Communications and Network





www.scirp.org/journal/cn

Journal Editorial Board

ISSN: 1949-2421 (Print) ISSN: 1947-3826 (Online)

http://www.scirp.org/journal/cn

Editor-in-Chief

Dr. Yi HUANG

The University of Liverpool, UK

Editorial Board

Prof. Photios Anninos	Democritus University of Thrace, Greece
Prof. Ruay-Shiung Chang	National Dong Hwa University, Taiwan (China)
Dr. Wiani Jaikla	Suan Sunandha Rajabhat University, Thailand
Dr. Xiaohong JIANG	Tohoku University, Japan
Prof. Hussein Mouftah	University of Ottawa, Canada
Prof. Jean-Frederic Myoupo	University of Picardie-Jules Verne, France
Prof. Francesco Zirilli	Sapienza Universita di Roma, Italy



TABLE OF CONTENTS

Volume 2 Number 1

February 2010

Using the Power Control and Cooperative Communication for Energy Saving in Mobile Ad Ho	c Networks
C. J. Chen, C. Jin, D. M. Li, J. C. Wang, J. N. Fang	1
A Comparison Study of Input ESD Protection Schemes Utilizing NMOS, Thyristor, and Diode	Devices
J. Y. Chio	11
On Possible a-Priori "Imprinting" of General Relativity itself on the Performed Lense-Thirrir LAGEOS Satellites	ıg Tests with
L. Iorio	26
Incorporating Heterogeneous Biological Data Sources in Clustering Gene Expression Data	
T. Wibg	31
Joint Power Control and Spectrum Allocation for Cognitive Radio with QoS Constraint	
Z. J. Zhao, Z. Peng, Z. D. Zhao, S. L. Zheng	
A Historical Narrative of Study of Fiber Grating Solitons	
X. L. Li, Y. S. Jiang, L. J. Xu	44
ADPF Algorithm for Target Tracking in WSN	
C. H. Song, H. Zhao, W. Jing, D. Liu	50
Designing Intrusion Detection System for Web Documents Using Neural Network	
H. Om, T. K. Sarkar	54
On Solvable Potentials, Supersymmetry, and the One-Dimensional Hydrogen Atom	
R. P. Martínez-y-Romero, H. N. Núñez-Yépez, A. L. Salas-Brito	62
How to Measure in the Near Field and in the Far Field	
T. Dlugosz, H. Trzaska	65
Proposed Model for SIP Security Enhancement	
M. B. Sayyad, A. Chatterjee, S. L. Nalbalwar	69
A Model for Cu-Se Resonant Tunneling Diodes Fabricated by Negative Template Assisted Electrodeposition Technique	
M. Chaudhri, A. Vohra1, S. K. Chakarvarti	73
Live Video Services Using Fast Broadcasting Scheme	
S. Chand	79

Communications and Network (CN)

Journal Information

SUBSCRIPTIONS

Communications and Network (Online at Scientific Research Publishing, www.SciRP.org) is published quarterly by Scientific Research Publishing, Inc., USA.

E-mail: service@scirp.org

Subscription rates: Volume 2 2010

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

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

SERVICES

Advertisements Advertisement Sales Department, E-mail: service@scirp.org

Reprints (minimum quantity 100 copies)

Reprints Co-ordinator, Scientific Research Publishing, Inc., USA. E-mail: service@scirp.org

COPYRIGHT

Copyright© 2010 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: cn@scirp.org



Using the Power Control and Cooperative Communication for Energy Saving in Mobile Ad Hoc Networks

Chunjie Chen¹, Chao Jin¹, Demin Li¹, Jiacun Wang², Jianan Fang¹

¹College of Information Science and Technology, Donghua University, Shanghai, China ²Department of Computer Science and Software Engineering, Monmouth University, West Long Branch, NJ, USA E-mail: chunjiechen@mail.dhu.edu.cn, dazhilou@163.com, deminli@dhu.edu.cn, jwang@monmouth.edu Received October 10, 2009; accepted November 11, 2009

Abstract: In this paper, we investigate the energy saving problem in mobile ad hoc network, and give out an improved variable-range transmission power control algorithm based on minimum spanning tree algorithm (MST). Using previous work by Gomez and Campbell [1], we show that in consider of node's mobility, the previous variable-range transmission power control based on minimum spanning tree algorithm can not support nodes' mobility in mobile ad hoc network. For this reason, we give out an improved variable-range transmission power control algorithm to support node's mobility and solve asymmetric graph problem. To save more energy without changing the topology of the network, we give out two new data transmission mechanisms based on the idea of cooperative communication. The results of this paper enhance the possibility of using variable-range transmission power control in mobile ad hoc networks.

