. Cross traffic packets are set to be 1500 bytes long .We sent probing packet every 50 ms, with the 1000 bytes probing packet. Like other PGM tools such as Spruce [13], we average individual samples using a sliding window over 360 packets. That is also why the first 18 second we get no output. From equation (15), we can estimate bandwidth output will be within boundary  $\pm 0.5217$ M/s with confidence level to 90%.

On grouped probing tests, we analyze how the factors such as probing packet size have effect on the measurement accuracy of PGM. We also study the accuracy of sample size prediction on a network path, and look into a related issue of the sample size prediction. There are 4 groups 32 combined tests conducted. In each group, we change different characteristics of probing and crossing traffic to evaluate its impact on **CV**.Sample size is calculated from equation (15). We chose the confidence interval to 0.10. From normal distribution table we get z=1.65, e=0.5217M/s,.The sample size is just 10 times the value of  $D(b'_c)$ .So the sample size calculation is simplified.

#### 5.2. Test Results

Figure 6 illustrates typical segments of continues probing results, plots the cross traffic over a period of 200 seconds measured by PGM probing. Form the figure we can find that the probe results  $\mathbf{b'_e}$  is a nice match for the actual competing traffic of 3M/s. The statistical robustness of PGM is verified. From the output data we can see the only a few point are found out of the boundary 0.5217M. That matches the preset confidence Level 90%.

Grouped probing tests result is shown in Table 2 and Figure 7.

In all 32 tests, 4 results (Dots in probe result curve) are found to be out of the boundary (Solid line in probe result curve) 0.5217M/s .That matches the preset confidence level 90% which assumes about 3 out of the boundary results.

The variance of a PGM probing is the variance to the

From these two experiments, the mathematical abstraction of section 2, 3 is validated.



# Figure 7. Grouped probing result

# 6. Conclusions

PGM probing method features a strong interplay between network and prober. Contributions of this paper are providing the engineers with both descriptive and analytical methods to dealing with the variability in PGM probing. We develop a more precise single-hop gap model that captures the relationship between the

# ON MODELING AND ACCURACY ANALYSIS OF THE AVAILABLE BANDWIDTH MEASUREMENT BASED-ON PACKET-PAIR SAMPLING

competing traffic throughput and the change of the packet pair gap for a single-hop network. We can use the model to understand, describe, and quantify important aspects of the PGM and predict the response from inputs. We explore how network parameters can affect the measurement accuracy. In the future we will work on how to use these analysis methodologies to conduct monitoring efforts on both research and commodity infrastructure.

# 7. References

- J. Hoe, "Improving the Start-up Behavior of a Congestion Control Scheme for TCP," Proceedings of ACM SIGCOMM, September 1996.
- [2] "Real-time SLA monitoring tools," http://www.nwfusion.com/news/tech/0115tech.html
- [3] J.C. Bolot and T. Turletti, "A Rate Control Mechanism for Packet Video in the Internet," Proceedings of IEEE INFOCOM, pp. 1216–1223, 1994.
- [4] "GGF Network Measurement Working Group (NMWG),"
  - http://nmwg.internet2.edu/
- [5] "EGEE R-GMA, Relational Grid Monitoring Architecture," http://www.r-gma.org/
- [6] W.W. Li, D.F. Zhang, and J.M. Yang, "Performance Analysis of End-to-end Path Capacity Measurement Tools," Journal of Computer Applications, October 2006.
- [7] W.W. Li, D.F. Zhang, and J.M. Yang, "On Evaluating the Differences of TCP and ICMP in Network Measurement," Computer Communications, pp. 428–439, February 2007.
- [8] K. Xie, D.F. Zhang, J.G. Wen, and G.G. Xie, "A Real-Time Network Monitor System Based on WinPcap," Journal of Hunan University (Natural Sciences), February 2006.
- [9] C. Fan, G.G. Xie, D.F. Zhang, and Z.C. Li, "Performance Analysis of HTTP Service Based on Network Active Measurement," Journal of Computer Research and Development, March 2005.
- [10] G.X. Zhang, D.F. Zhang, and G.G. Xie, J.H. Yang, "Internet Traffic Measurement and Characteristic Analysis on Output Link of Metro Area Network," Acta Electronica Sinica, November 2007.
- [11] V. Ribeiro, "Path Chirp: Efficient Available Bandwidth Estimation for Network Path," In: PAM., 2003.

- [12] V. Ribeiro, M. Coates, R. Riedi, S. Sarvotham, B. Hendricks, and R. Baraniuk, "Multifractal Cross-traffic Estimation," Proceedings of ITC Specialist Seminar on IP Traffic Measurement, September 2000.
- [13] J.Strauss, D. Katabi, and F. Kaashoek, "A measurement study of available bandwidth estimation tools," Proceedings of ACM SIGCOMM conference on Internet measurement, 2003.
- [14] M. Jain and C. Dovrolis, "Pathload: A Measurement Tool for End-to-End Available Bandwidth," In Passive and Active Measurements, Fort Collins, CO, March 2002.
- [15] P. Billingsley, "The Lindeberg-Levy theorem for martingales," The Proceeding of American Mathematical Society, pp. 788–792, 1961.
- [16] N. Hu and P. Steenkiste, "Evaluation and Characterization of Available Bandwidth Techniques," IEEE JSAC Special Issue in Internet and WWW Measurement, Mapping, and Modeling, 2003.
- [17] B. Melander, M. Bjorkman, and P. Gunningberg, "A New End-to-End Probing and Analysis Method for Estimating Bandwidth Bottlenecks," in Global Internet Symposium, 2000.
- [18] L. Lao, C. Dovrolis, M.Y. Sanadidi, "The probe gap model can underestimate the available bandwidth of multihop paths," ACM SIGCOMM Computer Communications Review, vol. 36, pp. 29–34, 2006.
- [19] W.W. Li, J.F. Wang, G.G. Xie, and D.F. Zhang, "An IPDV Measurement Method Based-on Packet-Pair Sampling," Journal of Computer Research and Development, August 2004.
- [20] J. Navratil, "ABwE: A Practical Approach to Available Bandwidth," in: PAM., 2003.
- [21] J.P. Bi, Q.L. Wu, and C. Zhong, "Measurement and Analysis of Internet Delay Bottlenecks," Chinese Journal of Computers, April 2003.
- [22] D. Anick, D. Mitra, and M. Sondhi, "Stochastic Theory of a Data Handling System with Multiple Sources," J STJ, 61(8), pp. 1871–1894, 1982.
- [23] W.E. Leland, M.S. Taqqu, W. Willinger, and D.V. Wilson, "On the Self-Similar Nature of Ethernet Traffic," IEEE/ ACM Transactions on Networking, 2(1): pp. 1–15, February 1994.
- [24] NS-2, http://www.isi.edu/nsnam/ns.
- [25] K. Claffy, G. Polyzos, and H. Braum, "Application of Sampling Methodologies to Network Traffic Characterization," Computer Communication Review, 23(4): pp. 194–20, 1993.

# J. LIU ET AL.

# Table 2. Grouped probing result

Group	Probing	Bottlenec	Traffic	Cross	Variance	Sample	Variance	Probing	Absolute
	packet	k Capacity	packet	traffic		size	probe	result	error
	(Mb)	(Mb/s)	length	(Mb/s)			result	(Mb/s)	(Mb/s)
			(Mb)						
1	0.008	10	0.0008	3	0	1	0	3	0
1	0.008	10	0.0024	3	0	1	0	3	0
1	0.008	10	0.004	3	6	60	5.9167	3.083	0.083
1	0.008	10	0.0056	3	12	120	11.975	2.975	0.025
1	0.008	10	0.0072	3	18	180	16.95	2.65	0.35
1	0.008	10	0.008	3	21	210	21.381	3.095	0.095
1	0.008	10	0.0096	3	27	270	23.933	2.4889	0.5111
1	0.008	10	0.012	3	36	360	34.875	2.875	0.125
2	0.008	10	0.012	1	14	140	13.536	0.964	0.036
2	0.008	10	0.012	2	26	260	23.038	1.73	0.27
2	0.008	10	0.012	3	36	360	34.875	2.875	0.125
2	0.008	10	0.012	4	44	440	44.398	4.0568	0.0568
2	0.008	10	0.012	5	50	500	47.5	4.5	0.5
2	0.008	10	0.012	6	54	540	53.5	5.833	0.167
2	0.008	10	0.012	7	56	560	55.91	6.91	0.09
2	0.008	10	0.012	8	56	560	56.0447	7.955	0.045
3	0.004	10	0.012	3	81	810	97.482	3.687	0.687
3	0.0056	10	0.012	3	55.28	553	45.535	2.368	0.632
3	0.0064	10	0.012	3	47.25	472	45.467	2.86	0.14
3	0.0072	10	0.012	3	41	410	42.333	3.125	0.125
3	0.008	10	0.012	3	36	360	34.875	2.875	0.125
3	0.0096	10	0.012	3	28.5	285	29.526	3.158	0.158
3	0.0104	10	0.012	3	25.61	256	27.401	3.322	0.322
3	0.012	10	0.012	3	21	210	19.286	2.571	0.429
4	0.008	4	0.012	3	9	90	9	2.966	0.034
4	0.008	5	0.012	3	13.5	135	13.75	3.1667	0.1667
4	0.008	8	0.012	3	27	270	26.6	2.933	0.067
4	0.008	10	0.012	3	36	360	34.875	2.875	0.125
4	0.008	12	0.012	3	45	450	37.156	2.346	0.654
4	0.008	15	0.012	3	58.5	585	55.962	2.846	0.154
4	0.008	18	0.012	3	72	720	77.514	3.261	0.261
4	0.008	20	0.012	3	81	810	103.5	3.938	0.938



# **Streaming Multimedia over Wireless Mesh Networks**

David Q. LIU, Jason BAKER

Department of Computer Science Indiana University – Purdue University Fort Wayne Fort Wayne, IN 46805, USA E-mail: {liud, bakejm01}@ipfw.edu

# Abstract

Wireless mesh network (WMN) research is an emerging field in network communications. However, WMNs pose several difficulties in the transmission of information, especially time critical applications such as streaming video and audio. In this paper, we provide an overview of several research papers which utilize mesh networks for streaming multimedia. We compare the results of the research and the significance they bring to the field of wireless mesh networks. We then provide possible directions for future research into wireless mesh networks as they apply to streaming multimedia.

Keywords: Wireless Mesh Networks, Multimedia, Video Quality

# 1. Introduction

Wireless mesh network (WMN) [1] research is an emerging field of study in network communications. Wireless mesh networks are dynamically self-organized and self-configured where the node in the network automatically establishing an ad hoc network and maintain the mesh connectivity. According to [1], there are three classes of WMNs:

- Infrastructural Backbone WMNs: mesh routers form an infrastructure for mesh clients to connect to them.
- Client WMNs: mesh clients constitute the actual network to perform routing and configuration functionalities as well as supporting end-user applications.
- Hybrid WMNs: the mesh router infrastructure is combined with client meshing as shown in Figure 1[1].
  - WMNs have the following characteristics:
- Support multi-hop wireless networking
- Support for ad hoc networking
- Self-form, self-heal, and self-organize
- Support minimal mobile routers and stationary or mobile clients
- Support both backhaul access to the Internet and peer-to-peer communications
- Provide power efficient protocols for mesh clients

• Interoperate with existing wireless networks such as Wi-Fi, WiMAX, Zig-Bee, and cellular networks

Recently research has bee done in the area of multimedia stream over wireless mesh networks. A multi-source multipath video streaming system [2] is proposed to support concurrent video-on-demand over wireless mesh network. Network coding [3] may be used to increase throughput over a wireless mesh network. Quality of Service (QoS) support is surveyed in [4]. Paper [5] proposes QoS routing in WMNs. Real-time video stream aggregation is studied in [6]. Results from a real testbed is reported in [7] about multimedia over wireless mesh networks. The core questions for multimedia streaming over wireless mesh networks is how to establish connection, maintain transmission of multimedia data, and achieve suitable quality across such networks. This paper, in Section 2, presents several techniques that will cover the topics of path determination. adaptive quality, and cross-layer information gathering. In Section 3 we compare the varied techniques across several dimensions to determine such attributes such as usefulness and efficiency. This is to be followed by areas for future research in Section 4. Finally, we will provide a brief conclusion stating the resultant findings in Section 5.

# 2. Existing Techniques for Multimedia Streams over Wireless Mesh Networks

There are a variety of techniques for transmitting multimedia streams across a given wireless mesh network. However, we will limit the discussion to three general areas: path determination techniques, adaptive quality, and cross layer information gathering. These techniques try to solve the problems that occur with improving the quality and performance of multimedia transmissions across wireless mesh networks.

#### 2.1. Path Determination Techniques

To transmit the data from the source to the destination requires the use of some form of path determination or routing algorithm. One such technique utilized in [8] is congestion-minimized routing. Congestion-minimized routing tries, as its namesake suggests, minimizing the congestion which is defined as the average delay per link. This is accomplished by dividing the total transfer rate into K sub-transfers of equal rate  $\Delta Rk$ , and assigning each sub-transfer to a given route. They minimize the equation

$$\min_{\rho_k} \sum_{(i,j)\in\rho_k} \frac{C_{ij}}{\left(C_{ij} - F_{ij}\right)^2} \Delta R_k$$
(1)

where  $\rho_k$  is the path from source to destination,  $C_{ij}$  is the capacity of the link from node *i* to node *j*, and  $\vec{F}_{ij}$ denotes the existing network traffic plus the prior subtransfer rates. To find the route for the sub-flow they make use of the capacity of the link and the current rate of traffic across that link. The routing problem in (1) is solved by utilizing the distributed Bellman-Ford algorithm, where each node has to store the minimum path cost from itself to the source and the link cost to all the neighbors. It is stated that solving (1) using Bellman-Ford will converge to the solution in a network of Nnodes and having diameter D within D rounds of information exchange. After the path has been found the destination sends a message down the reverse path to the source, thus giving the source the routing path for route Κ.

Paper [9] demonstrates three different path selection algorithms with varying levels of estimation utilized: end-to-end, localized, and estimation based. For all algorithms they assume that the topology does not change during video transmission, and that there is no contention for access to the medium. To accomplish this they employ the principles of HCCA protocol as applied to an IEEE 802.11e network. The application of this technique allows the scheduling of multiple flows creating an average transmission rate for each flow. The end-to-end algorithm uses complete information of the network so it must first generate the connectivity structure P for each node which has data to transmit. P is defined as all possible paths from source to destination without loops. Using P they find the minimum cost path using an algorithm which exhaustively searches all possible paths. The algorithm then finds the path which has the smallest estimated timing requirements for transmission. The information needed for this is transmitted using separate logical communication links between nodes. The next approach presented is a localized estimation, where only the link information of neighboring nodes is known. The rest of the path information is estimated based upon an approximation of several low layer data characteristics. The third technique presented is purely estimation based and uses an estimation technique similar to the localized version; however, it does not even keep track of the information on links between the nodes. Three path determination algorithms are discussed in [10]. One is a centralized algorithm which will not be discussed due to the fact it is almost exactly the same as the prior end-to-end algorithm from [9]. This is most likely due to the majority of authors being the same and the later publishing date of [10]. The other two are new peer-topeer (P2P) algorithms: distributed collaborative and distributed non-collaborative. Both approaches assume that there is enough available bandwidth for an overlay layer for communication of the network conditions between peers, and they both only run when a new flow is admitted, an existing flow leaves, or the topology of the network changes. The first is a collaborative distributed algorithm, which utilizes a local greedy approach. All nodes in the network must collectively sort the sub-flows. Sorting is done to satisfy the utility function for maximizing the total quality (MTQ) or the utility function for maximizing the minimum quality (MMQ). If they are utilizing the MTQ approach, then the sub-flows are sorted in descending order of  $\lambda_x/B_x$ , where  $\lambda_x$  is a flow specific parameter that depends on video characteristics and  $B_x$  is the rate requirement for the sub-flow. For the MMQ tactic they group the subflows by the quality layer and then sort by  $\lambda_x/B_x$ . After deciding and ordering all sub-flows the sender determines if there is any path to the destination. If no path exists, do not admit the video stream transmission to send, and cancel all sub-flows that depend on the denied sub-flow. If only one exists the transmission is started. If multiple paths exist, the path that leads to the smallest amount of introduced congestion is selected.

Two methods are offered for estimating congestion: bottleneck air-time congestion and mean end-to-end airtime congestion. Bottleneck congestion works using:

$$\varepsilon_x^i = \max_{v_a \in \rho_x^i} \{1 - \rho^a\}$$
(2)

where  $\rho_x^i$  is the *i*-th potential path for sub-flow *x*,  $v_a$  is node *a* along the path, and  $\rho^a$  is the fraction of total listening at  $v_a$ . Mean end-to-end works using the equation

$$\phi_x^i = \frac{1}{\left| p_x^i \right|} \sum_{v_a \in \rho_x^i} (1 - \rho^a)$$
(3)

Copyright © 2008 SciRes.

179

where  $|\rho_x^i|$  is defined as the number of nodes in the path and the rest of the parameters are the same as in (2). However, it defines the non-collaborative approach has each node sorting its own flows, before scheduling any other peers are allowed to schedule their flows to pass through the node.

Paper [11] presents no noteworthy developments for path determination; as it uses the UDP/IP protocol suite for routing and transmission of data. Therefore, it can be thought to be representative of the base line.

#### 2.2. Adaptive Quality

Adjusting the quality of the video is a must for utilizing available bandwidth, and for reducing congestion throughout the network. A technique used in [8] is to minimize the distortion of the encoded video, while limiting the congestion introduced. They achieve this by estimating the tradeoff between an increase in the rate and the decrease in the distortion. The given equation balances the rate versus the distortion

$$-\Delta D^{s} \approx \frac{\theta^{s}}{\left(R^{s} - R_{0}^{s}\right)^{2}} \Delta R^{s}$$
(4)

where  $-\Delta D^s$  is the distortion reduction for the stream,  $R^s$  is the encoding rate, and  $\theta^s$  and  $R_0^s$  are determined from trial encodings. At time intervals of size k, the source node increases the rate allocated to the stream and monitors the congestion versus the reduction in distortion. If the increase in rate isn't worth the reduction in distortion then the rate stays where it is.

There is a pre-agreed scaling factor  $\lambda$  for the congestion. This allows for a reduction in distortion as long as that is less than  $\lambda$  times the increase in congestion. This increase occurs until it reaches the optimal point, at which point it merely responds to the congestion present in order to raise or lower the distortion. This is guaranteed to converge for networks with fixed rates and link capacities.

In [9], we are presented with no adaptive quality changes for video transmission. Once a video is desired to be sent, the path provisioning scheme tries to find a path which meets the necessary requirements for the video transmission. If one can not be found, it does not send the video stream.

In [10], they demonstrate an interesting technique based on logical flows of data. Each sub-flow represents a partition in the actual quality of the video. This is represented by the equation

$$\hat{Q}_{y}(P_{y}) = Q_{y}^{0} + \sum_{\substack{f_{x} \in \Psi_{y} \\ a_{x} \neq -1}} \omega_{x} \lambda_{x} \log(B_{x})$$
(5)

For (5)  $Q_y^0$  is a parameter dependent upon the video, encoding parameters, etc.  $\omega_x$  is 1 for admitted sub-flows, 0 for non-admitted sub-flows, and -1 for rejected subflows.  $\lambda_x$  is a parameter like  $Q_y^0$  and it is dependent on the quality layer.  $B_x$  is the bit rate for the sub-flow. When a sub-flow is not admitted, all "enhancement layers" that depend upon that flow are not admitted as well. The parameter  $\omega_x$  allows ease of detection for a given non-admitted sub-flow. When a sub-flow is rejected  $\omega_x$  is set to -1 recursively for all dependent



Figure 1. Wireless Mesh Network [1]

$$U_{MTQ} = \sum_{y=1}^{N_p} \hat{Q}_y(\mathbf{P}_y)$$
(6)

where  $N_p$  is the total number of aggregate flows. By utilizing (5) MMQ can also be defined

$$U_{MMQ} = \min_{y} \left\{ \hat{Q}_{y}(\mathbf{P}_{y}) \right\}$$
(7)

These utility functions act as constraints, which allow the system to make quality of service decisions based upon individual sub-flows through the previously discussed path determination.

Paper [11] makes use of the standard MPEG video standard. The standard makes provisions for a quantization scale parameter (QSP), which allows the quality of the codec to be adjusted dynamically during the encoding process. A feedback formula is presented to achieve adaptive quality, which is defined as:

$$q = \Phi(s_q, a, j) = \min_{\overline{q} \in [1,31]}(\overline{q} : s_q + R_{a,j}(\overline{q}) \le \theta_F) \qquad (8)$$

where  $R_{a,j}(q)$  is the rate curve for the expected number of packets to encode for each frame type,  $s_q$  is the length of the transmission buffer before encoding the packet, and  $\theta_f$  is the queue length of the transmission buffer. Equation (8) allows us to make use of the rate curve to minimize the encoding QSP and maximize the PSNR, while maintaining a full transmission buffer.

#### 2.3. Cross Layer Information

Most of the techniques discussed in [8] utilize information from the bottom three to four layers of the OSI model. The problem is how to gather this information so that the path determination and adaptive quality algorithms can make use of it. Paper [8] gathers the busy, block and idle times which are referred to as  $T_{busy}$ ,  $T_{block}$ , and  $T_{idle}$  respectively. They also keep a running average of the video payload size for each stream over the time period known as  $B^s$ . The rate for a given stream on a node is:

$$F_n^s = \frac{B^s}{T_{busy} + T_{block} + T_{idle}}$$
(9)

The estimation for bandwidth capacity is:

$$C_n = \sum_{s} \frac{B^s}{T_{busy} + T_{block}}$$
(10)

Employing (9) and (10), an estimation of available bandwidth for a given stream at the given node is defined as:

$$C_n^s = C_n - \sum_{s' \neq s} F_n^{s'}$$
(11)

Estimation of the congestion increment, denoted as  $\Delta X_s$ , at a given node can be estimated using (9) and (10)

resulting in:

$$\Delta X_s = \sum_{n \in \mathbb{P}^s} \frac{C_n}{\left(C_n - \sum_s F_n^s\right)^2} \Delta R_s$$
(12)

which correlates the congestion increment against a given increase in the encoding rate.

Paper [9] gathers information for each link in the network. It utilizes the modulation, the bit error rate (BER), and the guaranteed bandwidth denoted as  $m(l_{i,j})$ ,  $e(l_{i,j})$ , and  $g(l_{i,j})$  respectively. The modulation is defined by the physical medium, but there is no mention of how the modulation is gathered. The probability of a BER is defined as:

$$e_{l_{i,j}}(L_{\nu}) = 1 - \left(1 - e(l_{i,j})\right)^{L_{\nu}}$$
(13)

where  $L_{\nu}$  is defined as the size of the MSDU in bits. This is due to the assumption that the BER is normally and randomly distributed. The guaranteed bandwidth is set up so as to follow the HCCA reservation rules. These parameters are used to estimate queuing delay for each link as:

$$d_{queue}(l_{i,j}) = \sum_{\forall u \in V_{queue}(l_{i,j})} d_{l_i,j}(L_u, t_{l_{i,j}}^{mean}(T_{p_i}^{max}))$$
(14)

where  $d_{l_{i,j}}(L_u, t_{l_{i,j}}^{mean}(T_{p_i}^{max}))$  defines the delay for the link between nodes *i* and *j* with MSDU of size  $L_u$ , with a mean retransmission time at the link and max retransmission time along the path denoted by  $t_{l_{i,j}}^{mean}(T_{p_i}^{max})$ . Equation (14) results in the queuing delay needed for the path determination algorithms in Section 2.1.

Information is gathered at the physical and network layers in paper [10]. The BER is also used in this paper, and is defined as

$$e\left(\theta_{x}^{a}\right) = \frac{1}{1 + e^{\mu(s-\delta)}}$$
(15)

where  $\mu$  and  $\delta$  are constants, *s* is the signal to noise ratio (SINR), and  $\theta_x^{a}$  is the physical layer mode at the given node *a* on stream *x*. This gives the expected goodput (throughput without errors) defined as:

$$\overline{G}_x^a = (1 - \mathcal{E}_x^a(L_x, \theta_x^a))T_x^a(\theta_x^a)$$
(16)

where the function  $\varepsilon_x^a$  is the probability of a bit error on a packet of size  $L_x$  bits (same equation as (14), and the function  $T_x^a$  is the physical layer transmission rate on the physical mode  $\theta_x^a$ . They then utilize the transmission service interval, listening service interval, and transmission service interval denoted as  $t_{SI}$ ,  $t_{SI}^{(RX)}$ , and  $t_{SI}^{(TX)}$  respectively. The assumption is made that the transmission service is much greater than the receiving service, implying that congestion is due to transmissions. Also recorded is the fraction of the listening time given

Copyright © 2008 SciRes.

180

to a stream for a node denoted as  $r_x^a$ . Using the gathered cross-layer information the admission of sub-flows is decided based on the inequality:

$$\overline{G}_{x}^{a}r_{x}^{a}\frac{t_{SI}^{(RX)}}{t_{SI}} < B_{x}$$
(17)

where if this is true the sub-flow can be allowed, otherwise it must be dropped.

## 3. Comparison

Given the wide variety of techniques previously described, this section will be devoted to comparing the overhead, benefits, and disadvantages of the techniques as they apply to situations dealing with streaming media.

#### **3.1.** Path Determination Techniques

The variety of path determination techniques explained in Section 2.1 have distinct advantages and disadvantages. To more objectively discuss these techniques they will be compared by the attributes over the domains: information complexity, algorithmic complexity, network overhead, congestion avoidance, and dynamic adaptation. These domains will be assigned ratings qualitatively from very bad, bad, average, and good, to very good.

Information complexity covers the storage and updateability of information, which widely varies depending on the technique used for each algorithm. Congestion-minimized routing is very good with regarding to informational complexity. This is due to the fact that it only requires estimations on the available bandwidth and utilized bandwidth on links between nodes as shown in (1). This information can easily be stored with the link information, and can be automatically updated by traffic passing through the node. End-to-end requires perfect knowledge of the entire network at each node and therefore has a very bad information complexity rating. This requires a massive amount of information to be stored at each node. Also, when the information changes all nodes must be updated. Localized attains a rating of good due to the fact that it only needs to store the information associated with the nodes one hop away in the paper, although since this grows as it looks farther away until the complexity reaches end-to-end. If it is just 1 hop away, the information is easily stored as in congestion-minimized routing. The estimation based approach receives a rating of very good for the fact that it stores only one piece of information regarding the link between nodes and estimates everything off of that one piece of information. Distributed collaborative path determination obtains an average rating. This is due to the fact that it must store information for each sub-flow passing through it; including listening times, and keep this information in

sorted order to determine when to drop quality layers. The distributed non-collaborative approach also receives an average rating, because it must also store the same amount of information as the non-collaborative. UDP/IP receives a very good rating due to the fact that no additional information has to be stored or updated.

Algorithmic complexity encompasses the run-time and ease of implementation, which again varies widely depending on the technique. Congestion-minimized routing obtains a rating of good; because of the distributed nature of the algorithm it can easily replace the distance metric used Bellman-Ford algorithm. However, it requires the use of an overlay network to keep nodes current on changes in topology, and this must be implemented to update the routes and their tables. End-to-end achieves a rating of very bad for algorithmic complexity due to the fact that it performs an exhaustive search of the path space from source to destination, and performs computations along the complete path. This results in a triple summation along different dimensions of activity to find the minimum path. The localized algorithm collects an average rating. This is due to the fact that it requires an exhaustive search of all neighbor nodes. However, it too requires an overlay network to convey information between the nodes, which is used to relay the information of the path up to the current determined node while determining which path to take. Estimation based gets a good rating due to the fact that it only requires the use of the overlay network to pass messages between nodes in the wireless mesh network. It is a simple algorithm of estimating the total path deadline and then chooses the route that minimizes the deadline criteria. The distributed collaborative algorithm gets a rating of good. This is due to the fact that the path provisioning algorithm only runs to meet significant changes in the network. However, each time a path is determined they have to see if the path can support the bit flow, incurring a tremendous amount of overhead. The ease of implementation is about the same between this and other algorithms presented due to the overlay layer used. The distributed collaborative algorithm scores a good rating. This is caused by the imposed requirement of an overlay network for notifications and path admission. However, it does not incur the complexity of communicating with its neighbor about what it is doing. UDP/IP receives a very good rating because it follows the simple pre-established algorithms provided on all routers.

Network overhead is related to informational complexity for any information transferred over the network. Congestion-minimized routing obtains an average rating for network overhead due to the fact that again it requires an overlay network which performs status message passing. It also requires that the Bellman-Ford algorithm be re-run each time a path is introduced. In addition, the algorithm takes a number of iterations of the Bellman-Ford algorithm equal to the diameter of the

graph to converge. End-to-end routing is the most costly in terms of overhead for the network. It must obtain the complete connectivity structure and all link information at each node for path determination. The amount of network overhead introduced by this is phenomenal and grows at an astonishing rate as nodes are added to the network, leading to its low rating of very bad. Localized only needs to transmit information between the neighboring nodes, so the amount of network overhead is low and the use of the overlay network does not significantly impact this technique. Estimation based acquires a favorable score for network overhead due to the fact that it transfers no information between nodes. It utilizes only link information that is set a priori to determine what path to take. Thus, estimation based gets a very good rating for network overhead because there is no overhead aside from the actual transfer. Distributed collaborative procures a rating of average for network overhead, because it is necessary to retrieve information from the network to find viable paths, and admit them into the network. This process has a high overhead in network message passing using the proscribed overlay network. The distributed non-collaborative technique receives a good rating because it is similar to the distributed collaborative in many respects, such as the use of an overlay network and path determination sans congestion avoidance. However, the individual nodes don't spend time messaging back and forth to allocate the flows reducing the overhead. UDP/IP has no network overhead aside from normal transmission costs; it also has no retransmission due to failed packets.

Congestion avoidance is a necessary feature for streaming video; if you stream video through a congested node then you will have jittery video or even lost streams. Congestion-minimized routing boasts a very good rating for congestion avoidance. The algorithm is designed to avoid the inherent problem by estimating the network congestion and finding the path through the network, described in (1), that minimizes congestion. End-to-end has a side-effect of congestion avoidance since it tries to minimize the queuing delay along the possible paths. By thus avoiding the heavily queued nodes, it avoids congestion. Therefore, end-toend receives a very good rating for congestion avoidance. The localized estimation uses the same technique as endto-end; however, it does not possess perfect knowledge of the network, and thus bad estimations may result in a quick build up in congestion. This method will nevertheless work for local congestion avoidance, so it receives a rating of good. The estimation based approach has no knowledge of network conditions and makes estimations for all links set up and congestion from (2) and mean end-to-end congestion from (3). The bottleneck congestion is useful for networks with bottlenecks, and mean end-to-end is useful in all other cases, both techniques provide excellent feedback on the network from the point view of reserving air space. This

leads to routing into congested areas quite easily, so the congestion avoidance for the estimation approach has a rating of very bad. The distributed collaborative technique receives a rating of very good for congestion avoidance. It gives two techniques for congestion avoidance bottleneck. Therefore, the necessity for a priori knowledge about bottlenecks detracts from the possible perfect score in congestion avoidance. Distributed non-collaborative gets a rating of very bad for implementation due to the fact that it does not implement any kind of congestion avoidance mechanism. UDP/IP also gets a low rating for congestion avoidance; by itself UDP/IP supports no inherent means for congestion avoidance. Since no other congestion avoidance means were discussed in [11], UDP/IP receives a rating of very bad.

Dynamic adaptation encompasses how responsive and accurate the algorithm works with changes to topology and data flows. Congestion-minimized routing adjusts the paths of video streams over time based on the conditions in the wireless mesh network. This automatic adjustment earns congestion-minimized routing a very good rating for the dynamic adaptation of the environment. End-to-end, localized and estimation based all assume that network topology is fixed while transmitting and that time is reserved for communication before transmission starts. These three algorithms are sorely lacking in the ability to respond to dynamic changes in the wireless mesh network. They do not reduce the current stream's bandwidth to accommodate new ones, and if a node becomes unresponsive during transmission there is no recovery for this happenstance. However, the overlay network provides a possibility of recovery if implemented in the future. As a result we assign all three of these techniques a rating of very bad Distributed collaborative and non-collaborative both do not perform dynamic adaptation of packet routing. Again, although the algorithms used and the overlay network potentially allow for extensibility into this area. Therefore, distributed collaborative and noncollaborative both receive a rating very bad also. UDP/IP inherently routes around areas where nodes are malfunctioning or have left the wireless mesh network; however, it does not take into consideration data flows already going through the network. Therefore UDP/IP receives a rating of good.

Table 1 shows the overall associated ranks for each technique and the critiqued areas. As can be easily seen some path determination techniques faired better than others. However, congestion-minimized routing, distributed collaborative, and the standard UDP/IP suite faired the best overall.

#### 3.2. Adaptive Quality

To compare the adaptive quality techniques from Section 2.2, we utilize the dimensions of scalability, algorithmic

	Critiqued Areas									
Technique	Information Complexity	Algorithmic. Complexity	Network Overhead	Congestion Avoidance	Dynamic Adaptation					
Congestion- Minimized Routing	very good	good	average	very good	very good					
End-To-End	very bad	very bad	very bad	very good	very bad					
Localized	good	average	good	good	very bad					
Estimation based	very good	very good	very good	very bad	very bad					
Distributed collaborative	average	good	average	very good	very bad					
Distributed non-collaborative	good	good	good	very bad	very bad					
UDP/IP	very good	very good	very good	very bad	good					

 Table 1. Path determination overview

complexity, and smoothness. Again they will be appraised on the same scale as used in Section 3.1.

Scalability refers to how well the technique scales in terms of the network size and the amount of traffic present in the wireless mesh network. The technique presented in [8] is extremely scalable well due to the fact that the technique only worries about its own traffic so each transmitting node is self monitoring. There is an overhead incurred by using the overlay network to transmit the needed information about congestion along the path. This is somewhere between constant and linear in terms of the growth of the network. This can be said because the addition of a node may or may not impact the path taken. If the path is lengthened then, the impact is linear; if the path stays the same, then it is constant. Also, (4) provides a way to calculate the given decoding distortion, which can then be utilized to optimize the system against the increase in congestion. This leads to even the network layer adapting to the introduction of streams of data as described in (1). This appears to be extremely scalable because it is also done per sending node, with overhead in the overlay network. Therefore, this technique receives a rating of good for scalability. Paper [9] presents no techniques for adaptive quality and thus obtains a rating of very bad for this section. In [10], distributed collaborative algorithm only, the flows are split into hierarchical layers equivalent to layers of quality or "enhancement layers." To make use of this two possible utility functions are described in (6) and (7) which do determine the MMT and MMQ respectively. Both of these equations utilize (5), which allows the wireless mesh network to determine the aggregate quality of a video stream across the network. Then, when utilizing the MMT or MMO utility functions the system can possibly self-optimize the flows through utilization of the overlay network. However, the overhead incurred on the overlay network would be quite extreme. This is because that both utility functions require perfect knowledge of the number of video streams through the network and all rejected "enhancement layers" need to be kept track of to reintroduce these video streams into the system. This can be deduced from both utility functions requiring perfect knowledge of the number of video streams moving through the network to make any decisions and from the fact that all rejected "enhancement layers" need to be kept track of to reintroduce them into the system. The first problem grows tremendously with respect to the network, whereas the second problem is only the responsibility of the sender nodes. This leads to a rating of average for scalability. Paper [11] presents a feedback approach local only to the system encoding video through (8). This leads to each individual encoder only adjusting its encoding rate without communicating with other nodes in the wireless mesh network. This problem scales perfectly as there is no communication between the nodes, therefore it receives a rating of very good.