Keywords: minimum spanning tree, traffic capacity, energy savings, cooperative communication

1. Introduction

Energy saving problem is a very important issue in wireless ad hoc networks, because the transmission power impacts not only the connectivity but also the traffic capacity of the network. Choosing a higher transmission power can increase the connectivity and performance of the network, but reduce traffic capacity on the physical layer and energy on the network layer. Obviously, it's a trade-off problem. Today, the design of protocols for wireless ad hoc networks is primarily based on common-range transmission control, such as the work by Santi et al. [2,3]. But systems based on common-range transmission control [3] usually assume nodes are homogeneously distributed. For some nodes, the topology will be too sparse with the risk of having network partitions. For other nodes, the topology will be too dense, resulting in many nodes competing for transmission in a shared medium. This problem is discussed in [4], where the authors propose a method to control the transmission power levels in order to control the network topology. In [1], Gomez and Campbell show how variable-range transmission control can improve the overall network performance and suggests the design of MAC and routing protocols for wireless ad hoc networks should base on variable-range power control, not on common-range transmission control which is prevalent today. But in consider of node's mobility, using previous variablerange

power control still have some problems such as the traffic capacity tends to be zero because the area of overlapping region is zero. Furthermore, because in variable-range control, all nodes use different transmission powers, this also may lead to asymmetric graph problem.

In this paper we give out an improved power control algorithm based on minimum spanning tree to solve these problems and provide the possibility of using variable-range transmission power control in wireless ad hoc networks.

Furthermore, we investigate the effect of cooperative communication mechanism on energy saving for ad hoc networks. Many recent cooperative communication mechanisms change node's transmission power such as in [5] to save energy. But as Cardei said, using cooperative communication to control transmission power is an NP-Complete problem. Then we give out two new data transmission mechanisms based on the idea of cooperative communication. We prove these two mechanisms can solve the energy saving problem by reducing the transmission time rather than changing node's transmission power so it can be realized easily.

The structure of this paper is as follows. In Section 2 we give out an improved variable-range transmission power control algorithm to support node's mobility and solve asymmetric graph problem. Given the topology of the network, in Section 3 we will describe two new data transmission mechanisms Sequence Data Check Trans-

mission mechanism (SDCT) and Nearest Data Check Transmission mechanism (NDCT) based on our cooperative communication mechanism in detail to save more energy without changing the topology. To show the advantage on energy saving of our new data transmission mechanisms along with numerical simulation, in Section 4, we give out their mathematical model and some basic suppose. Through these mathematical models, we compute the average hop of traditional data transmission mechanism (TDT), SDCT and NDCT. Finally, we present some concluding remarks in Section 6.

2. Improved Variable-Range Power Control Algorithm

According to Gomez and Campbell [1], we find whether the algorithm can support node's mobility, depending on how large the area of overlapping region is. So can we enlarge the transmission power value in variable range control to a certain value? How to choose the bound of this value? How to deal with the asymmetric graph problem?

2.1 Traffic Capacity Using Previous Variable-Range Transmission Control in Mobile Network

In a mobile ad hoc network, nodes always move in a fast speed, this produces extra signaling overhead which consume a large part of network resources. In [1], Javier give out an equation to compute the signaling overhead of route maintenance. Javier suggested using previous variable-range transmission control based on minimum spanning tree lead to the time of a node remains in overlapping region b tends to zero. Figure 1 highlights one of overlapping regions.

It is obvious, because the network is built by minimum spanning tree, sol-hop nodes are always on the edge of their source nodes' coverage range and the area of the overlapping region b is zero. We prove this point by (15).

$$\lim_{h \to 0} T(R) = \frac{\pi R^2 \arccos\left(\frac{R-h}{R}\right) - \pi (R-h)\sqrt{2Rh-h^2}}{2\nu R \arccos\left(\frac{R-h}{R}\right)}$$
(1)
$$= \frac{\pi R}{2\nu} - \lim_{h \to 0} \frac{\pi (R-h)\sqrt{2Rh-h^2}}{2\nu R \arccos\left(\frac{R-h}{R}\right)} = 0$$

Because the average number of route-repair events persecond per route, proportional to $\frac{1}{T}$, J(R) is very large and because the traffic capacity of variable-range control $\lambda(R,t)$ tend to be a const when $n \to \infty$, so the capacity available to nodes for data transmission $\overline{\lambda}(R,t)$ is very small. This means the topology of the network is always changing and the whole wireless channel is occupied by



Figure 1. Overlapping region between two nodes



Figure 2. Topology change when node moves

the signaling overhead caused by route maintenance.