Algorithmic complexity is defined as ease of implementation and speed of response when performing adaptive quality. Paper [8] has a quick initial response due to the nature of how the congestion and distortion converge. In the beginning, the decoding distortion rapidly decreases as you raise the encoding rate. Whereas, the introduced congestion in the network increases slowly in the beginning and increases rapidly after the initial low rates. Thus, the network quickly converges upon the goal encoding rate given the maximal network congestion. Ease of implementation is again a problem with the techniques described in [8]. For the algorithm to efficiently and quickly converge, the overlay network must be implemented and quickly pass the requisite knowledge to the appropriate nodes. This information passing overlay network is a significant burden to overcome, which is why it receives a rating of good for algorithmic complexity. Again, [9] obtains a meager rating of very bad for not having implemented any kind of adaptive quality. Paper [10], with respect to the collaborative approach, has a quick response due to several key features dealing from queue admittance discussed in path determination to the actual use of MMQ and MTQ to ensure quality of service. This quick convergence is due to the fact that for admittance into the network, a given sub-flow is ordered by the benefit to the aggregate video stream. Thus, all base flows will be added; then, "enhancement layers" will be added by order of contribution. However, this is detracted by the need for an overlay network to enforce the utility functions, ensure admittance of new sub-flows, and update all nodes in the wireless mesh network. These detractions result in a rating of good for algorithmic complexity. Paper [11] presents an adaptive quality technique which is apparently simplistic to implement and extremely responsive to the criteria presented. It is simple to implement because the algorithm runs at each sender and encodes at the given quality rate based upon the current transmission buffer capacity as determined by (8). It responds instantaneously to any change in the buffer. These attributes give this technique a rating of very good.

Before we can discuss smoothness, it must be defined with respect to streaming media over wireless mesh networks. Smoothness will be defined as the continuity of the quality as it changes over time. In [8] the smoothness is somewhat ideal. This can be attributed to the small changes in (4). The congestion is only incremented by a given amount for each rate increase. This algorithm works on k time steps, such that it only increases or decreases the rate over time based on the current congestion in the network. This leads to a very smooth curve of quality for small k and a very discontinuous curve for large k. For most applications, though, it would be safe to assume that k is much smaller than 2 seconds. Therefore, it is safe to say that a rating of very good is very accurate and deserved for this technique. Again, [9] receives a rating of very bad due to not implementing any techniques. For [10] - we once more only consider the distributed collaborative - it appears to have an ordered but non-uniform smoothness. This can be said because the system implicitly orders the quality layers by contribution to the quality dependent upon the chosen utility function. Ergo, the quality is as smooth as the quality layer added or removed based upon network conditions. The problem associated with this is the varying delta for the improvement in quality. For more important base flows, the delta is much higher than for less important "enhancement layer" flows. This is still a good technique; it's just more prone to nonuniform increases in visual quality. Thus it has a resultant rating of very good. Paper [11] implements the feedback formula in (8). However, this leads to problems in quality where the algorithm tries to ensure the transmission buffer capacity is constant. By ensuring the buffer capacity is constant, the quality of the frame may change suddenly if the buffer rapidly increases and decreases in size radically in an alternating fashion. This will leads to an oscillating digital waveform where the resulting video will have extremely high quality images

and then drop to low quality images. Because of this severe deficiency the algorithm scores a rating of bad.

Table 2 presents the final results for the papers adaptive techniques versus the critiqued areas. Most of the adaptive quality techniques are of above average quality overall.

	Critiqued Areas						
Technique	Scalability	Algorithmic Complexity	Smoothness				
Media-aware [1]	good	good	very good				
Cross-layer [9]	very bad	very bad	very bad				
Resource-Ex [10]	average	good	very good				
Multipath-Rt [11]	very good	very good	average				

Table 2. Adaptive quality

#### **3.3.** Cross Layer Information Gathering

The criteria used to compare the cross layer information gather techniques are estimation accuracy and ease of implementation.

Estimation accuracy involves the underlying assumptions and the accuracy of the actual estimation formula. Paper [8] gathers information from nodes and congestion. Equation (9) utilizes the busy, block, and idle times of the wireless card along with the payload size of the stream to calculate estimation on the rate for a given stream. This seems like a reasonable estimation due to the fact that it represents the amount of actual transmitted data over the entire time it took to transmit said data. The equations (10) and (11) deal with capacity of actual nodes. For (10), we assume that the estimated capacity is extremely accurate because the blocking time isn't normally large, and the busy time represents the time to send. This is a fairly accurate estimate on the capacity, as long as the packet sent is of an average size. Equation (12) provides a direct correlation between rate and congestion. This common observation is that an increase in rate leads to an increase in congestion. This is a very apt model for congestion increments. If the amount of expected bandwidth left is minuscule then the congestion blows up quickly for even small rate increments. However, if a plethora of bandwidth is left which results in a much smaller increase in congestion. The only possible problem in this formula is that the sum of the stream rates is larger than the expected bandwidth, which is impossible due to the lack of idle time being present in the estimated capacity. The validity of assumptions and accurate estimates results in a rating of very good. Paper [9] gathers information about low level transmission and buffer issues with respect to cross layer information gathering. They utilize the BER and the packet size in (13) to determine the probability of an error happening during transmission across a given link. This is an accurate assessment of an error happening on the link for a packet of size  $L_v$  bits. Due to the equation

calculating the probability of an error not occurring for one bit, and taking that to a power equal to the number of bits, this is not a problem. This gives the probability of an error not occurring for the transmitted data. The assumption they made was that the BER was normally and randomly distributed which is a fair assumption for normal network transactions. Using (13) they derive (14) which represents expected queuing delay on a given link. This formula sums all the queuing delays for all MSDU passing through the given link, which gives an approximate estimation of the queuing delay. However, this assumes that the delay per packet is accurately calculated. Given the amount of resources poured into devising this calculation in the paper, we assume it is valid without further verification. Therefore, this paper obtains a rating of very good for estimation accuracy also. Paper [10] gathers information on expected throughput and the BER. Equation (15) calculates the expected error rate given the physical layer mode. Again, this equation will have to be assumed correct due to the lack of resources for further verification. Equation (16) is much more feasible in its intentions. It calculates the expected valid throughput for the packet size, and multiplies this against the expected physical rate on the given physical mode. The assumptions that the error rate is accurate make this a good estimation of the available goodput. However, this is detracted by the fact that it does not subtract bandwidth for retransmissions caused by errors in the estimation. As for (17), this is a simple admission metric which accurately represents the ability to transmit a given packet. This assumes that goodput, service intervals, and the listening time given to the stream are known values. Due to the high number of assumptions and some flaws in the estimation we would assign this paper a rating of good.

Ease of implementation covers the ease of gathering this information and the availability of the actual information using modern systems. For the information that [8] gathers it is safe to say that not all of it is accessible by hardware. It is challenging to believe that the busy, idle, and wait time for each packet is accessible via a hardware interface. If the hardware does not support retrieval of this information, it has to be estimated from a custom device driver sitting on top. This device driver could perform a timing operation when called to send a message across the network interface, which can then be stored for future retrieval. The packet size is obviously easily calculated. Thus, the cross layer information for this paper is easily implemented. Due to the one major deficiency, the score obtained is only a rating good. Paper [9] has issues about where the majority of the cross layer information comes from. The modulation for the transmission at the physical layer is needed for many calculations. However, where and how the modulation is obtained is never discussed. This is a serious shortcoming for the algorithm. The error rate and throughput are easily

calculated and implemented. Thus, I would say there are not any issues with the rest of the low level gathered statistics. Therefore, this paper obtains a rating of good. The last information gathered is from [10]. Here they gather the BER from the physical layer mode. Again there is no mention of how the physical layer mode is determined as in [9]. If this mode is determined a priori then all nodes need to be set up in advance. Otherwise, the system needs to be able to identify the physical layer mode automatically and adjust accordingly. The MSDU size is easily calculated from known information at the application level. Other low level information like the length of the reservation time is easily transmitted across the overlay network and must already be known to actually send the video. Thus, we assign a rating of good also for the same reasons as the other sections.

Table 3 provides a quick overview of how cross layer information gathering faired for each paper. All papers essentially estimate accurate low level information or gather it directly. There do not appear to be any major problems with how the information was gathered or utilized.

	Critiqued Areas					
Technique	Estimation Accuracy	Ease of Implementation				
Media-aware [1]	very good	good				
Cross-layer [9]	very good	good				
Resource-Ex [10]	good	good				
Multipath-Rt [11]	N/A	N/A				

Table 3. Cross layer information gathering

# 4. Possible Future Research

There are a variety of areas for potential research for streaming media over wireless mesh networks, due to the new and their emerging nature and the problems they present.

#### 4.1. Dynamically Adapted Path Determination

From the path determination techniques presented in this paper it is clear that dynamically adapted routing algorithms need to be researched for wireless mesh networks. The majority of the presented algorithms had problems when faced with working on topologies that changed.

#### 4.2. Congestion Avoidance

Congestion in a wireless mesh network leads to dropped packets, reduced throughput, and dropped video streams. To avoid this cross layer techniques need to implement congestion avoidance. Once implemented, this will help increase the total network utility and save on wasted transmissions.

# 4.3. Adaptive Quality to Algorithms

Adaptive quality is requisite to meet the demands of the changing network. Without adjustments to quality, video transmissions may be jittery or even dropped. All that is needed is a dynamic adjustment of the video quality with respect to the state of the network. Several algorithms presented in this paper are good starting points for continuing adaptive quality algorithms.

## 4.4. Network Overlays

Several algorithms utilized network overlays in this paper. However, the overhead incurred by the network overlays is extreme in some cases. Optimizations need to be made to existing algorithms, or overlay networks need to be redesigned to reduce the overhead incurred by overlay networks.

## 4.5. Information Reuse

Cross-Layer optimized methods were the bulk of the presentation in this paper. A majority of the papers here only used the information gathered for one very specific purpose and not as many ways as could possibly be done.

#### 4.6. Expansion of Described Algorithms

The current algorithms presented in these papers could be improved in a variety of ways such as congestion avoidance. A majority of the papers had the capability for congestion avoidance and simply not yet implemented it.

# 5. Conclusions

This paper compared five papers which streamed video across wireless mesh networks. All papers were into sections dealing decomposed with path determination, adaptive quality, and cross laver information gathering. Commonly an overlay network was used to ensure transmission quality, and route information to nodes in the network. Path determination was lead by the newer congestion-minimized, localized, and distributed collaborative algorithms on average. Although, these techniques are harder to implement and have a higher overhead they are more flexible in the long run. The congestion-minimized routing approach is a good example which is more flexible than UDP/IP, especially in areas of congestion avoidance. Adaptive quality — which is important for the fact that quality must be dynamically adjusted during transmission to deal with the state of the network - was planned into the algorithms in some way. Wireless mesh networks are new and emerging field of study, but this entails that strong research needs to be focused on this area. If this is done it will allow for new robust algorithms for transmitting large amounts of data wirelessly independent of any centralized control.

# 6. References

- [1] Akyildiz, X. Wang, and W. Wang, "Wireless mesh networks: a survey," Computer Networks, vol. 47, pp. 445–487, 2005.
- [2] D. Li, Q. Zhang, C.N. Chuah, and S.J. Ben Yoo, "Multi-Source Multi-Path Video Streaming over Wireless Mesh Networks," IEEE International Symposium on Circuits and Systems (ISCAS), May 2006.
- [3] H.Seferoglu and A.Markopoulou, "Opportunistic Network Coding for Video Streaming over Wireless," the 17th Packet Video Workshop, November 2007.
- [4] P.S. Mogre, M.Hollick, and R. Steinmetz, "QoS in wireless mesh networks: challenges, pitfalls, and roadmap to its realization," 17th International workshop on Network and operating systems support for digital audio & video, June 2007.
- [5] V. Kone, S. Das, B. Y. Zhao, and H. Zheng, "QUORUM– Quality of Service Routing in Wireless Mesh Networks," IEEE/ICST International Conference on Heterogeneous Networking for Quality, Reliability, Security and Robustness (QShine), August 2007.
- [6] V. Navda, A. Kashyap, S. Ganguly, and R. Izmailov, "Real-time video stream aggregation in wireless mesh network," IEEE 17th International Symposium on Personal, Indoor and Mobile Radio Communications, pp. 1–7, September 2006.
- [7] V. Chavoutier, D.Maniezzo, C.E. Palazzi, and M. Gerla, "Multimedia over wireless mesh networks: results from a real testbed evaluation," The Sixth Annual Mediterranean Ad Hoc Networking WorkShop, pp. 56–62, Corfu, Greece, June 12–15, 2007.
- [8] X.Q. Zhu and B. Girod, "Media-Aware Multi-User Rate Allocation over Wireless Mesh Network," IEEE First Workshop on Operator-Assisted (Wireless Mesh) Community Networks, pp. 1–8, September 2006.
- [9] Y. Andreopoulos, N. Mastronarde, and M. Van Der Schaar, "Cross-layer optimized video streaming over wireless multihop mesh networks," IEEE, vol. 24 pp. 2104–2115, November 2006.
- [10] N. Mastronarde, D.S. Turaga, and M. Van Der Schaar, "Collaborative resource exchanges for peer-to-peer video streaming over wireless mesh networks," IEEE, vol. 25, pp. 108–118, January 2007.
- [11] F. Licandro, A. Lombardo, and G. Schembra, "Applying multipath routing to a video surveillance system deployed over a wireless mesh network," the 2nd ACM International Workshop on Wireless Multimedia Networking and Performance Modeling (Terromolinos, Spain, October 06–06, 2006), WMuNeP'06, ACM, New York, NY, pp. 19–26, 2006.



# A Novel Adaptive Hybrid Error Correction Scheme for Wireless DVB Services

# **Guoping TAN<sup>1</sup>**, Thorsten HERFET<sup>2</sup>

<sup>1</sup> Student Member, IEEE, <sup>2</sup> Senior Member, IEEE FR 6.2 – Telecommunications Lab, Saarland University Campus Building C6 3, Floor 10, 66123 Saarbrücken, Germany E-mail: {tan, herfet}@nt.uni-saarland.de

## Abstract

Real-time applications usually not only have a certain Packet Loss Ratio (PLR) requirement but also can have strict delay constraints. In the past, we proposed a Hybrid Error Correction (HEC) scheme with Packet Repetition (PR) technique for guaranteeing a certain PLR requirement under strict delay constraints. Unfortunately, the HEC-PR scheme can only work efficiently in multicast scenarios with small group size and small link PLR. Our further studies show that better performance can be obtained by combining the HEC-PR scheme with other traditional HEC schemes such as Type I HARQ and Type II HARQ techniques. Based on this idea, in this paper, a novel Adaptive HEC (AHEC) scheme combining the HEC-PR scheme with Type I and Type II HARQ techniques is proposed to satisfy a certain PLR requirement for delay bounded multicast services. Furthermore, the performance of the AHEC scheme is optimized by choosing the scheme with the least needed redundancy information automatically among the three HEC schemes. Finally, by applying the AHEC scheme in a typical wireless DVB scenario, we analyze the performances of the AHEC scheme and compare it with the HEC-PR scheme and an Adaptive Forward Error Correction (AFEC) scheme. The results show that the proposed AHEC scheme outperforms both the AFEC scheme and the HEC-PR scheme.

Keywords: ARQ, Forward Error Correction (FEC), Hybrid ARQ (HARQ), Multicast, Wireless Networks

# 1. Introduction

With the rapid development of broadband wireless networks, more and more attention has turned to distributing real time multimedia services over wireless networks. Many classes of mobile commerce applications require or can benefit from real-time multicast support in wireless networks: mobile auction or reverse auction, mobile entertainment services and interactive games, mobile distance educations etc. [1]. As a major example, most of our personal digital assistants (PDAs) and laptops are factory-equipped with a Wi-Fi interface. In recent years, more and more places are covered by wireless LANs with the IEEE 802.11 [2] family of protocols in hotspots like hotels, airports or conference locations. This will allow travelers to use their PDAs or laptops for watching television broadcastings, enjoying games or participating in video

conferences etc. All these new real-time multicast applications are very likely to appear soon with upcoming WiMAX or DVB-H [3] enabled devices. In the following, to show the packet loss issue in wireless real-time multicast systems, we will take the practical Digital Video Broadcasting (DVB) services over wireless LANs as the example for illustration.