2.2 Improved Minimum Spanning Tree Algorithm for Upper Bound

We can see the key problem of previous variable-range transmission control is caused by route maintenance, so we must first discuss in which conditions the route of the network will be updated. Figure 2 shows that the change of topology in an ad hoc network whose nodes are distributed randomly. When node d moves in that direction shown by the arrow, the network will form a new minimum spanning tree shown by doted lines. 1) Forwarding nodes move out the region which is covered by transmitter node at the upper bound of transmission power.

2) Current structure of the network is no longer a minimum spanning tree and need to rebuild a new minimum spanning tree.

We can see the problem is to choose a suitable upper bound for node's transmission power.

By analyzing, we find when a node moves, the length of edges which connected to this node will change, and when this length beyond the shortest edge of its 1-hop nodes, which results in the nodes are not connected in original minimum spanning tree, the network needs to rebuild and the route needs to be updated.

Highlight from this point, we get the conclusion that the upper bound for a node's transmission power should be set to the value of the shortest edge of its 1-hop nodes which are not connected in original minimum spanning tree. But, there is an extreme condition that the upper bound may be much larger than original transmission power. To keep the superiority of minimum spanning tree algorithm, there must be a threshold to limit the upper bound. If we use $p_{i\min}$ to denote the minimum transmission power of node *i* in minimum spanning tree, to every node *i*, $\exists C_i$

$$C_i p_{i\min} \le \min\left\{p_{jk}\right\} \tag{2}$$

where *j* is 1-hop node of node *i* and *k* is 1-hop node of node *j*, p_{jk} is the transmission power between node *j* and node *k*. Because $p_{i\min}$ is the minimum transmission power of node *i* in minimum spanning tree. We set

$$C_i = \frac{\min\left\{p_{jk}\right\}}{p_{i\min}}$$

and there are two cases:

$$\begin{cases} C_i = 1 , \ k = j \\ C_i > 1 , \ k \neq j \end{cases}$$
(3)

where $C_i = 1, k = j$ means node *j* has only one 1-hop node, and this 1-hop node is node *i*.

Then we get the threshold

$$C = \min \left\{ \begin{array}{ccc} C_1 & C_2 & \cdots & C_n \end{array} \right\}$$
(4)

Finally we set the upper bound of the node's transmission power to:

$$p_{upper} = C p_{i\min} \tag{5}$$

where p_{upper} stands for the upper bound of transmission power according to the improved algorithm.

2.3 Adaptive Transmission Power Control for Lower Bound

To support node's mobility, we know node's transmission power should larger than that got by minimum spanning tree. But the upper bound $Cp_{i\min}$ may be much larger than $p_{i\min}$, considering that the mobility of mobile nodes is limited by physical restrictions, it is not necessary to enlarge the transmission power to $Cp_{i\min}$ in just one step. It can be a gradually increment process.

Considering that a mobile node's future location and velocity are likely to be correlated with its past and current location and velocity, according to Ben Liang and Zygmunt J. Haas [6], we use a 2-D Gauss-Markov mobility model to estimate node's future distance from its back-warding node.

In this 2-D Gauss-Markov mobility model, the velocity of a mobile is represented by the random vector:

$$\vec{V}_n = \vec{V} \left(nT_{APT} \right) = \left[V_n^x, V_n^y \right]^T \tag{6}$$

where T_{APT} is the time interval used to sense and compute. Define memory level :

$$\overline{\alpha} = \left[\alpha_n^x, \alpha_n^y\right]^T = \left[e^{-\beta_x T_{APT}}, e^{-\beta_y T_{APT}}\right]^T$$
(7)

where β_x, β_y are parameters correlated to T_{APT} .

Then, the 2-D velocity process can be expressed as follows:

$$\vec{V}_n = \overline{\alpha} \otimes \vec{V}_{n-1} + (1 - \overline{\alpha}) \otimes \overline{\mu} + \overline{\sigma} \otimes \sqrt{1 - \overline{\alpha}^2} \otimes W_{n-1}$$
(8)

where \otimes denotes element-by-element multiplication,

 $\overline{\sigma} = \left[\sigma_n^x, \sigma_n^y\right]^T$ and $\overline{\sigma}^2$ is the variance of \overline{V}_n , $\overline{\mu} = \left[\mu_n^x, \mu_n^y\right]^T$ is its mean. $\{W_n\}$ is an uncorrelated Gaussian process with zero mean and unit variance and is independent of $\{V_n\}$.