We know that the IEEE 802.11 has been expected to be used for DVB services over home and nomadic networks. Moreover, since IP multicast provides a scalable and efficient means for distributing datagram to a group of receivers [4], IP-based networks were proposed for delivering DVB services [5]. The DVB systems based on IP multicast typically employ an application-level protocol to provide some information about the set of receivers and reception quality statistics. The Real-time Transport Protocol (RTP) [6] is usually used for this purpose. RTP or MAC layer of the IEEE 802.11 does not, however, guarantee any Quality of Service (QoS) for real-time multicast applications, although the amount of lost packets varies during the day and depends on the multicast data rate [7]. Therefore, it is essential to employ some error control techniques at application layer to guarantee a certain Packet Loss Ratio (PLR) requirement (e.g.  $10^{-6}$ , refer to [5]) needed for DVB services. In this paper, we thus address the packet loss issue in wireless real-time multicast systems.

Traditionally, the packet loss issue can be treated as erasure errors. As it is well known, there are mainly two categories of erasure error control techniques: Automatic Repeat Request (ARQ) that retransmits the lost packets and packet-level Forward Error Correction (FEC) that transmits redundant packets. Recently, many researchers have been studying how to efficiently improve the reliability of real-time multimedia multicast over wireless networks by using these two techniques. Some approaches focus on only one of the two schemes, ARQ alone [8] or FEC alone [7]. Many other approaches consider the combination of both to improve the performance (see, e.g., [9-13]). The integrated FEC/ARQ schemes are referred to as Hybrid Error Correction (HEC) schemes in this paper. The studies indicate that HEC schemes are much more efficient for recovering data packets than the schemes with either FEC or ARQ alone. In HEC schemes, many authors employ powerful FEC erasure coding techniques (e.g. Graph Codes [9] or Reed-Solomon (RS) codes [10,11,12]). In addition, different retransmission-based schemes for error control in multicast protocols geared toward real-time multimedia applications are analyzed in [14]. It is found that retransmission schemes are appropriate for such applications, and actually can be quite effective. In fact, the studies have shown that the retransmission based error control schemes for point-topoint communication or single receiver in multicast scenario can outperform all the existing point-to-point schemes [15]. Therefore, using retransmission based error control mechanism with a Packet Repetition (PR) technique; we developed an HEC-PR scheme for satisfying the target PLR requirement under strict delay constraints and optimized its performance in [13] for DVB distribution in home networks. Even though the scheme works perfectly in the in-home scenario and additionally has the merit of being backward compatible (so that conventional receivers with input buffers can benefit from this scheme without modifications), it's not fully scalable for applications with larger group sizes. To overcome its shortage, we proposed a Type I HARQ scheme in [12] for those multicast scenarios with large number of receivers.

However, the previous works mentioned above only show that the good performance can be obtained by combining the HEC-PR scheme and the Type I HARQ scheme, while the question on how to combine them has been left unanswered. In this paper, we thus try to answer this question by proposing an adaptive HEC (AHEC) scheme combining the HEC-PR scheme with traditional Type I and Type II HARQ scheme. Following the idea, we focus on developing one framework for the AHEC scheme and then optimizing its performance. In this paper, our main contributions include: (i) A novel Adaptive HEC scheme combing the HEC-PR scheme with other two traditionally HARQ schemes is proposed. This scheme is suitable for any delay bounded multicast application. (ii) By building a general mathematical framework for the AHEC scheme under strict delay constraints, we optimize its performance by minimizing the total needed Redundancy Information (RI). To the best of our knowledge, no general frameworks combining those HEC schemes have been proposed before for optimizing the performance of those schemes under strict delay constraints.

The rest of the paper is organized as follows. In Section 2, the performance of the packet level FEC scheme is introduced. In Section 3, we present a general mathematical framework on our novel AHEC scheme and introduce a method to optimize its performances. Applying the AHEC scheme in a typical DVB scenario over wireless LANs, we analyze its performances and compare it with the HEC-PR scheme and an Adaptive FEC scheme in Section 4. Finally, conclusions are given in Section 5.

Notation: Throughout of the paper, E(X) denotes the expected value of a random variable X; and we keep in mind that  $\binom{m}{h}$  is the number of ways h objects can be

chosen from among m objects without repetition.

# 2. Performance of Packet Level FEC

In this paper, it is assumed that a perfect forward erasure error correcting code (e.g. RS code) is used for the AFEC scheme and the AHEC scheme. More recent codes like LDPC [16] or Fountain / Raptor-codes [17] do scale well for long blocks and offer advantages concerning computational efficiency; in this paper, however, a perfect erasure code is taken as the "upper anchor", and the block sizes for the used application scenario (DVB over WLAN) can well be solved by RS codes with acceptable computational complexity. For the convenience of description, the perfect FEC code is denoted by (n, k) code here, where k denotes the number of data symbols per codeword and n denotes the code word length. Figure 1 shows the structure of the coding block transmitted within packets protected by the ideal (*n*, *k*) code.

As shown in Figure 1, the source data packet stream is divided into blocks each consisting of k consecutive data packets with a length of l bytes. The (n, k) code is applied to each row containing k data packets in order to produce a group of (n-k) parity packets. Without loss of generality, it is assumed in this paper that the symbol

size is always one byte for the FEC code. The coding block is transmitted by packets in the form of columns. Assuming exactly one packet per column, the receiver only needs to correctly receive any k of these n columns to be able to recover all the k data packets. Therefore, the PLR performance of this scheme is exactly the same as the performance of the ideal (n, k) code.



Figure 1. Applying ideal (n, k) code at the packet level, which forms an FEC coding block in n packets

To simplify the analysis, in this paper, the packet loss channel in wireless networks is modeled as the erasure error channel with independently and identically distributed (*i.i.d.*) losses with uniform distribution. Firstly, we define the probability of *b* packets lost in a sequence of *n* packets in the erasure channel with link PLR of  $P_e$  as  $P(b,n,P_e)$ . Since all of the *n* packets have the same loss probability of  $P_e$ , the probability  $P(b,n,P_e)$ is given by:

$$P(b,n,P_e) = \binom{n}{b} (1-P_e)^{n-b} (P_e)^b$$
(1)

In the following, let the random variable  $I_k$  represent the number of data packets lost in a group of k data packets after decoding using the (n,k) code. Upon the definition of  $I_k$ , the PLR performance of the (n,k) code actually can be computed by  $E(I_k)/k$ . That means we only need to calculate the expected value of  $I_k$ . To obtain  $E(I_k)$ , we firstly have to find out the Probability Distribution Function (PDF) of  $I_k$ . For the convenience of description, here we assume that the value of  $I_k$  is *i* and there are b packets lost in a group of n packets. If b is not more than *n*-*k*, the number of packets received in a group of n packets will be at least k so that all of the kdata packets can be recovered. Obviously, the value of  $I_k$ is zero in this case. On the other hand, if the value of  $I_k$ is more than zero, there are exactly *i* (where  $1 \le i \le k$ ) data packets lost in the group of *n* packets after decoding with the (n,k) code. It indicates that at least max (n-k+1,i)and at most (n-k+i) packets are lost in this group. That is, the value of b is in the range of  $[\max(n-k+1,i), n-k+i]$ . Let  $P_d(i,b)$  denote the probability of *i* data packets lost under the condition of b packets lost in a group of npackets. In other words: Among all of the b packets lost in the group, there are i data packets lost among all of the k data packets and b-i parity packets lost among all of the *n*-*k* parity packets. Let  $P_d(i,b)$  denote the

probability of i data packets lost among all of those b packets lost. Note that all packets in a group of n packets have the same loss probability in the *i.i.d* erasure error channel, we thus have:

$$P_{d}(i,b) = \frac{\binom{k}{i}\binom{n-k}{b-i}}{\binom{n}{b}}$$
(2)

Based on the analysis above, using (1) and (2), we then obtain the PDF of  $I_k$  as follows:

$$\Pr(\mathbf{I}_{k}=i) = \sum_{b=\max(n-k+1,i)}^{n-k+i} P(b,n,P_{e}) P_{d}(i,b), i = 1,2,...,k.$$
(3)

Following (3), the expected value of  $I_k$  thus can be calculated by:

$$E(\mathbf{I}_{k}) = \sum_{i=1}^{k} \sum_{b=\max(n-k+1,i)}^{n-k+i} i \times P(b,n,P_{e}) P_{d}(i,b)$$
(4)

Finally, when the ideal (n,k) code is applied in the erasure error channel with link PLR of  $P_e$ , the PLR achieved will be:

$$PLR_{(n,k)} = \frac{E(\mathbf{I}_k)}{k}$$

$$= \frac{1}{k} \sum_{i=1}^k \sum_{b=\max(n-k+1,i)}^{n-k+i} i \times P(b,n,P_e) P_d(i,b)$$
(5)

From the analysis above it follows that (5) is also the PLR performance of the packet level FEC scheme with an ideal (n,k) code over the erasure error channel with link PLR of  $P_e$ .

#### 3. Proposed AHEC Scheme

In this section, first, we introduce the system model of the proposed Adaptive Hybrid Error Correction (AHEC) scheme combing the HEC-PR scheme and the traditional Type I and Type II HARQ schemes. Then, we present a mathematical framework for the AHEC scheme. Based on the mathematical framework, we finally present how to design the optimum parameters for the AHEC scheme guaranteeing a certain PLR requirement under strict delay constraints.

At the beginning, for the AHEC scheme using retransmission technique, rather than focus on a particular transport protocol, we shall consider a generic retransmission based scheme with the following features:

- A selective repetition, NACK-only retransmission scheme is used;
- The transmitter multicasts the required packets immediately to all receivers upon getting a NACK.

In addition, to simplify the analysis we make the

189

following assumptions for the retransmission based schemes:

- The feedback channel for NACKs is assumed to be error-free. Since NACKs are control messages and real systems usually provide mechanisms to guarantee the reliable transmission of control signals, this assumption is realistic in many cases. Or alternately, the effect of NACKs loss can be overcome by setting a margin for the PLR performance of the AHEC scheme.
- All of the receivers experience erasure error channel with *i.i.d* of uniform distribution. This means we do not consider the effect of temporal correlation of the channel and spatial correlation among different receivers in this paper; this, however, is actually ongoing work, which will come soon in [18].

Now, the essential symbols are defined and summed up in Table 1.

Symbol	Definition
PLR <sub>target</sub>	target PLR requirement
D <sub>target</sub>	target Delay requirement
$P_{e}(j)$	the PLR for the <i>j</i> -th receiver
RTT(j)	average round trip time for the <i>j</i> -th receiver,
	one way delay is $RTT(j)/2$
N <sub>recv</sub>	number of receivers in the multicast scenario <sup>1</sup>
t <sub>s</sub>	the average interval between two continuous original data packets at the transmitter <sup>2</sup>
t <sub>rw</sub>	the average waiting time at each receiver, which is the time between the detection of a packet loss and the time when the corresponding NACK is sent <sup>3</sup>
t <sub>sw</sub>	the average waiting time at the transmitter, which is the time between receiving a NACK message and the time when the corresponding packets required by the NACK message are retransmitted
$t_{lp}(j)$	the duration from the time the <i>j</i> -th receiver detected packets lost to the time it possibly receives the required packets, which is $RTT(j) + t_{sw} + t_{rw}$

**Table 1. Symbols definitions** 

### 3.1. System Model

From [13] we know that the HEC-PR scheme has a major drawback: the total needed RI raises with the increase of the group size linearly in a multicast scenario, because the receivers can not share common retransmission packets for repairing different missing data packets. For overcoming this shortage, we proposed a Type I HARQ scheme in [12] for those multicast

scenarios with large number of receivers. However, the works mentioned above only show that good performance can be obtained by combining the HEC-PR scheme and the Type I HARQ scheme, while the question on how to combine them has been left unanswered. In this paper, we thus propose an adaptive HEC (AHEC) scheme combining the HEC-PR scheme with traditional Type I and Type II HARQ scheme. The system model of the AHEC scheme is shown in Figure 2.

As shown in this figure, the transmitter firstly transmits encoding blocks to all receivers using the packet level FEC code. Here it is assumed that perfect forward erasure error correction code (e.g. Reed-Solomon code.) is used and the number of source data packets is k in one encoding block. That is, upon received any k packets of one encoding block, the receiver can recover all the data packets. Otherwise, the receiver will send Negative-Acknowledgments (NACKs) to the transmitter for repairing the missing data packets. Now we explain the AHEC scheme in more detail:

- 1) First, the sender sends a certain amount of redundant packets or only k data packets to all of the receivers immediately with the first transmission. Especially, when k is set to one then the redundant packets during all of the retransmissions are always multiple copies of source data packets; this scheme acts actually as the HEC-PR scheme proposed in [13].
- 2) If any k packets of one encoding block are received, the receiver then can recover all of the k data packets and forward them to the application immediately. Otherwise, the receiver will transmit a NACK message to the sender to require essential redundant packets for recovering all of the missing data packets.
- 3) Upon getting the first NACK message for one encoding block during each retransmission round, the sender will multicast a certain number of redundant packets to all of receivers immediately with one copy (or multiple copies) of these retransmission packets; afterwards, if other NACKs for the same encoding block are received, the sender will decide if multicast more redundant packets to all of the receivers according to the requirements of NACKs. That is, if those later NACKs require more redundant packets than the fist NACK message, the sender will multicast further redundant packets to all of receivers immediately; otherwise, the sender will neglect those NACKs. Similarly, all of receivers can do suppression of NACKs by this rule if those NACKs are transmitted by multicasting mode.

From above introduction, we know that the performance of the AHEC scheme mainly depends on three parameters: the number of retransmission rounds; the number of redundant packets with the first transmission and retransmissions and the number of copies of redundant packets with retransmissions. Note that if the redundant packets are parity packets, they

<sup>&</sup>lt;sup>1</sup>The parameter  $N_{recv}$  is also viewed as the group size in a multicast scenario in this paper.

<sup>&</sup>lt;sup>2</sup> In this paper, it is assumed that the interval is same to the retransmission interval for different copies of retransmission packets. 3 T

<sup>&</sup>lt;sup>3</sup> The average waiting time at each receiver is identical due to the same process for all of the receivers.

should not be repeated, but possibly more new parity packets should be retransmitted. The remaining task is to find out those suitable parameters of the AHEC scheme for satisfying the strict QoS requirements of real-time services. In the following section, we will present a mathematical framework for analyzing the performances of the AHEC scheme.

#### 3.2. Performance Analysis

Theoretically, we should also design parameters for the AHEC scheme with each receiver separately as for HEC-PR scheme as in [13]. However, it is very hard to implement for practical systems if different FEC codes used for different receiver with the first transmission. To simplify the implementation, therefore, the AHEC scheme will adopt the same parameters for every receiver. Since the assumed erasure code is perfect, a suitable choice of the code rate can guarantee that all receivers with their different channel conditions can be served, so this simplification doesn't negatively influence the overhead.



Figure 2. System model of the AHEC scheme

The parameters for the AHEC scheme with retransmissions are defined as follows:

- *k*: the number of source data packets in one encoding block;
- N<sub>p</sub>: the number of redundant packets in one encoding block with the first transmission;
- *N<sub>cc</sub>*: a constant coefficient, which is the number of additional new redundant packets for one encoding block with different retransmission rounds;
- N<sub>blk</sub>: the number of packets in one block with the first transmission, which is k+N<sub>p</sub>;
- *N<sub>rr,max</sub>*: the maximum possible number of retransmission rounds;
- $N_{rt}^q$ : the number of copies for each retransmission packet at the sender during the *q*-th retransmission round, where  $1 \le q \le N_{rr \max}$ ;
- *N<sub>rt,max</sub>*: the maximum possible number of copies for each retransmission packet at the sender, which is:

$$N_{rt,\max} = \sum_{q=1}^{N_{rr,\max}} N_{rt}^q$$

For the convenience of description, two additional random variables are defined as follows:

• *I<sub>k</sub>(j,w)*: a random variable representing the number of missing data packets for the *j*-th receiver in one encoding block of *k* source data packets after

experiencing *w* retransmission rounds, where  $1 \le w \le N_{rr,max}$ ;

•  $N_{req}(j)$ : a random variable representing the number of redundant packets required to receive for recovering all of the *k* data packets in one block for the *j*-th receiver in the first retransmission round, where  $0 \le N_{rea}(j) \le k$ .

Based on above definitions, we now begin to analyze the PLR performance for one receiver (without loss of generality, it is assumed to be the *j*-th receiver) with the AHEC scheme. To derive the PLR performance of the AHEC scheme for the *j*-th receiver, we need to calculate the expected value of the number of missing data packets in one encoding block of k source data packets after the retransmission packets experiencing w  $(0 \leq w \leq N_{rr,max})$ retransmission rounds. The PLR performance of the AHEC scheme for the *j*-th receiver then can be calculated as:  $E(I_k(j,w))/k$ , which is the final PLR at the *j*-th receiver after all of the retransmission packets experienced w retransmission rounds. Note that if the w is set to zero, the AHEC scheme acts as the AFEC scheme. In the following, it is always assumed that the *w* is more than zero for the AHEC scheme.

First of all, in order to recover all of the missing data packets for each receiver that received fewer than k packets for one block, at least  $N_{req,max}$  redundant packets need to be retransmitted at the sender:

$$N_{req,max} = max(N_{req}(1), N_{req}(2), ..., N_{req}(N_{recv}))$$
 (6)

Obviously,  $N_{req,max}$  is also a random variable. Since it is assumed that the feedback channel is error-free, the random variable  $N_{req,max}$  always reflects the true maximum number of lost packets in one block for the worst receiver. Before calculating the average number of lost data packets in one block for the *j*-th receiver, here we define two useful probabilities: one is the PDF of  $N_{req,max}$  (i.e.  $Pr(N_{req,max}=i)$ ), which is denoted by  $P^{i}_{Nreq,max}$ ; the other is the probability of  $N_{req,max}$  of *i* in the condition of  $N_{req}$  (*j*) of *c* (i.e.  $Pr(N_{req,max}=i|N_{req}(j)=c)$ ), which is denoted by  $P_{req}(i,c,j)$ . The detail derivation on these two important probabilities is attached in the appendix.