For simplicity of presentation, one may further assume that the velocity has the same memory level, the same asymptotic mean, and the same asymptotic standard deviation in both dimensions. In this isotropic case (8) becomes:

$$\vec{V}_{n} = \alpha \vec{V}_{n-1} + (1 - \alpha) \mu + \sigma \sqrt{1 - \alpha^{2}} W_{n-1}$$
(9)

If we use $\vec{r}_n = \vec{r} (nT_{APT}) = [r_n^x, r_n^y]^T$ to denote the distance from back-warding node to its 1-hop node, we get:

$$\begin{cases} \vec{r}_0 = r_{\min} \\ \vec{r}_n = \vec{r}_{n-1} + \vec{V}_n \end{cases}$$
(10)

Then we choose the increment of transmission range, denoted by τ_n shown in Figure 3, as the step of adaptive mechanism:

$$\tau_{n} = \|\vec{r}_{n}\| - \|\vec{r}_{n-1}\| = \sqrt{\left(r_{n-1}^{x} + V_{n}^{x}\right)^{2} + \left(r_{n-1}^{y} + V_{n}^{y}\right)^{2}} - \sqrt{r_{n-1}^{x^{2}} + r_{n-1}^{y^{2}}}$$
(11)

So within a slot of T_{APT} , 1-hop node is still in the transmission range of transmitter node to ensure network connectivity until the transmission range achieve the upper bound. When a node's transmission power approximates to its upper bound, it can broadcast to the whole network to ready for rebuild the minimum spanning tree. So the adaptive and heuristic algorithm can meet the networks switch for route and the requirement for QoS.

Now compare the transmission energy and traffic capacity between common-range control, previous variable-range control and improved variable-range control.

2.3.1 Transmission Energy

We know in common-range control, the transmission power, denoted by p_{COM} must meet the restriction:

$$p_{COM} \ge \max\{ p_{1\min}, p_{2\min}, \cdots, p_{n\min} \}$$
(12)

to make sure that every node is connected to the network.

Here we use $p_{i MST}$ to stand for the transmission power of improved variable-range control. By using adaptive

mechanism, we get

$$p_{i\min} + p_i(\tau) \le p_{i\text{ IMST}} = p_{i\min} + kp_i(\tau) \le p_{upper} = Cp_{i\min}$$
 (13)

where
$$p_i(\tau)$$
 is the extra transmission power of node *i* that needed to enlarge the transmission range of minimum spanning tree with one step.

The total transmission energy of one route in previous variable range control is

$$P_{\min} = \sum_{i=1}^{n} p_{i\min}$$
(14)

We can use (12) and (13) to compute the total transmission energy of one route in the network. We get the total transmission energy of common-range control

$$P_{COM} = \sum_{i=1}^{n} p_{COM} = n p_{COM}$$
(15)

and the total transmission energy of improved variablerange control

$$P_{IMST} = \sum_{i=1}^{n} p_{i \ IMST} \tag{16}$$

From (12), (13), (14), (15) and (16), compare the transmission energy between common-range control, previous variable-range control and improved variable-range control, we get

Theorem 2.1 The total transmission energy of im-

proved variable-range control is between common-range control and previous variable-range control: $P_{min} \leq P_{MST} \leq P_{COM}$.

$$P_{\min} = \sum_{i=1}^{n} p_{i\min} \leq \sum_{i=1}^{n} \left\{ p_{i\min} + p_i \left(\tau \right) \right\} \leq P_{IMST}$$

$$\leq C \sum_{i=1}^{n} p_{i\min} \leq \sum_{i=1}^{n} C_i p_{i\min} \leq \sum_{i=1}^{n} \max \left\{ p_{1\min} , p_{2\min} , \cdots, p_{n\min} \right\}$$

$$\leq n p_{COM} = P_{COM} \qquad (17)$$

2.3.2 Traffic Capacity of Improved Minimum Spanning Tree Algorithm

We can see, after enlarge the transmission power of minimum spanning tree to the upper bound of $Cp_{i\min}$, the parameter h is increased from 0 to (C-1)R, using (6) we deduce the signaling overhead of the improved variable-range control.

From (3) and (4) we know C is a constant larger than 1, so $J(R_{IMST})$ is a constant too. Compare to the signaling overhead of previous variable-range control which is tend to infinite, it's much smaller.

From (2), (4) and (12) we know

$$h_{\rm IMST} = (C-1)R_{\rm min} \le R_{\rm COM} - R_{\rm min} = h_{\rm COM}e \qquad (19)$$

and J(R) decreases as the parameter h increases.

According to the analysis and comparison above, we know the improved variable-range control can balance the transmission energy and signaling overhead in mobile ad hoc networks and it is an optimization of energy-saving and traffic capacity, it improves the mobility of the network greatly at the cost of a little more transmission power.

2.4 Asymmetric Graph

Different from common-range control, in variable-range control, all nodes use different transmission power, this may lead to the final structure of the network is an asymmetric graph, like Figure 3. To solve this problem, we select the sequence of transmission ranges from source node to the destination node to be an incremental sequence Here we use the method of "Incremental" which means if the transmission range of forwarding node is smaller than the source node, we use the transmission range of the sequence of transmission range of the source node to replace the transmission range of forwarding node. By this way, we ensure the sequence of transmission ranges from node to the destination node is monotone-up and the structure of the network in Figure 4 changes to Figure 5 by the monotone-up selecting method.

There are two extreme situations shown in Figure 6a and Figure 6b. In Figure 6a, the sequence of transmission ranges from source node to the destination node is monotone-up and In Figure 6b the sequence of transmission ranges from source node to the destination node is monotone-down. After "Incremental", we get the new transmission ranges shown by Figure 6c and Figure 6d. Then we get that:

L

Theorem 2.2 The total transmission power of the network after "Incremental" is between improved variablerange control and common-range control.

$$\left(\frac{2\nu R \arccos\left(\frac{R-(C-1)R}{R}\right)}{\pi R^{2} \arccos\left(\frac{R-(C-1)R}{R}\right)-\pi \left(R-(C-1)R\right)\sqrt{2R(C-1)R-(C-1)^{2}R^{2}}}\right)$$
(18)

Proof: If we use $\{p_{1-IMST}, p_{2-IMST}, \dots, p_{n-IMST}\}$ to stand for the sequence of transmission ranges from source node to the destination node after adaptive power control, then the sequence of transmission ranges from source node to the destination node after coverage is $\{p_1, p_2, \dots, p_n\}$,

In the condition of $p_{i \text{ IMST}} \le p_{i+1 \text{ IMST}}$ we get the transmission power of forwarding node is



Figure 3. Increment step of transmission range



Figure 4. Asymmetric structure of the network



Figure 5. Monotone-up transmission range

$$p_{i+1} = p_{i+1 \ IMST}$$
 (20)

And on the other hand, when the condition is $p_{i \ IMST} > p_{i+1 \ IMST}$ the transmission power for next hop will be

$$p_{i+1} = p_{i \ IMST} \tag{21}$$

And (21) stand for once change on the sequence of transmission ranges from source node to the destination node after adaptive power control. Because this change occurs when $p_{i \ IMST} > p_{i+1 \ IMST}$, so the total transmission power will increase after this change.

If original sequence is monotone-up, the number of change is zero which is the smallest. On the other hand, if original sequence is monotone-down, the number of chan-



Figure 6. (a) Condition of monotone-up; (b) Condition of monotone-down; (c) Evolution from Figure 6a after "Incremental"; (d) Evolution from Figure 6b after "Incremental"

ge is n-1 which is the largest. And if the original sequence is monotone-down, the final sequence of transmission power is

$$\left\{ p_{1 \ IMST} , p_{1 \ IMST} , \cdots , p_{1 \ IMST} \right\}$$

where

$$p_{1 IMST} = \max \left\{ p_{1 IMST}, p_{2 IMST}, \cdots, p_{n IMST} \right\}$$
 (22)

From (12) we can see $p_{1 IMST}$ is smaller than p_{COM}

When the original sequence is neither monotone-up nor monotone-down, the number of change is between zero and n-1, so the total transmission power after "Incremental", denoted by P_{cov} is:

$$P_{IMST} \le P_{COV} < P_{COM} \tag{23}$$

Thus, by enlarging the transmission range to the upper bound $p_{upper} = Cp_{i\min}$, adaptive transmission rang and the "Incremental" of transmission range, we have solved the traffic capacity tends and asymmetric graph problem of previous variable-range control mentioned above.

3. Cooperative Communication for Energy Saving

We all know that energy is the product of transmission power and transmission time, so the transmission energy is related to not only the transmission power but also the transmission time which is decided by the length of data. We usually choose changing transmission power to optimize the network, but changing transmission power also lead to the change of network topology. In [5], Cardei proves it's a NP-complete problem. So after get the topology of the network by improved variable-range control, the transmission power can not be changed further, one way to save more transmission power is to control the time of transmission. Here we use the cooperative communication mechanism because it can save transmission power without changing the topology of the networks (related to transmit time or length of data).