Secondly, according to the definitions above, the parameter  $N_{cc}$  is always constant for any random value of  $N_{req,max}$  in each retransmission round. Now let symbol "r" denote the number of received redundant packets within w retransmission rounds. Note that multiple copies of a retransmission packet received are counted as one redundant packets received. We then define the probability of r different redundant packets received after all of the  $m=N_{req,max}+N_{cc}$  different redundant packets experiencing w retransmission rounds for the *j*-th receiver as  $P_{recv}(r,m,w,P_e(j))$ . Note that the loss probably of each retransmission packet within w

retransmission rounds will be  $(P_e(j))^{\sum_{q=1}^{n} N_n^q}$ . We thus have:

$$P_{recv}(r,m,w,P_{e}(j)) = \binom{m}{r} \left(1 - (P_{e}(j))^{\frac{w}{q-1}}\right)^{r} \left((P_{e}(j))^{\frac{w}{q-1}}\right)^{m-r}$$
(7)

For the convenience of description, in the following, we define some temp symbols as follows: symbol 'i' denotes the number of data packets lost in a group of  $N_{blk}$  packets for the *j*-th receiver in the first transmission; symbol 'b' denotes the total number of packets lost in a group of  $N_{blk}$  packets for the *j*-th receiver in the first transmission; symbol 's' denotes the total number of packets lost in a group of  $N_{blk}$  packets for the *j*-th receiver in the first transmission; symbol 's' denotes the number of redundant packets sent at the sender in the first retransmission round; symbol 'm' denotes the number of redundant packets received during the retransmission round for the *j*-th receiver. Now, as introduced in Section 2, we also adopt  $P_d(i,b)$  to denote the probability of *i* data packets lost under the condition of *b* packets lost in a group of *n* packets. The conditional probability  $P_d(i,b)$  also can be calculated by (2).

Finally, we can derive the PDF of  $I_k(j,w)$  based on those probabilities introduced above. Now we assume that the value of  $I_k(j,w)$  is *i* (where  $1 \le i \le k$ ), which means that there are i data packets lost after experiencing wretransmission rounds. Obviously, it indicates that there are b (where  $\max(N_p+1, i) \le b \le N_p+i$ ) packets lost in the block of  $N_{blk}$  packets for the *j*-th receiver in the first transmission. According to the AHEC scheme, this receiver will require  $b-N_p$  redundant packets for retransmission at the sender for recovering the missing *i* data packets. However, at the same time, the sender will possibly send s (where  $b-N_p+N_{cc} \le s \le k+N_{cc}$ ) parity packets due to combining all of the NACKs from overall receivers in the multicast scenario. Finally, note that the receiver obtained  $k+N_p-b+m$  packets at the end of the w retransmission rounds. Note that the data packets lost only happen under the condition of  $k+N_p-b+m$  being less than k, which means the value of m will be less than  $b-N_p$ . Based on above analysis, using (2), (7) and the probability  $P_{req}(m,c,j)$  (See Appendix), the PDF of  $I_k(j,w)$ then can be expressed as this form:

$$Pr(I_{k}(j,w)=i) = \sum_{b=\max(N_{p}+l,i)}^{N_{p}+i} \sum_{s=b-N_{p}+N_{cc}}^{k+N_{cc}} P_{d}(i,b) \times P_{req}(s-N_{cc},b-N_{p},j) \left( \sum_{m=0}^{b-N_{p}-1} P_{recv}(m,s,w,P_{e}(j)) \right)$$
(8)

where *i*=1,2,...,*k*.

Following (8), we then obtain the expected value of  $I_k(j,w)$ :

$$E(\mathbf{I}_{k}(j,w)) = \sum_{i=1}^{k} i \times \Pr(\mathbf{I}_{k}(j,w) = i)$$
(9)

Relying on (9) and substituting (7) and (8) into (9), we then obtain the PLR performance of the AHEC scheme for the *j*-th receiver with  $N_{rr,max}$  retransmission rounds immediately:

$$PLR_{AHEC}(j, N_{rr, \max}) = \frac{E(I_{k}(j, N_{rr, \max}))}{k}$$

$$= \frac{1}{k} \sum_{i=1}^{k} \sum_{b=\max(N_{p}+1,i)}^{N_{p}+i} \sum_{s=b-N_{p}+N_{cc}}^{k+N_{cc}} i \times$$

$$P_{d}(i,b)P_{req}(s - N_{cc}, b - N_{p}, j) \left( \sum_{m=0}^{b-N_{p}-1} P_{recv}(m, s, N_{rr, \max}, P_{e}(j)) \right)$$

$$= \frac{1}{k} \sum_{i=1}^{k} \sum_{b=\max(N_{p}+1,i)}^{N_{p}+i} \sum_{s=b-N_{p}+N_{cc}}^{k+N_{cc}} X_{r}$$

$$P_{d}(i,b)P_{req}(s - N_{cc}, b - N_{p}, j) \times$$

$$\left( \sum_{m=0}^{b-N_{p}-1} \binom{s}{m} (1 - (P_{e}(j))^{\sum_{q=1}^{N_{rr}} N_{r}^{q}})^{m} ((P_{e}(j))^{\sum_{q=1}^{N_{rr}} N_{r}^{q}})^{s-m} \right)$$

$$= \frac{1}{k} \sum_{i=1}^{k} \sum_{b=\max(N_{p}+1,i)}^{N_{p}+i} \sum_{s=b-N_{p}+N_{cc}}^{k+N_{cc}} X_{r}$$

$$P_{d}(i,b)P_{req}(s - N_{cc}, b - N_{p}, j) \times$$

$$\left( \sum_{m=0}^{b-N_{p}-1} \binom{s}{m} (1 - (P_{e}(j))^{N_{rr,\max}})^{m} ((P_{e}(j))^{N_{rr,\max}})^{s-m} \right)$$
(10)

For simplifying the description, we define a vector as follows:  $\vec{P}_e^{j} = [P_e(1), P_e(2), ..., P_e(j-1), P_e(j+1), ..., P_e(N_{recv})]$ . By observing (10), we can find that it is actually a function with these parameters:  $k, N_p, N_{cc}, \vec{N}_{rt}, P_e(j)$  and  $\vec{P}_e^{j}$ , which is denoted by  $f_{PLR,AHEC}(k, N_p, N_{cc}, N_{rt,max}, N_{recv}, P_e(j), \vec{P}_e^{j})$  in this paper.

In the following, let's consider the total needed RI with the AHEC scheme, which includes two parts: one is the common part for all of the receivers in the first transmission, which is  $N_p/k$ ; the other is the part in the retransmissions, which is caused by the retransmissions of redundant packets for all of the receivers. For the convenience of calculation, we divide the second part into two subparts: one is the needed RI in the first retransmission round (denotes by *RI*<sub>AHEC-I</sub>); the other part is the needed RI in the retransmission rounds (denotes by RI<sub>AHEC-II</sub>) of from the second retransmission round to the  $N_{rr,max}$  retransmission round. First, considering the calculation of  $RI_{AHEC-I}$ , note that the value *i* (where  $0 \le i \le k$ ) of  $N_{req,max}$  means that  $N_{rt}^1 \times (i + N_{cc})$  redundant packets will be retransmitted in the first retransmission round at the sender. Using the PDF of  $N_{req,max}$  (See Appendix), thus, the  $RI_{AHEC-I}$  is given by:

$$RI_{AHEC-I} = \frac{1}{k} \sum_{i=1}^{k} N_{rt}^{1} \times (i + N_{cc}) \times P_{N_{req,max}}^{i}$$
(11)

Then, considering the calculation of  $RI_{AHEC-II}$ , to simplify the analysis, here it is assumed that there is only one Equivalent Receiver (ER) in the multicast scenario with  $N_{recv}$  receivers and the loss probability of each

Copyright © 2008 SciRes.

I. J. Communications, Network and System Sciences, 2008, 2, 105-206

retransmission packet for the ER is  $\overline{P}_e$ , where  $\overline{P}_e$  is the average link PLR and defined as:

$$\overline{P}_{e} = \frac{\sum_{j=1}^{N_{recv}} P_{e}(j)}{N_{recv}}$$
(12)

To derive  $RI_{AHEC-II}$ , let's note the following fact: if the ER requires i  $(1 \le i \le k)$  redundant packets for repairing missing data packets in the q-th  $(2 \le i \le N_{rr,max})$  retransmission round, it indicates that the ER required j  $(i \le j \le k)$  redundant packets in the first retransmission round and received only j-i redundant packets in the group of  $N_{cc}$ +j retransmission redundant packets in all of the previous q-1 retransmission rounds. Therefore, using the PDF of  $N_{req,max}$  and (7), we can obtain the probability of the ER requiring i redundant packets in the q-th retransmission round, i.e.:

 $\sum_{j=i}^{k} P_{N_{req,\max}}^{j} P_{recv}(j-i, N_{cc}+j, q-1, \overline{P}_{e})$ . Then the calculation

of *RI*<sub>AHEC-II</sub> can be written as this form:

$$RI_{AHEC-II} = \frac{1}{k} \sum_{q=2}^{N_{rr,max}} \sum_{i=1}^{k} N_{rt}^{q} \times (i + N_{cc}) \times \left( \sum_{j=i}^{k} P_{N_{req,max}}^{j} P_{recv}(j-i, N_{cc}+j, q-1, \overline{P}_{e}) \right)$$
(13)

As a result, by substituting (7) into (13) and then combining (11) and (13), the total needed RI for the AHEC scheme with  $N_{rr,max}$  retransmission rounds is given by:

$$RI_{AHEC} = \frac{N_p}{k} + RI_{AHEC-I} + RI_{AHEC-II}$$

$$= \frac{N_p}{k} + \frac{1}{k} \sum_{i=1}^k N_{rt}^1 \times (i + N_{cc}) \times P_{N_{req,max}}^i + \frac{1}{k} \sum_{q=2}^{N_{rr,max}} \sum_{i=1}^k N_{rt}^q \times (i + N_{cc}) \times \left( \sum_{j=i}^k P_{N_{req,max}}^j \binom{N_{cc} + j}{j - i} \left( (\overline{P_e})^{\sum_{g=1}^{q-1} N_{rt}^g} \right)^{N_{cc} + i} \left( 1 - (\overline{P_e})^{\sum_{g=1}^{q-1} N_{rt}^g} \right)^{j - i} \right)$$
(14)

Now we define two vectors  $\vec{P}_e = [P_e(1), P_e(2), ..., P_e(N_{recv})]$  and  $\vec{N}_{rt} = [N_{rt}^1, N_{rt}^2, ..., N_{rt}^{N_{rr,max}}]$ . By observing (14), we can also find that the RI performance of the AHEC scheme actually is a function with such these parameters: k,  $N_p$ ,  $N_{cc}$ ,  $\vec{N}_{rt}$ ,  $\vec{P}_e$  and  $\vec{P}_e$ , which is denoted by:

 $f_{RL,AHEC}(k, N_p, N_{cc}, \overline{N}_{rt}, N_{recv}, \overline{P}_e, \overline{P}_e)$  in this paper.

#### 3.3. Optimization of the AHEC Scheme

In this section, we will propose a method to design suitable parameters for the AHEC scheme. Note that all of the receivers share identical parameters for this scheme. Therefore, if the AHEC scheme can guarantee the QoS requirements for the worst receiver, it can also guarantee the same QoS requirements for every receiver in the multicasting scenario. Without loss of generality, it is assumed the first receiver is the one with the worst situation in a multicasting scenario. In other words, the first receiver has the largest link PLR and the largest RTT. Our remaining task is to design suitable parameters of the AHEC scheme, which will satisfy a certain PLR requirement for the first receiver under strict delay constraint with minimum total needed RI.

First of all, it is known that the delay requirement will limit the number of data packets in one block and the number of retransmissions. In the following, we will derive the boundary of the two parameters  $N_{rr,max}$  and k based on the strict delay constraint. For those retransmission packets in the first receiver, the maximum possible end-to-end delay includes four parts: the one-way delay in the first transmission (which is RTT(1)/2;  $N_{rr,max}$  of two-way delays in the retransmission rounds (which is  $N_{rr,max} \times t_{lp}(1)$ ); the decoding delay caused by the length of encoding block (which is  $(k \times t_s)$ ) and total number of intervals for the copies of retransmission packets (which is  $(N_{rt,max} \times t_s)$ ). Thus, for those retransmission packets of the first receiver, the maximum possible end-to-end delay must satisfy:

$$\frac{RTT(1)}{2} + N_{rr,\max} \times t_{lp}(1) + (k + N_{rt,\max}) \times t_s \le \mathbf{D}_{\text{target}}$$
(15)

Because the value of  $(k + N_{rt,max})$  is at least 2 for the AHEC scheme, the maximum allowable number of retransmission rounds is constrained by:

$$\hat{N}_{rr,\max} = \left| \frac{D_{\text{target}} - \frac{RTT(1)}{2} - 2 \times t_s}{t_{lp}(1)} \right|$$
(16)

where  $\lfloor x \rfloor$  denotes the largest integer not greater than *x*. Therefore, for the AHEC scheme the parameter  $N_{rr,max}$  will be limited in the range of between one and  $\hat{N}_{rr,max}$ . Then, we define the length of *k* with *w* retransmission rounds and maximum *v* copies of retransmission packets as k(w,v), which is given by:

$$k(w,v) = \left[\frac{D_{\text{target}} - w \times t_{lp}(1) - \frac{RTT(1)}{2} - v \times t_s}{t_s}\right]$$
(17)

where  $1 \le w \le \hat{N}_{rr, \max}$  and  $v \ge w$ .

Note that in (17) the parameter k will only rely on the parameters w and v if  $t_{lp}(1)$ ,  $t_s$ , RTT(1) and  $D_{target}$  are fixed. To simplify the design, we assume that the parameters  $t_{lp}(1)$ ,  $t_s$  and RTT(1) is always fixed in a scenario and  $D_{target}$  is also constant for the given QoS requirements. Therefore, the length of k will only depend on the parameters w and  $N_{rt,max}$ .

Considering practical implementation, actually, we can set an upper band of the maximum possible value for  $N_{rt,max}$  (denotes by  $\tilde{N}_{rt,max}$ , note that it is not less than  $\hat{N}_{rr,max}$  due to the practical consideration). Theoretically, the value of  $\tilde{N}_{rt,max}$  can be set to infinite. However, as shown in (14), the total needed RI increase with the increasing of  $\tilde{N}_{rt,max}$  significantly. Therefore, the minimum total needed RI usually can be acquired under a small value of  $\tilde{N}_{rt,max}$ . Now the parameter  $\bar{N}_{rt}$  can be limited in a small finite space  $\Phi^{N_{rr,max}}$  (i.e.  $\bar{N}_{rt} \in \Phi^{N_{rr,max}}$ ), which is:

$$\Phi^{N_{rr,\max}} = \left\{ \begin{bmatrix} r_1, r_2, \dots, r_{N_{rr,\max}} \end{bmatrix} \middle| \begin{array}{c} \sum_{k=1}^{N_{rr,\max}} r_k \leq \widetilde{N}_{rl,\max} \\ r_i \geq 1, N_{rr,\max} \geq i \geq 1 \end{bmatrix} \right\}$$

where:  $1 \le N_{rr, \max} \le \hat{N}_{rr, \max}$ 

Note that the parameter  $N_{rt,max}$  of the AHEC scheme is actually the sum of all of the elements in the vector  $\bar{N}_{rr}$ , which is denoted by sum $(\bar{N}_{rr})$  here. Depending on the PLR and RI performance analyzed for the AHEC scheme in Subsection 3.2, our optimization problem thus can be written as the following form:

$$\begin{split} RI_{AHEC,opt} &= \mathop{\arg\min}_{\bar{N}_{rt} \in \Phi^{N_{rr},\max}} \\ f_{RI,AHEC} \left( k(N_{rr,\max}, \operatorname{sum}(\bar{N}_{rt})), N_{p}, N_{cc}, \bar{N}_{rt}, N_{recv}, \overline{P}_{e}, \bar{P}_{e} \right) \end{split}$$

Subject to:

$$1 \le N_{rr,\max} \le \hat{N}_{rr,\max}$$
$$f_{PLR,AHEC} \left( k(N_{rr,\max}, \operatorname{sum}(\bar{N}_{rt})), N_p, N_{cc}, \operatorname{sum}(\bar{N}_{rt}), N_{recv}, P_e(1), \bar{P}_e^1 \right)$$
$$\le \operatorname{PLR}_{\operatorname{target}}$$

(19)

(18)

By solving (19) with traversing the full space, we thus can obtain the optimal parameters for the AHEC

scheme:  $k, N_p, N_{cc}, N_{rr,max}$  and  $\bar{N}_{rt}$ .

Remarks: If k is set to one, the AHEC scheme acts as a pure HEC-PR scheme; If k is set to more than one and  $N_p$  is set to more than zero, the AHEC scheme acts as the traditional Type I HARQ scheme; If k is set to more than one and  $N_p$  is set to zero, the AHEC scheme acts as the traditional Type II HARQ scheme. As a result, the AHEC scheme can choose the best scheme automatically among the HEC-PR scheme, traditional Type I and Type II HARQ scheme by solving (19).