3.1 Cooperative Communication Mechanism

In [7], Agarwal introduced two parameters related with SNR: γ_p , which is the threshold needed to successfully decode the packet payload, and γ_{acq} , which is the threshold required for a successful time acquisition. In [5] Cardei assumes a packet received with a SNR γ , is: 1) fully received, if $\gamma_p \leq \gamma$; 2) partially received if $\gamma_{acq} \leq \gamma < \gamma_p$, and 3) unsuccessfully received, if $\gamma < \gamma_{acq}$. And in Cardei's cooperative communication model, consider the packet is fully received when

$$\sum_{k} \frac{p_k}{d_{kj}^{\alpha}} \ge 1 \tag{24}$$



Figure 7. Node c receive k packets in m packets from node a

where p_k is the transmission power of node k, d_{kj} is the distance between node k and node j, and α is a communication medium dependent parameter. But if there are overlapping parts of partial data, can (24) stand for full reception of the packet?

Based on this question, we give out our own cooperative communication mechanism.

We consider that messages are divided into several packets. If forwarding nodes are in the area of source node's transmission range, all packets can be transmitted to the 1-hop node of source node completely and correctly. But when forwarding nodes, such as 2-hop node and 3-hop node, are out of the transmission range of source node, only a part of packets can reach 2-hop node and 3-hop node, and among these partial data, some packets will fail because of distortion. We can see from Figure 7 if there are *m* packets to be transmitted from node a to node c, there will be k packets reach node cwhen node *a* transmits all these packets to node *b*. After validation, we can know which packets are correct. And when node b transmit data to node c, we do not need to retransmit these correct packets, and only transmit those packets which are not received or incorrect.

This mechanism can save more transmission energy without change the topology of network.

3.2 Mathematical Model

To simplify, we just consider the influence of source node on its 2-hop nodes.

Noting that wave transmits at the same speed in the same medium, from Figure 7, we know $\exists k \leq m$

$$kt_3 \le mt_1 \tag{25}$$

so

$$k = \left\lceil m \frac{t_1}{t_3} \right\rceil = \left\lceil m \frac{d_{ab}}{d_{ac}} \right\rceil$$
(26)

where t_1 and t_3 is the time need to transmit per packet from node a to node b and from node a to node c.

Though there will be k packets reach node c in Figure 7,

because of noise, some packets will fail because of distortion and how many packets can reach node *c* correctly is depend on the symbol-error-rate (SER) which is relate to the signal-to-noise ratio (SNR).

In [8], Ahmed and Weifeng Su give out a model of the relationship between symbol-error-rate (SER) and signal-to-noise ratio (SNR):

$$P_{SER}(m) = \sum_{i=0}^{2^{N}-1} \left\{ E_{CSI} \left[\Psi_{PSK} \left(SNR_{i} \right) \right] \prod_{k=1}^{N} E_{CSI} \left[P_{k,i}^{m} \right] \right\}$$
(27)

The model in [8], always choose the source node as the most reliable node, but in multi-hop network, the distance between source node and destination node is usually very long, this lead to the direct link between source node and destination node can not meet the requirement of high enough SNR, or there is even no direct link between source node and destination node. Another important point is that in this model, there are $2^{N} - 1$ possible network states, this means in most states, there are some nodes on the route are not in state and they do not relay the copy of data to other nodes. But as the nodes on the route form a chain, if there is one node do not relay the copy of data to its next node, the data can not reach the destination node correctly. So we should improve this model according to our improved routing strategy.

In our cooperative communication model, all packets can be transmitted to the 1-hop node of source node completely and correctly, so after the packets reached the 1-hop node, we can treat it as a new source node for the rest route. Compared to the source node, 1-hop node is closer to the destination node so the SNR of 1-hop is higher than the SNR of the source node. So there is only one network state in our cooperative communication model that all nodes are in state and relay the copy of data to other nodes.

So according to our model, we can simplify (27) to:

$$P_{SER}(1) = E_{CSI} \left[\Psi_{PSK} \left(SNR_d \right) \right] E_{CSI} P_{i,d}^1$$
(28)

where i is 1-hop node of the source node and d is the destination node. When just consider the influence of source node on its 2-hop node, the destination node is also the 2-hop node of the source node.

We assume that whether a packet can be received corrected is independent on other packets, and use function I(k) to denote whether the k th packet is received correctly.

$$I(k) = \begin{cases} 1, \ received \ correctly \\ 0, \ received \ not \ correctly \end{cases}$$
(29)

thus the number of packets E(n) that can be received correctly is

$$E(n) = \sum_{k} I(k) = k \left(1 - P_{SER}(1)\right)$$
(30)

So for every hop in the route, we can save the transmission energy used to transmit $k(1-P_{SER}(1))$ packets.

3.3 Two New Data Transmission Mechanisms Based on Cooperative Communication

After we prove the cooperative communication mechanism can save energy by shorting the time of transmission, we now give out two new data transmission mechanisms based on it.

3.3.1 Sequenced Data Check Transmission Mechanism

In this mechanism, the whole transmission process is sequenced according to the order of nodes from source to destination. The concrete process can be described as follow:

Step 1: The source node start the transmission process by sending route requiring signal (RRS) to the destination through forwarding nodes.

Step 2: Then the destination node return a route confirm signal (RCS).

Step 3: Source node will first send an inquire signal to its next hop to ask which packets does next hop need.

Step 4: After next hop return the packets ID which it has not received correctly, the source node begins to transmit these packets to next hop. Other nodes on the route receive these packets at the same time. They can know which packets are correct by checking.

Step 5: After all packets reach next hop return an ACK to the transmitter node and get ready to transmit data to its next hop.

Step 6: All forwarding nodes repeat Step 4-Step 5.

Step 7: If all packets reach the destination correct, the destination node will sent out a signal to acquire all nodes on the route to end this transmission process.

This mechanism is similar to recent point to point communication system, so it can be realized easily.

3.3.2 Nearest Data Check Transmission Mechanism

Though SDCT can save energy a lot, bet because the transmission process is ordered by nodes, the destination node can end the transmission process only after every node on the route received all data packets. Sometimes it is unnecessary and may lead to extra energy waste.

Thus we give out another data transmission mechanism:

Nearest Data Check Transmission (NDCT) to avoid this kind of energy waste. The main difference between SDCT and NDCT is: in NDCT, every node can send a route require signal to destination after it receive a correct packet and the destination node will choose the most reliable node (usually the nearest node) to relay. But in SDCT, for the k th relay, only the k th node on the route can send a route require signal to destination after it receive all packets correctly. The concrete process is as follow:

Step 1: Source node start the transmission process by sending route requiring signal with first packet ID to the destination through forwarding nodes.

Step 2: When the destination node receives several route require signal, it will choose the nearest node to return a route confirm signal.

Step 3: Node which receives the confirm signal will send one packet to its next hop.

Step 4: After next hop receive this packet correctly, it will send out a route require signal to ask for transmitting this packet. This route requiring signal will keep until the destination node reply a clear signal to start the transmission process of next packet.

Step 5: All forwarding nodes repeat Step 2-Step 4.

Step 6: If this packet reach the destination node correctly, the destination node will send out a clear signal to acquire all nodes on the route to clear their route requiring signal and turn to the transmission process of next packet.

This mechanism can save more energy than SDCT but it is much more complex.

4. Mathematical Models of SDCT and NDCT

Given the whole process of data transmission, now we can use Markov process to build mathematical models for SDCT and NDCT.

Here we consider a route with K nodes. The route status at moment n is defined as

$$S(n) = \left\{ s_1(n), s_2(n), \cdots, s_K(n) \right\}$$
(31)

here $s_i(n)$ is the status of node i at moment n and

$$s_i(n) = \begin{cases} 1 \text{ node } i \text{ has the correct packet by moment } n \\ 0 \text{ node } i \text{ has not the correct packet by moment } n \end{cases}$$

т раске l

(32)

Obviously, $S(0) = \{1, 0, \dots, 0\}$ and during the whole transmission process, there will be $N = 2^{K-1}$ possible statuses from $S^1 = \{1, 0, \dots, 0\}$ to $S^N = \{1, 1, \dots, 1\}$.

State transition probability matrix is:

$$P_{N\times N} = \begin{bmatrix} p_{11} & \cdots & p_{1N} \\ \vdots & \ddots & \vdots \\ p_{N1} & \cdots & p_{NN} \end{bmatrix}$$
(33)

here p_{ii} is the probability the route status transmit from S^i to S^j .

According to the feature of Markov process, we know

$$S(n) = S(0)P^n$$
 then $P(S(n) = S^i) = [P^n]_{1, i}$.

Here we first do some basic suppose: Suppose (1) The possibilities of every forwarding node receive a correct packet from transmitter node is independent of each other.

Then
$$p_{ij} = \prod_{m=1}^{K} p\left\{ s_m(n) = S_m^{j} \middle| s_1(n-1) = S_m^{i} \right\}$$
 (34)

Suppose (2)

$$p\left\{s_{i}\left(n\right)=1\middle|s_{i}\left(n-1\right)=0\right\}=p_{transmitter node, i}$$
(35)

where $p_{transmitter node, i}$ is the probability of node I receive a correct packet from the transmitter node. Then

$$p\left\{s_{i}\left(n\right)=0\middle|s_{i}\left(n-1\right)=0\right\}=1-p_{transmitter node, i}$$
(36)

$$p\left\{s_{i}\left(n\right)=1\middle|s_{i}\left(n-1\right)=1\right\}=1$$
(37)

$$p\left\{s_{i}(n)=0 \middle| s_{i}(n-1)=1\right\}=0$$
(38)

thus $[P_{1,i}^n] = f_{i,j}^n$ where $f_{i,j}^n$ is the probability that the route status transmit from S^1 to S^N for the first time after n steps.