#### 4. Analysis Results

In this section, we firstly analyze the performances of three schemes over an erasure error channel: the HEC-PR scheme, the AFEC scheme and the AHEC scheme, respectively; and then compare them with each other. Then, we study the effect of the parameter k in the AHEC scheme. For the convenience of comparing different schemes fairly, we make some assumptions as follows: the entire receivers experience the *i.i.d* erasure channel with the same level of original link PLR; the three schemes use identical system parameters with the same QoS requirements for the PLR and the same latency. In this case, we consider DVB services over wireless home networks with a group size of less than 7, RTT of less than 15ms and a wireless link PLR of up to 10% when the video multicast data rate is more than 500Kbps [7]. The target PLR requirement is set to  $10^{-6}$ under the strict delay constraint of 100ms (refers to [5]). However, it should be clear that the proposed AHEC scheme is suitable for any wireless multicasting scenario under strict delay boundary. Now we apply the three schemes in a typical scenario with the common system parameters, which are summarized in Table 2.

Table 2. System parameters

PLR Requirement: PLR <sub>target</sub>	10-6
Latency Constraint: D <sub>target</sub>	100ms
Multimedia Data Rate:	4Mbps
Packet Loss Model	GE Model
RTT	15ms
$t_{rw}+t_{sw}$ :	0ms
Encoding Packet Length: N <sub>symb</sub>	1250bytes
Original Average Link PLR: P <sub>e</sub>	$10^{-3} \sim 10^{-1}$

Actually, the following theoretical analysis results have been accompanied by simulations with ns-2 [19]. However, those simulation results are not further explained in this paper due to the fact theory matches simulation very well.

#### 4.1. Optimization Results

In the following, we will design the optimum parameters for the AHEC scheme for this typical scenario and then

Copyright © 2008 SciRes.

compare it with the HEC-PR scheme and the AFEC scheme. First, to meet the latency requirement we can obtain the maximum allowable number of retransmission rounds  $\hat{N}_{rr,max}$  with the AHEC scheme by solving (16), which are at most five. Then, since it is found that the optimum value for the parameter  $N_{rt,max}$  is usually much less than five by our numerical calculations for this case, the maximum value of  $N_{rt,max}$  is set to five (i.e.  $\tilde{N}_{rt,max} = 5$ ) to make sure that the optimum results of the AHEC scheme can be acquired by the searching algorithm. Note that in the traditional Type I and Type II HARQ schemes, multiple copies of a parity packet only can be counted as one redundant packet at the receivers. For the consideration of efficiency, we should only adopt new parity packets instead of multiple copies of redundant packets at each retransmission round. For this reason, the parameters of the AHEC scheme will always satisfy  $N_{rt,max} = N_{rt,max}$  (i.e.  $N_{rt}^q \equiv 1, 1 \le q \le N_{rr,max}$ ) in the case of k>1. Based on these boundaries, the optimum parameters of AHEC with different link PLR and small group size can be obtained by solving (19). Parts of the optimum parameters are shown in Table 3 and 4.

Table 3. Optimum parameters of the AHEC scheme with  $$N_{recv}{=}2$$ 

Average	Optimum Parameters of the AHEC scheme						
Link PLR	k	N.,	N.,		Ñ	rt rt	
	r	1 <b>•</b> p	1 * cc	$N_{rt}^1$	$N_{rt}^2$	$N_{rt}^3$	$N_{rt}^4$
0.001	23	0	0	1	1	-	-
0.01	16	0	0	1	1	1	-
0.03	9	0	0	1	1	1	1
0.05	9	0	0	1	1	1	1
0.07	1	0	0	1	1	1	2
0.09	1	0	0	1	1	1	2
0.10	16	2	1	1	1	1	-

 Table 4. Optimum parameters of the AHEC scheme with average link PLR of 0.06

N <sub>recv</sub>	Optimum Parameters of the AHEC scheme						
	k	Nn	N <sub>ee</sub>		$\vec{N}$	rt	
		110	1,00	$N_{rt}^1$	$N_{rt}^2$	$N_{rt}^3$	$N_{rt}^4$
1	1	0	0	1	1	1	1
2	1	0	0	1	1	1	1
3	23	2	1	1	1	-	-
4	23	2	1	1	1	-	-
5	23	3	1	1	1	-	-
6	23	3	1	1	1	-	-
7	23	3	1	1	1	-	-

From these two tables, we can see that the AHEC scheme can automatically choose the most suitable scheme according to current group size and average link PLR among the HEC-PR scheme, the Type I HARQ scheme and the Type II HARQ scheme. For example, as shown in Table 3, the AHEC scheme will act as the

Type II HARQ scheme if the real-time multicast scenario is with  $N_{recv}$ =2 and average link PLR of 0.05. Similarly, as shown in Table 4, the AHEC scheme acts as the HEC-PR scheme if the scenario is with  $N_{recv}$ =2 and average link PLR of 0.06; and it will act as the Type I HARQ scheme if the scenario is with  $N_{recv}$ =7 and average link PLR of 0.06.



Figure 3. The total needed RI with different schemes

#### 4.2. Performance Comparisons

Upon those optimum results shown in Table 3 and Table 4, now we obtain the total needed RI by (14) of the AHEC scheme with different group size, which are shown in Figure 3.

For comparing the performances among different schemes, Figure 3 also show the total needed RI with the HEC-PR scheme and the AFEC scheme. Note that the AFEC scheme is actually a special case of the AHEC scheme with  $N_p=0$  and  $N_{rr,max}=0$ . From this figure, we can see that the total needed RI of the HEC-PR scheme increases with the increase of the number of receivers significantly but not strict linearly. The reason is clear: although all of the receivers are independent, they also can recover some common missing packets by retransmitting a small part of identical packets with the HEC-PR scheme. As a matter of fact, the probability of sharing identical packets among different receivers will increase with the increase of the number of receivers. This leads to the total needed RI of the HEC-PR scheme increase with the number of receivers significantly but not strict linearly. In other words, the speed of the increase of the total needed RI will slow down with the increase of the number of receivers.

Additionally, as shown in this figure, the AHEC scheme always outperform the HEC-PR scheme and the AFEC scheme, because it can choose the best scheme automatically among different HEC schemes. However, from this figure, we also can see that the performance of

the HEC-PR scheme is very close to the AHEC scheme when the number of receivers is no more than two or the average link PLR is less than 10<sup>-2</sup>. Because the implementation of the HEC-PR scheme is very simple due to without any encoding and decoding algorithm, the HEC-PR scheme should be considered for such as those real-time multicast services with small group size and small average link PLR.

## 4.3. The Effect of the Parameter k

Finally, we study the effect of the parameter k in the AHEC scheme in this section. For the convenience of analysis, we searched for the optimum parameters for the AHEC scheme with the fixed parameters  $N_{recv}=5$  and  $N_{rr,max}=2$  under different average link PLR of 0.01,0.05 and 0.10. Actually, on the effect of the k, the tendency is similar for any average link PLR. Here we only take three typical examples to demonstrate the tendency of its effect. Part of those optimum parameters is listed in the Table 5.

Table 5. Optimum Parameters of the AHEC Scheme with  $N_{recv}$ =5 and  $N_{rr,max}$ =2

k		20	40	80	120	160	200
PLR=0.01	$N_p$	1	1	1	1	2	2
	$N_{cc}$	0	0	0	0	0	0
PLR=0.05	$N_p$	2	3	6	8	10	12
	$N_{cc}$	1	1	1	1	1	1
PLR=0.10	$N_p$	4	7	12	16	21	25
	$N_{cc}$	2	2	2	2	2	2

Note that although here we only show the results for the AHEC scheme with the length of k being less than 200, it should be clear that the AHEC scheme is suitable for any length of k upon requirements; moreover, the higher k is employed by the ideal erasure error code, the more efficient code rate can be adopted. From this table, we can see that the parameter  $N_{cc}$  is always invariable under certain average link PLR. That is, we only need to change the parameter  $N_p$  for the AHEC scheme according to different value of k under certain average link PLR. Note that the variable value of k means different multimedia data rate under certain delay constraints if the packet size is constant. Obviously, this feature of the AHEC scheme can simplify its implementation in real-time multicast scenarios with variable source data rate.

Upon those optimum parameters with different length of k, the total needed RI of the AHEC scheme is obtained and shown in Figure 4.

From this figure, we can see that the total needed RI decreases significantly when the parameter k increases from 10 to 60. When k is more than 60, however, this parameter has only a little effect on the performance of the AHEC scheme. Note that different k values mean different delay constraints or different source data rates if the packet size is fixed. Therefore, under certain delay

constraints with fixed packet size, the higher the multicast source data rate is, the better performance can be achieved in the AHEC scheme. Moreover, since the stable good performance can be obtained if the parameter k is more than 60, a suitable fixed short length of k ( $\geq 60$ ) can be always adopted when the data rate is high enough to provide good delay performance. On the other hand, the parameter k is only associated with the delay constraints if both the source data rate and the packet size are fixed: to guarantee a certain PLR requirement, shorter delay constraints the system has, shorter length of the parameter k in the AHEC scheme has to be adopted so that more RI needed. Therefore, the AHEC scheme also provides a good way for the tradeoff between the total needed RI and delay constraints by choosing different k.



Figure 4. The total needed RI of the AHEC scheme with  $N_{recv}=5$  and  $N_{rr,max}=2$ 

# 5. Conclusions

In this paper, we propose a novel Adaptive Hybrid Error Correction (AHEC) scheme by choosing the most suitable HEC scheme among HEC-PR scheme, Type I HARQ scheme and Type II HARQ scheme under strict delay constraints. Using the proposed mathematical framework for the AHEC scheme, we can design the most suitable parameters of the AHEC scheme for guaranteeing a certain PLR requirement under strict delay constraints with minimum needed RI. By applying the proposed AHEC scheme in a typical Wireless DVB scenario for performance analysis and comparisons, we have found:

 The AHEC scheme outperforms the HEC-PR scheme and AFEC scheme in all cases. However, when either the group size or the average link PLR is small enough, the performance of the HEC-PR scheme is very close and even equal to the best performance of the AHEC scheme. Due to the simplicity of the HEC-

#### A NOVEL ADAPTIVE HYBRID ERROR CORRECTION SCHEME FOR WIRELESS DVB SERVICES

PR scheme without any encoding or decoding algorithm, this scheme should be considered in the real-time multicast scenarios with small group size or small average link PLR.

- 2) In most cases, the best performance of the AHEC can be obtained with variable network and channel conditions by only varying the parameter  $N_p$ . Thus, the AHEC scheme is usually robust and simple to implement for practical systems.
- 3) The length of k has a great effect on the performance of the AHEC scheme. The performance increases in case the scenario allows for the choice of a bigger k. It indicates that the AHEC scheme is very suitable for the real-time multicast systems with high data rate. Also, it provides a good way for the tradeoff between the total needed RI and the strict delay constraints.

In this paper, for simplifying the analysis, we have made a strong assumption: all of the receivers are independent and experience *i.i.d* channel with uniform distribution. That is, we do not consider the effect of temporal and spatial correlation in real wireless channels. For future works, we will analyze the performance of the AHEC scheme based on accurate Gilbert-Elliot [20,21] channel model for practical wireless channels. First results, however, do confirm that all conclusions made in this paper remain valid and that spatial and temporal correlation shifts the architectural choice in the parameter space but do not change the conclusions.

# 6. Appendix: PDF of N<sub>req,max</sub> and Derivation of P<sub>req</sub>(*i*,*c*,*j*)

First, to derive the PDF of  $N_{req,max}$ , we define two basic probabilities: one is the probability of  $N_{req}(j)$  being *i*, which is denoted by  $P_{req}^{=i}(j)$ ; the other is the probability of  $N_{req}(j)$  being less than *i*, which is denoted by  $P_{req}^{<i}(j)$ .

Using (1), these two probabilities can be calculated as

follows, respectively:

$$P_{req}^{=i}(j) = \Pr(N_{req}(j) = i) = P(N_p + i, N_{blk}, P_e(j))$$

$$P_{req}^{
(20)$$

where  $1 \le i \le k$  and  $1 \le j \le N_{recv}$ 

To simplify the analysis in this paper, it is assumed that all of the receivers have the same PLR level of  $\overline{P}_e$  as defined by (12). Thus, following (20) we define:

$$P_{req}^{=i} = P_{req}^{=i}(j) = P(N_p + i, N_{blk}, \overline{P_e})$$

$$P_{req}^{
(21)$$

where :  $\forall j \in \{1, 2, ..., N_{recv}\}$ 

Now let  $P_{N_{req,max}}^{i}(h)$  be the probability of h receivers lost  $N_{p}+i$  packets and the other  $N_{recv}-h$  receivers lost less than  $N_{p}+i$  packets. Based on the definitions above and using (21), this probability can be obtained:

$$P_{N_{req,\max}}^{i}(h) = \binom{N_{recv}}{h} (P_{req}^{=i})^{h} (P_{req}^{(22)$$

Then, Let  $P_{N_{req,max}}^{i}$  denotes the probability of  $Pr(N_{req,max}=i)$ . Upon (22), we can obtain the PDF of  $N_{req,max}$ :

$$P_{N_{req,\max}}^{i} = \Pr(N_{req,\max} = i) = \sum_{h=1}^{N_{rec}} P_{N_{req,\max}}^{i}(h), 1 \le i \le k$$
 (23)

Secondly, let's consider the probability  $P_{req}(i,c,j)$  for the *j*-th receiver in the following. Similarly, to simplify the analysis in this paper, we assume that all of the other  $N_{recv}$ -1 receivers except for the *j*-th receiver have the same PLR level of  $\overline{P_e}^{j}$  as defined as follows:

$$\overline{P}_{e}^{j} = \frac{\sum_{i=1}^{j-1} P_{e}(i) + \sum_{i=j+1}^{N_{recv}} P_{e}(i)}{N_{recv} - 1}$$
(24)

Based on this assumption for those  $N_{recv}$ -1 receivers, we define:

$$\hat{P}_{req}^{=i} = P_{req}^{=i}(d) = P(N_p + i, N_{blk}, \overline{P}_e^j)$$

$$\hat{P}_{req}^{
(25)$$

where:  $\forall d \in \{1, 2, ..., N_{recv}\} - \{j\}$ 

Now concerning those  $N_{recv}$ -1 receivers except for the *j*-th receiver, let  $P_{N_{req,max}}^{i}(h, j)$  be the probability of *h* receivers lost  $N_{p}$ +*i* packets and the other  $N_{recv}$ -*h*-1 receivers lost less than  $N_{p}$ +*i* packets. Using (25), then, the probability can be calculated by:

$$P_{N_{req,\max}}^{i}(h,j) = \binom{N_{recv} - 1}{h} (\hat{P}_{req}^{=i})^{h} (\hat{P}_{req}^{(26)$$

Actually, the calculation of  $P_{req}(i,c,j)$  should be divided into two parts according to two different cases:

1) One part is the probability of  $Pr(N_{req,max} = i, N_{req}(j) = c)$ with i=c, in which case the number of missing packets in one block are no more than  $N_p+c$  for any receiver among those  $N_{recv}-1$  receivers except for the *j*-th receiver. Using (26),  $P_{req}(i,c,j)$  can be expressed as:

Copyright © 2008 SciRes.

I. J. Communications, Network and System Sciences, 2008, 2, 105-206

$$P_{req}(c,c,j) = \Pr(N_{req,\max} = c, N_{req}(j) = c)$$
  
=  $P(N_p + c, N_{blk}, P_e(j)) \left( \sum_{h=0}^{N_{recv} - 1} P_{N_{req,\max}}^c(h, j) \right)$  (27)

2) The other part is the probability of  $Pr(N_{req,max} = i, N_{req}(j) = c)$  with i > c, in which case at least one receiver among those  $N_{recv}-1$  receivers except for the *j*-th receiver lose  $N_p+i$  packets in one block and all of the other receivers lose less than  $N_p+i$  packets in the block. Similarly, in this case,  $P_{req}(i,c,j)$  should be calculated by:

$$P_{req}(i, c, j) = \Pr(N_{req, \max} = i, N_{req}(j) = c)$$
  
=  $P(N_p + c, N_{blk}, P_e(j)) \left( \sum_{h=1}^{N_{recv} - 1} P_{N_{req, \max}}^i(h, j) \right)$  (28)

To integrate (27) and (28) for the expression of  $P_{req}(i,c,j)$ , we define a function  $f_{cmr}(x1,x2)$  (where  $x1 \ge x2$ ) as follows:

$$f_{cmr}(x1, x2) = \begin{cases} 0, x1 = x2\\ 1, x1 > x2 \end{cases}$$
(29)

As a result, based on (27), (28) and (29), the calculation of  $P_{req}(i,c,j)$  can be expressed as the following form:

$$P_{req}(i,c,j) = P(N_p + c, N_{blk}, P_e(j)) \times \left(\sum_{h=f_{cmr}(i,c)}^{N_{recv}-1} P_{N_{req,max}}^{\max(i,c)}(h,j)\right)$$
(30)

#### 7. References

- U. Varshney, "Multicast support in mobile commerce applications," Computer, vol. 35, no. 2, pp. 115–117, February 2002.
- [2] IEEE 802.11, "Wireless LAN medium access control (MAC) and physical layer (PHY) specifications," 1999.
- [3] ETSI EN 302 304 V1.1.1 (2004-11), Digital Video Broadcasting (DVB); Transmission System for Handheld Terminals (DVB-H).
- [4] C. Perkins, et al., "A Survey of Packet Loss Recovery Techniques for Streaming Audio", IEEE Network, vol. 12, no. 5, pp. 40–48, September–October 1998.
- [5] Draft ETSI TS 102 034 v0.14, "Digital Video Broadcasting (DVB); Transport of DVB Services over IP-based Networks; Part 1: MPEG-2 Transport Streams," May 2003.
- [6] H. Schulzrinne, S. Casner, et al., "RTP A Transport

Protocol for Real – time Applications," RFC 1889, January 1996.

- [7] H. Fujisawa, K. Aoki, et al., "Estimation of Multicast Packet Loss Characteristic due to Collision and Loss Recovery using FEC on Distributed Infrastructure Wireless LANs," IEEE WCNC, pp. 21–25 March 2004.
- [8] J. Ott, S. Wenger, N. Sato, et al., "Extended RTP profile for RTCP-based feedback," draft-ietf-avt-rtcp-feedback-11.txt, August 2004.
- [9] Q.H Du and X. Zhang, "Adaptive Low-Complexity Erasure Correcting Code Based Protocols for QoS Driven Mobile Multicast Services," IEEE Second International Conference on Quality of Service in Heterogeneous Wired/Wireless Networks, pp. 22–24, August 2005.
- [10] J. Nonnenmacher, E.W. Biersack, et al., "Parity-Based Loss Recovery for Reliable Multicast Transmission," IEEE/ACM Trans. Networking, vol. 6, August 1998.
- [11] B. Adamson, C. Bormann, M. Handley, and J. Macker, "Negative-Acknowledgment (NACK) – Oriented Reliable Multicast (NORM) Protocol," RFC 3940, November 2004.
- [12] G. Tan and T. Herfet, "Application Layer Hybrid Error Correction with Reed-Solomon Code for DVB Services over Wireless LANs," the 3<sup>rd</sup> IEEE International Conference on Wireless Communications, Networking and Mobile Computing (WiCOM), Shanghai, China, September 2007.
- [13] G. Tan and T. Herfet, "Optimization of an RTP Level Hybrid Error Correction Scheme for DVB Services Over Wireless Home Networks Under Strict Delay Constraints," IEEE Trans. Broadcasting, vol. 53, no 1, Part 2, pp.297–307, March 2007.
- [14] S. Pejhan, M. Schwartz, and D. Anastassiou, "Error control using retransmission schemes in multicast transport protocols for real-time media," IEEE/ACM transaction on networking, vol 4, no 3, pp.413–427, June 1996.
- [15] S.R. Chandran and S. Lin, "Selective-repeat-ARQ schemes for broadcast links," IEEE Trans. Commun., vol.40, no.1, pp.12–19, January1992.
- [16] R.G. Gallager, "Low Density Parity-Check Codes," Cambridge, MA: MIT Press, 1963.
- [17] Shokrollahi, "Raptor Codes," IEEE Trans. Information Theory., vol. 52, no.6, pp.2551–2567, June 2006.
- [18] G. Tan and T. Herfet, "On the Architecture of Erasure Error Recovery under Strict Delay Constraints," in preparation for submitting to IEEE Trans. On Information Theory, 2008.
- [19] http://www.isi.edu/nsnam/ns/.
- [20] E.N. Gilbert, "Capacity of a burst-noise channel," Bell Syst. Tech. J., vol.39, pp.1253–1265, September 1960.
- [21] E.O. Elliott, "Estimates of error rate for codes on burstnoise channels," Bell Syst. Tech. J., vol.42, pp.1977– 1997, September 1963.

198



# A Co-verification Method Based on TWCNP-OS for Twoway Cable Network SOC

Chong LI<sup>1</sup>, Xiaotong ZHANG, Yadong WAN, Qin WANG

Department of Computer Science and Technology University of Science and Technology Beijing, Beijing, P.R.China E-mail: <sup>1</sup>lichong0564@163.com

# Abstract

Co-verification is the key step of software and hardware codesign on SOC. This paper presents a hw/sw coverification methodology based on TWCNP-OS, a Linux-based operating system designed for FPGA-based platform of two-way cable network (TWCNP) SOC. By implementing HAL (hardware Abstraction level) specially, which is the communications interface between hardware and software, we offer a homogeneous Linux interface for both software and hardware processes. Hardware processes inherit the same level of service from kernel, as typical Linux software processes by HAL. The familiar and language independent Linux kernel interface facilitates easy design reuse and rapid application development. The hw/sw Architecture of TWCNP and design flow of TWCNP-OS are presented on detail. A software and hardware co-verification method using TWCNP-OS is proposed, through the integrated using of Godson-I test board and TWCNP, which realizes the combination of design and verification. It is not a replacement of the coverification with generic RTOS modeling, but is complementary to them. Performance analysis of our current implementation and our experience with developing this system based on TWCNP-OS will be presented. Most importantly, since the introduction of TWCNP-OS to our FPGA-based platform, we have observed increased productivity among high-level application developers who have little experience in FPGA application design.

Keywords: CM, SOPC, MIPS, Co-verification, FPGA, HAL

# 1. Introduction

Now CATV mainly transmits via analog signals in some regions and countries. In order to realize A/D transmission, two-way network chip based on HFC (Hybrid Fiber-Coax) is becoming the core technology and key equipment in interactive digital television and other integrated digital services. SOC enables the realization of all these functions in one chip and has reusable technology of IP module to develop the SOC which has the independent intellectual property has important meaning.

The hardware and software codesign is commonly used in SOC design. It is different from traditional design methods that simulation and verification is used to find and rectify mistakes in time and optimized the system in the end. With the complexity increasing in SOC design, scientific design method and high efficient

verification are growing more important [1]. Hw/sw coverification does not just occur at the system integration point but rather throughout the design process. A high efficient verification platform should provide simple and scientific verification environment to ensure the performance of whole system. Co-verification includes hardware verification. software verification and software-hardware interface verification. Using concurrent design and co-verification to shorten project period and combine hardware behavior model with software running environment, construct a high efficient co-verification platform has been the research hotspot recently.

While traditional hw/sw codesign researches have produced encouraging results in the area of hw/sw partitioning, cosimulate, cosynthesis, and co-verification, most of them rely on self-contained design environments that are based on their specific input languages or library API's [2]. As a result, migrating existing software designs to two-way cable network platform using these traditional hw/sw codesign methodologies would have incurred major re-engineering efforts, including learning a new language and API, getting familiar with a new design environment and reimplementing existing designs in the new language environment.

So, an easy to use hw/sw interface that allows rapid application development and migration should be (1) familiar and intuitive to both software and hardware engineers; and (2) language independent. We achieve this goal by setting hw/sw boundary at the embedded operating system kernel level.

Embedded operating systems (EOS) have become one of the most important components of embedded systems due to the growing complexity of the system functionalities as well as the increasing time-to-market pressures. With this trend, a demand for fast cosimulation of hardware and embedded software including an EOS is also becoming stronger in order to validate the functionality of the overall systems. With this trend, a demand for fast cosimulation of hardware and embedded software including an EOS is also becoming stronger in order to validate the functionality of the overall systems.

In this paper, we present TWCNP-OS, an operating system designed specifically for FPGA-based platform of two-way cable network (TWCNP) SOC, which is key device based on DOCSIS (Data-Over-Cable Service Interface Specifications) on HFC network whose CPU is Godson-I. Under TWCNP-OS, hardware and software share the same familiar Linux interface and the same level of support from the OS kernel. We use the concept of hardware process [3], which is the same as a normal Linux process except its "program" is an FPGA hardware design instead of software program. HAL is implemented, and communications between hardware and software are accomplished through it, which uses conventional Linux inter-process communication (IPC) mechanisms, such as shared file, pipe, shared memory, signal, and message passing. Hardware processes have access to system resources as their software counterparts, such as the general file system, standard input, standard output.

By building a cross GNU/GCC compiler for Godson-I, we can develop software in host machine (PC).Software designs can be developed in C/C++ language development environment a designer is familiar with. For hardware designs to communicate with the kernel, TWCNP-OS defines a standard message passing network that resembles the software system call interface by maintaining the hw/sw interface-HAL at the kernel level. This standardized network allows hardware designs be developed in any hardware language environment of choice.

The remainder of this paper is organized as follows. Section 2 surveys related work on hw/sw cosimulate/coverification with RTOS supports, and analyses the difference all of them including our TWCNP. Section 3 introduces the system architecture, particularly the software and hardware architecture of CM platform in tow-way CATV. Section 4 elaborates on the co-verification approach based on TWCNP-OS, and its detailed design process is presented. Section 5, we give the system testing and performance comparison of the platform designed by this method.

## 2. Related Work

Hw/sw co-verification has been studied around the world for more than a decade. Simulator is an usual method of hw/sw co-verification, which must join with other parts of the system verification problem to create a complete verification solution. A number of commercial cosimulators have come onto the market. The cosimulators are currently one of indispensable CAD tools in the design of embedded systems and systemson-chip.

There is no single hardware-software co-simulation problem and as a result there is no single tool or technology that can successfully solve all the problems associated with co-verification of hardware and software. Each developer must trade off the desired performance against the level of accuracy. Each developer must trade off the desired performance against the level of accuracy. For hardware designers the trade-off is principally between software simulation and the relatively faster technologies of hardware acceleration, emulation or rapid prototyping. For software developers of embedded system application software the primary need is simulation speed and the trade-off is among a variety of simulation approaches which achieve greater degrees of speed at the cost of diminishing the accuracy of the modeling of the hardware.



Figure 1. Trade-off between simulation accuracy and speed

HW emulation is a very important alternative to overcome the SW speed problems [4]. A co-verification platform using C++ simulator and FGPA (Field-Programmable Gate Arrays) emulator is presented in [5], gets a good simulation speed, and provides an accuracy/efficiency tradeoff through various abstraction levels (Figure 1). When the C/C++ abstraction is at the algorithm level, the simulation speed can reach about 1 MHz, but the simulation results have no cycle accuracy. The speed of hardware emulator is typically up to 1 MHz, and they preserve cycle accuracy.

However, since the cosimulators do not feature explicit supports for RTOSs, a designer needs to use an ISS in many cases to run embedded software including an RTOS. Due to the slow execution speed, extensive simulation of large software is impossible.

In the recent years, several research efforts have been made to model RTOSs for system-level design and cosimulation. SoCOS presented in [6] is a system-level design environment where the OSAPI library provides generic RTOS system calls to application software. OSAPI is a virtual RTOS to enable native execution of embedded software. After simulation-based validation. the OSAPI calls are replaced with the system calls of the actual RTOS used in the final implementation to obtain the final software code. In [7], a similar approach is presented. A main difference is that the RTOS model in [7] is build upon an existing system-level design language, i.e., SpecC [8], so that existing CAD tools such as simulators can be used. Techniques presented in [9] also use the SpecC language to model the preemptive behavior. In [10], an OS model is proposed for fast and time-accurate cosimulation. The model focuses on accurately modeling the RTOS overhead during task execution as well as the preemptive behavior. In [11], a method which automatically generates RTOS-dependent software from SystemC description is presented. The method replaces SystemC's constructs for concurrency and communication with corresponding RTOS service calls. In this sense, it can be considered that SystemC involves a simple RTOS model in itself. With the development of Embedded OS, [12] gives us a RTOScentric hw/sw co-simulator, which features cosimulation with functional simulation models of hardware written in C/C++ and co-simulation with HDL simulators, supports a complete simulation model of an RTOS based on uITRON [13].

A common weakness of these OS models is that they support only a limited set of RTOS services in order to make the models generic and independent of specific RTOSs. for example, [12] have more than 80 service calls but only for µITRON-based RTOSs platform. However, the RTOS model in [7] supports only 16 service calls. These services may need to be fully utilized in order to write high quality software. Therefore, it is easily imagined that the quality of software automatically generated by these previous methods is lower than that of hand-crafted RTOSdependent software. Instead, we have developed fully supports RTOS services by porting the standard MIPS/Linux kernel to our platform, such high-quality software can be designed and simulated. But there are a weakness for hardware-dependent, and have to put a Virtual hardware level in HAL when some hardware component haven't finished.

TWCNP shares a similar design philosophy as

BORPH [3] in providing a unifying coarse-grain hardware (FPGA) and software component interface. They are not a complete system to perform typical hw/sw codesign tasks such as partitioning, cosynthesis, cosimulate, or verification. Instead, by providing basic Linux OS services, it acts as a platform on which these tasks can be carried out. In fact, our TWCNP is not a replacement of the cosimulators with generic RTOS modeling, but is complementary to them.

The main contribution of TWCNP-OS is that by leveraging conventional Linux semantics to FPGAbased platform of two-way cable network (TWCNP) SOC, it provides a unique, unified environment for both FPGA and software application designers. The Linux semantics is familiar to developers across many research domains, thus lowering the barrier-to-entry into FPGAbased SOC. Furthermore, since TWCNP-OS is implemented as an extended Linux kernel, a TWCNP-OS managed system may leverage all commodity Linux software applications for developing, testing, benchmarking, and deploying FPGA applications.

### 3. Two-way CM Platform Architecture

The CM (cable modem) platform is based on HFC network, which can increase the transmission quality efficiently, and integrate the CATV and internet digital network. Using the CM can expediently develop many services based on Internet protocol, such as VoIP, VOD, HDTV and high-speed web browsing. The design of CM includes three parts: evaluation board design, embedded software design and SOC design. Verification is used through the whole process of SOC development. The software and hardware co-verification is stressed in this paper.



(a) FPGA evaluation board (b) TWCNP SOC hardware frame

Figure 2. Overall frame of FPGA emulation board

The hardware architecture of CM is described in Figure 2. It includes three modules:

- 1) EuroDOCSIS protocol processor: It includes independent developed physical layer modem (Jupiter), EuroDOCSIS MAC protocol processor and radio frequency interface circuit.
- 2) BUS and peripheral equipment module: Interconnection through AMBA bus and peripheral equipment includes IIC, GPIO, UART, PWM and Ethernet MAC.

3) Godson-I CPU: Godson-I IP core, Based on MIPS instruction, can apply to both the universal and embedded system, has a good compatibility.

The software architecture which is described in Figure 3, includes three main parts:

- Operating System: based on MIPS/Linux configuration, and making Linux run on CM based on Godson-I, including Bootloader and OS kernel, It supports memory management, file management, task schedule, and modularization.
- Hardware driver: it ensures hardware working highefficiently and properly. The driver such as GPIO, IIC, network adapter can either be compiled into OS kernel or be loaded dynamically.
- 3) EuroDocsis protocol stack and CM application: realizing central functions defined in EuroDocsis, it provides software guarantee to bridge CM and CMTS (cable modem terminal system). Application software defines a set of user interface and ensures function maintenance and management to CM.

DOCSIS STACK					CM APPL	ICATION	A	pplication	n Level
Task Manager	Device Manager	File System	Clock Manager	HAL	Memory Manager	Loadable Module	Embedded TCP /IP Stack	Device Drivers	OS Level
			Hard	Ware of	TWNCP			Phys	ic Leve

Figure 3. Software frame of two-way cable network CM

OS is the cornerstone of the software design platform, which manages all software and hardware directly, provides interface to application and system call. It is also used to validate hardware and software function, guide the hardware design to be more reasonable.

# 4. Co-verification Based on TWCNP-OS

There are two common verification methods on soc design.

- The first is software verification which simulates the system behavior through establishing emulation model. Through this method, system's detailed status should be observed but the run time is much longer. ISS (Instruction set simulation) is an example of it.
- 2) The second is the hardware verification. Now the prevalent one is the SOPC based on FPGA, which can online program and has accurate clock cycle. This way has a higher emulation speed, but the system's status cannot be observed. It includes three aspects: software verification, hardware verification, and software-hardware interface verification.

Combining the merits of these two kinds of method, a co-verification method based on TWCNP-OS is proposed in this section, which is supposed that the hw/sw partitioning and hardware cosynthesis have finished. Firstly, we introduce the co-verification platform, give a general idea of porting the MIPS/Linux to our SOC. Secondly, Software and Hardware co-verification based on TWCNP-OS will be introduced.

The co-verification platform (Figure 4) is composed of host, two self-existent target machines and equipments which connects them each other, such as UART, USB and network. Host is a common PC, which provides cross-compilation environment, console and HDL simulator. Godson-I test board and FPGA development board (TWCNP) constitute target machine. Godson-I test board adopts Godson-I CPU and it8172G chip sets, MIPS instruction set and Linux OS can be run in it. This one is mainly used in the development of DOCSIS protocol stack. Being composed of FPGA vp200, peripherals and interfaces, FPGA board is the hardware verification platform of the whole SOC system, which is mainly used in the verification of EuroDOCSIS protocol processor and related hardware logic.



Figure 4. Co-verification platform of the TWCNP

#### 4.1. Realization of TWCNP-OS

TWCNP-OS is an operating system designed for TWCNP based on MIPS processor. It extends a standard Linux kernel to include support for FPGA's in our platform. Treating FPGA's as both coprocessors and processor, TWCNP-OS treats FPGA's as hardware process in the system as normal computational resources. The interface between hardware and software is HAL. All processes can therefore be either software programs, or hardware designs running on FPGA's. Therefore, to the rest of the system, communicating with a hardware process is no different from communicating with a normal Linux process. This homogeneous handling of hardware and software in the kernel forms the foundation of coarse grain hw/sw codesign boundary. Figure 5 depicts this conceptual block diagram.



Figure 5. TWCNP-OS extends a traditional Linux system with hardware process support. HW/SW processes share the same I/O interface and communicate by HAL

Like embedded OS, TWCNP-OS includes crosscompiler, Bootloader and OS Kernel too. Crosscompiler is necessary to design an embedded OS, which is a compiler capable of creating executable code for a platform other than the one on which the compiler is run. It is a tool that one must use for a platform where it is inconvenient or impossible to compile on that platform, like microcontrollers that run with a minimal amount of memory for their own purpose. Whole process of building a cross-compiler can find in [15].

Bootloader is used to boot the machine, such as initializing basic CPU registers, UART, SDRAM controller, and copying OS kernel from Flash to memory. Its design is derived from PMON [14], which is a freeware ROM-monitor developed for early LSI Logic MIPS R3000 evaluation boards. Since its creation, PMON has become a very common firmware for MIPS evaluation boards and development systems. It contains support for several boards, R4000-style exceptions, and uses the GNU make system. So the major revision is the address of hardware components. There are three important parts will be down:

- 1) CPU register in Godson-I have to be set up the normal value, such as interrupt, timer, cache, memory management;
- 2) The drivers of hardware components should be taken effect to use UART, serials and console etc. We can debug and interact with target machine.
- 3) Only after moving the kernel from Flash kernel space to memory, the kernel can boot and manage whole system.

Linux is one of the common embedded OS. For MIPS architecture of Godson-I, OS mainly adopts MIPS/Linux schema, using ANSI C, MIPS assembly language and cross-compiler-MIPS-GCC 3.23 as SDK. TWCNP-OS is divided into two parts. The hardware independent part is realized by MIPS assembly language. The primary design flow is described in figure 6. We will introduce the design of TWCNP-OS on the basis of MIPS/Linux kernel source code in the following part. Overall design steps as following:

- 1) Hello World! Get board setup, serial porting working, and print out "Hello, world!" through the serial port.
- 2) Get early printk working Make the first TWCNP-OS image and see the printk output from kernel.
- 3) Serial driver and serial console Get the real printk() function working with the serial console.
- 4) KGDB KGDB can be enormously helpful in our development.
- 5) CPU support Because Godson-I CPU is not currently supported, we need to add new code that supports it in arch directory of Linux source codes.
- 6) Board specific support Create your board-specific directory. Setup interrupts routing/handling and kernel timer services.

- 7) HAL It makes the upper layer software be independent of the bottom hardware and manages or simulates a large number of the bottom hardware.
- 8) Ethernet drivers We should already have the serial port working before attempting this. With ethernet driver working, we can set up a NFS (network file system) root file system which gives you a fully working Linux user space.
- 9) ROMFS root file system Alternatively you can create a user space file system as a ROMFS image stored in a ramdisk, or Flash root file system.

Here we will emphatically present HAL [19] (Figure 6(b)), which is the interface between Hardware and software, which includes device driver, configuration manager, control manager and data manager and VHL (virtual hardware level). It includes all the functions that manage and schedule the hardware (hardware process) uniformly for the SOC, and can support a many of different hardware devices or virtual hardware devices (VHL). So we can send messages to hardware using standard system call just as to a software process by HALIF (HAL interface). Overall control module is the key of HAL implementation. Configuring manager takes charge of configuration channel. Control manager supplies control channel. Data manager takes charge of two-way data channel. Driver module regards all hardware devices which are managed by HAL as one device, and its run driver is regard as a hardware process.



Figure 6. Design flow of the TWCNP-OS, comparing to generic Linux kernel, we give a HAL on kernel

### 4.2. Software Verification Based on TWCNP-OS

Software is an important part of the system, verification of software is necessary. Software verification mainly uses ISS, though it has accurate clock cycle, the simulation speed is comparatively slow. The software verification proposed in this paper simulates hardware on target machine. Host and Godson-I test board to validate the validity of software directly on function level. Software verification is mainly used to DOCSIS protocol stack module and the relative CM application program. According to DOCSIS specification, DOCSIS protocol stack module is divided into hardwaredependent part and hardware-independent part. The hardware-independent software realizes user level function such as VOD etc.

In the traditional software and hardware concurrent design, the design of hardware-dependent software can not progress until the accomplishment of the hardware design. VHL can efficiently resolve this problem. The VHL established on OS can be one or a group of virtual hardware, which simulate hardware to carry out function verification. VHL is extension of driver. it can not only apply uniform interface to application, drive hardware, but also response according to hardware's actual characteristic. If the hardware design has been accomplished, we only need to simply change the VHL to drivers and the software need not to be changed. Thus, hardware-dependent software design and verification can be processed without hardware environment which has great meaning.



Figure 7. The sketch map of software co-verification, the left part of dotted line denotes the hw-dependent software, the right one is hw-independent software. Each uses a different verification method

TWCNP-OS make it possible to use VHL in software verification. The DOCSIS protocol stack and CM application program validated on development board can drive hardware on TWCNP-OS. Godson-I test board is mainly used in software verification. Target host is the auxiliary equipment which provides convenience for the verification. The software verification flow is described in Figure 7, the broken line divided software into hardware-dependent part and hardware-independent part. The hardware-independent part is validated directly on development board and run on target FPGA board. The hardware-dependent part is mainly in DOCSIS protocol stack module which needs the support of MAC coprocessor. Using modularize mechanism to establish virtual MAC co-processor on VHL can simulate MAC co-processor to realize software-hardware interaction. Thus the software is independent on or weak dependent on hardware. Virtual MAC co-processor construct standard data, notify the hardware-independent part in DOCSIS module through interrupt to process data analysis. The verification of CM application program use C/S mode, the server responses to the client by simulating CMTS.

# 4.3. Hardware Verification Based on TWCNP-OS

Hardware is the support platform of embedded system. Hardware verification validates through simulating each unit's behavior and mainly has two modes. One mode is pre-simulation which simulate all hardware modules to validate if the hardware accord with the design requirement. The common simulation tool is ModelSim [17], NC-verilog and NC-VHDL [18]. The other mode is integrated post-simulation which can put composite latency files into integrated simulation model to estimate the effect brought out by gateway delay. The hardware verification introduced in the following base on the accomplishment of pre-simulation and is in the domain of FPGA-based post-simulation.

Any verification should carry out according to some standards. We adopt TWCNP-OS and DOCSIS protocol stack module as test program to validate hardware. TWCNP-OS which bases on Linux structure has high reliability and stability through long time verification. DOCSIS protocol stack module has passed the software verification. We can validate the MAC co-processor as long as modify the virtual MAC co-processor.

The TWCNP-OS-based hardware verification has high stability and simulation speed. At the time of the accomplishment of verification of hardware and software, the hw/sw interface has been validated. This verification method is on the basis of passing presimulation for each hardware module and validates the function of hardware system. For the complexity of the system, the hardware system is divided into two parts: MAC co-processor and other hardware logics. The hardware verification has three steps:

- The first step is the verification of MAC co-processor. We first connect FPGA with the bus of Godson-I test board, download the hardware logic of MAC coprocessor to FPGA, regard MAC co-processor as a peripheral. Then modify the virtual MAC coprocessor, delete its function of simulation hardware, make it be a driver to control hardware. At last, validate it through DOCSIS protocol stack. Thus we can easily find out errors in it. Furthermore, the concurrent process of hardware verification can improve the efficiency.
- 2) The second step is validating other hardware logic and software-hardware interface using TWCNP-OS and FPGA board. The hardware is validated respectively by calling the drivers compiled to drive hardware.
- 3) The third step is the system verification on FPGA board by using TWCNP-OS and software module after integrating the hardware logic. The function and performance test is processed in this step. The detailed process is introduced in next section.

# 5. System Testing and Performance Comparison

The whole system on the actual running environment for

testing is the last step of verification. Therefore, we will give a topology of cable network (Figure 8) and construct a test environment (Figure 9) accordingly.



Figure 8. The topology of cable network



Figure 9. Performance test environment of CM

Stability, time delay and network congestion are the performance target of CM. Stability is a basic target for any device. Here it means data can be transferred to CMTS safely without losing data. Time delay is another target to estimate performance, if it can't meet the some special requirement of the system, it will be no-good. We will test the response time between CPE (customer premises equipment) and PC. Network congestion is an important problem, it is good rule to evaluate whether the device is excellent or not. Because some equipment can work well without network congestion, but it will induce packages losing when it happens network congestion. So it is a complemental test to stability.

In the following parts in this section, we will give some tests to get the above performance targets respectively in TWCNP and SM5100. SM5100 is a mainstream CM produced by Motorola, and its performance can delegate a standard of CM in a certain sense. Special declaration, test results will be influenced by real environment.

#### Test 1

To get the response time of CM, We send 100000 IP data packages by ping command from CPE to PC via TWCNP and SB5100 CM respectively, each package is 74bytes, the results as table 1.

Table1. Functional test and comparison of CM emulation platform

Test platform	$T_{ave}$	T <sub>ans</sub>	$T_{ans}$ (max)	р	
TWCNP	12ms	10ms	17ms	0	
SM5100	9ms	6ms	15ms	0	
Parameters	$N_{sum} = 100000$ , packet size=74 bytes				

The following is the formula of calculate different response time. Here,  $T_{ans}$  is response time, p is rate of loss package,  $T_{ave}$  is average response time,  $T_{sent}$  is time

to reach destination,  $T_{resv}$  is response data package cost time,  $N_{lost}$  is number of loss package,  $N_{sum}$  is total number of package.

$$T_{ans} = T_{sent} + T_{resv} \tag{1}$$

$$P = \frac{N_{lost}}{N_{sum}} \times 100\%$$
 (2)

$$T_{ave} = \frac{\sum_{i=1}^{N_{som}} (T_{sent}[i] + T_{resv}[i])}{N_{sum} - N_{lost}} \times (1 - P)$$
(3)

From testing results of Table 1, we can get that the average network response performance is similar between TWCNP and SM5100.

#### Test 2

To get network congestion and stability results, we fabricate a mass of normal at the speed of 1.4 MB/s and congestion data packets at a speed beyond upper line by CommView on CPE, and send them to PC via TWCNP an SM5100 about 24 hours. At the same time, we will capture the throughput on PC.



Figure 10. The stability test of TWCNP



Figure 11. Data throughput analyses of TWCNP and Motor 5100 CM

Sampling 120 minutes normal packets (Figure 10)

and Sampling 90s' data on congestion time (Figure 11), we get that throughput of TWCNP is still up to 1.52M B/s, MOTOROLA 5100s' is 1.67 M B/s under network congestion. TWCNP can transfer all packets to PC at same speed. The result of this case study demonstrates that it gets a high stable throughput.

From the above tests, we can find that the performance of our TWCNP is equivalent or a little low to SM5100. The main reason is that the hardware of them is not on the same level, Godson-I CPU (only 26M Hz) frequency is far lower comparing to SM5100'processor, so throughput should be lower. But it still exerts a good performance.

#### 6. Conclusions

This paper presents a hw/sw co-verification methodology based on TWCNP-OS, a Linux-based operating system designed for FPGA-based platform of two-way cable network (TWCNP) SOC, and then particularly expatiates on the software and hardware coverification method. Through the implementation of TWCNP-OS on TWCNP and development on host/two targets, we realize the combination of design and verification. It is not a replacement of the co-verification with generic RTOS modeling, but is complementary to them.

The experiment makes it clear that the verification platform is simple-built and cost-effective. The scientific verification method shortens the development cycle and the validated system is stable and can be applied in the verification of other embedded system. Most importantly, since the introduction of TWCNP-OS to our FPGAbased platform, we have observed increased productivity among high-level application developers who have little experience in FPGA application design.

# 7. References

- [1] N. Ohba and K. Takano, "An SoC design methodology using FPGAs and embedded microprocessors," Proceedings of the 41st annual conference on Design automation, San Diego, CA, USA, June 07–11, 2004.
- [2] Habibi and S. Tahar, "Design and verification of SystemC transaction-level models," IEEE Trans. VLSI Syst., 14(1): pp. 57–68, January 2006.
- [3] H. So, A. Tkachenko, and R. Brodersen, "A Unified Hardware/Software Runtime Environment for FPGA-Based Reconfigurable Computers using BORPH," CODES+ISSS, 2006.
- [4] D. Atienza, P.G. Del Valle, G. Paci, et al., "A fast

HW/SW FPGA-based thermal emulation framework for multi-processor system-on-chip," Proceedings of the 43rd annual conference on Design automation, San Francisco, CA, USA, July 24–28, 2006.

- [5] Y. Nakamura, K. Hosokawa, I. Kuroda, K. et al., "A fast hardware/software co-verification method for system-ona-chip by using a C/C++ simulator and FPGA emulator with shared register communication," Proceedings of the 41st Design Automation Conference (DAC'04), San Diego, Calif., USA, pp. 299–304, June 2004.
- [6] D. Desmet, D. Verkest, and H. De Man, "Operating system based software generation for systems-on-chip," Proceedings of Design Automation Conference (DAC), 2000.
- [7] A. Gerstlauer, H. Yu, and D. Gajski, "RTOS modeling for system level design," Proceedings of Design Automation and Test in Europe (DATE), Embedded Software Forum, 2003.
- [8] SpecC Technology Open Consortium, http://www.specc.org/.
- [9] H. Tomiyama, Y. Cao, and K. Murakami, "Modeling fixedpriority preemptive multi-task systems in SpecC," Proceedings of Workshop on Synthesis and System Integration of Mixed Technologies (SASIMI), 2001.
- [10] Y. Yi, D. Kim, and S. Ha, "Virtual synchronization technique with OS modeling for fast and time-accurate cosimulation," Proceedings of International Conference on Hardware/Software Codesign and System Synthesis (CODES+ISSS), 2003.
- [11] F. Herrera, H. Posadas, P. Sanchez, and E. Villar, "Systematic embedded software generation from SystemC," Proceedings of Design Automation and Test in Europe (DATE), Embedded Software Forum, 2003.
- [12] S. Honda, T. Wakabayashi, H. Tomiyama, et al., "RTOScentric hardware/software cosimulator for embedded system design," Proceedings of the 2nd IEEE/ACM/IFIP international conference on Hardware/software codesign and system synthesis, Stockholm, Sweden, September 08–10, 2004.
- [13] ITRON, http://www.assoc.tron.org/itron/.
- [14] http://www.linux-mips.org/wiki/PMON.
- [15] http://cross-lfs.org/view/svn/mips/index.html.
- [16] http://linux.junsun.net/porting-howto/.
- [17] Mentor Graphics Corporation, http://www.mentor.com.
- [18] C. Wang, X.G. Xue, et al., "FPGA/CPLD Design TOOL-Specification for Xilinx ISE 5.X," Posts and Telecom Press, 2003.
- [19] S. Yoo, I. Bacivarov, A. Bouchima, et al., "Building fast and accurate SW simulation models based on hardware abstraction layer and simulation environment abstraction layer," Proceedings of the Design, Automation and Test in Europe (DATE'03), Munich, Germany, pp. 500–506, March 3–7, 2003.

# **Call for Papers**



# International Journal of

# Communications, Network and System Sciences (IJCNS)

ISSN 1913-3715 (Print) ISSN 1913-3723 (Online) Http://www.SRPublishing.org/journal/ijcns/

IJCNS is an international refereed journal dedicated to the latest advancement of communications and network technologies. The goal of this journal is to keep a record of the state-of-the-art research and promote the research work in these fast moving areas.

Guest Editors	
Prof. Tom Hou	Department of Electrical and Computer Engineering, Virginia Tech., USA

# Subject Coverage

This journal invites original research and review papers that address the following issues in wireless communications and networks. Topics of interest include, but are not limited to:

MIMO and OFDM technologies	Sensor networks
UWB technologies	Ad Hoc and mesh networks
Wave propagation and antenna design	Network protocol, QoS and congestion control
Signal processing and channel modeling	Efficient MAC and resource management protocols
Coding, detection and modulation	Simulation and optimization tools
3G and 4G technologies	Network Security

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 *IJCNS*.

# 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. For more details about the submissions, please access the website.

# Website and E-Mail

Http://www.SRPublishing.org/journal/ijcns

Jcnss@SRPublishing.org

# TABLE OF CONTENTS

Volume 1 June 20	008
Novel Joint Chip Sampling and Phase Synchronization Algorithm for Multistandard UMTS Systems	
Y. SERRESTOU, K. RAOOF, J. LIÉNARD	105
Performance Analysis of an AMC System with an Iterative V-BLAST Decoding Algorithm	
S.J. RYOO, K.W. LEE, I. HWANG	119
Impact of Depolarization Phenomena on Polarized MIMO Channel Performances	
N. PRAYONGPUN, K. RAOOF	124
A Quadratic Constraint Total Least-squares Algorithm for Hyperbolic Location	
K. YANG, J.P. AN, Z. XU	130
Analysis of Lifetime of Large Wireless Sensor Networks Based on Multiple Battery Levels	
R.H. ZHANG, Z.P JIA1, D.F. YUAN	136
Beacon-driven Leader Based Protocol over a GE Channel for MAC Layer Multicast Error Control Z. LI, T. HERFET	144
Routing and Wavelength Assignment in GMPLS-based 10 Gb/s Ethernet Long Haul Optical Networks with and without Linear Dispersion Constraints L.N. BINH.	154
On Modeling and Accuracy Analysis of the Available Bandwidth Measurement Based-on Packet-pair Sampling	
J. LIU, D.F. ZHANG, J.H. JIN	168
Streaming Multimedia over Wireless Mesh Networks D.Q. LIU, J. BAKER	177
A Novel Adaptive Hybrid Error Correction Scheme for Wireless DVB Services G.P. TAN, T. HERFET	187
A Co-verification Method Based on TWCNP-OS for Two-way Cable Network SOC C. LI, X.T. ZHANG, Y.D WAN, Q. WANG	199

Copyright©2008 SciRes

I. J. Communications, Network and System Sciences, 2008, 2, 105-206