Suppose (3) The probability of every transmitter node transmit a correct packet to its next hop is 1:

$$P\{S_{i+1}(n) = 1 | S_{i+1}(n-1) = 0, S_i(n-1) = 1\} = 1$$
(39)

So we can easily conclude the status of route will reach $S^{N} = \{1, 1, \dots, 1\}$ by most K - 1 steps.

Given these three basic suppose now we can compute the average hops of traditional data transmission (TDT), SDCT and NDCT.

4.1 Traditional Data Transmission Mechanism

In traditional data transmission mechanism

$$p_{transmitter \ node, \ i} = \begin{cases} 1 \ i \ is \ the \ next \ hop \ of \ transmitter \\ 0 \ i \ is \ not \ the \ next \ hop \ of \ transmitter \end{cases}$$

(40)

and the transmission process will be ended at S^N . Obviously

$$\left[P^{n}\right]_{1,n} = \begin{cases} 1 & n = K - 1 \\ 0 & n < K - 1 \end{cases}$$
(41)

Then the average hop of traditional data transmission mechanism is K - 1.

4.2 SDCT

In SDCT, the transmission process is sequenced by the order of nodes. So the transmitter can be denoted by

$$\min\{i | where \ s_i(n-1) = 1 \& s_{i+1}(n-1) = 0\}$$
(42)

Because the destination node can end the transmission process until every node on the route received all data packets. So the transmission process will be ended at S^N .

And

Then the average hop of SDCT is $\sum_{i=1}^{K-1} i \left[P_{SDCT}^{i} \right]_{1,n}$

Theorem3.1 The average hop of SDCT is less than TDT

Proof: Because

$$\sum_{i=1}^{K-1} \left[P_{SDCT}^{i} \right]_{1,n} = 1$$
 (43)

so we get

$$\sum_{i=1}^{K-1} i \left[P_{SDCT}^{i} \right]_{1,n} \le \left(K - 1 \right) \sum_{i=1}^{K-1} \left[P_{SDCT}^{i} \right]_{1,n} = K - 1 \quad (44)$$

Theorem 1 is proved.

4.3 NDCT

In NDCT, the destination node will choose the nearest node to which has the correct packet as the transmitter. So the transmitter can be denoted by

$$\max\left\{i\middle| where \ s_i\left(n-1\right)=1\right\}$$
(45)

and the destination will end current packet transmission process as soon as the correct packet reach the destination node. So the transmission process will be ended at every status

$$\left\{S^{j}\middle|where \ s_{k}\left(n-1\right)=1\right\}$$
(46)

Then the average hop of NDCT is

$$\sum_{i=1}^{K-1} i \sum_{S^j} \left[P^j_{NDCT} \right]_{1,S^j} = \sum_{i=1}^{K-1} i \sum_{k=1}^{\frac{1}{2}} \left[P^j_{NDCT} \right]_{1,2k}$$
(47)

Theorem3.2 The average hop of NDCT is less than SDCT

Proof: To compare the average hop of SDCT and NDCT, we should divide the transmission process into two cases:

1) The route status reaches S^N do not through $\{S^j\}$

2) The route status reaches S^N through $\{S^j\}$

In case 1), because both SDCT and NDCT end the transmission at S^N , so their average hop are the same.

In case 2), according to suppose (2) $\left[P_{1,i}^n\right] = f_{i,j}^n$

then we get

$$\sum_{i=1}^{K-1} i \sum_{s^{j}} \left[P_{NDCT}^{j} \right]_{1,s^{j}} = \sum_{m=1}^{K-1} \sum_{n=1}^{K-1-m} (m+n) f_{i,s^{j}}^{m} f_{s^{j},s^{N}}^{n}$$

$$\geq \sum_{m=1}^{K-1} (m+1) f_{i,s^{j}}^{m} \sum_{n=1}^{K-1-m} f_{s^{j},s^{N}}^{n}$$

$$= \sum_{m=1}^{K-1} (m+1) f_{i,s^{j}}^{m} > \sum_{i=1}^{K-1} i \left[P_{NDCT}^{i} \right]_{1,s^{j}}$$
(48)

Theorem 2 is proved.

Copyright © 2010 SciRes

5. Numerical Simulation

To verify our theory, we do some simulation on the average hop of TDT, SDCT and NDCT.

Here we use a route of 4 nodes, and the state transition probability matrix is as follow:

Table 1 shows the hops used to complete the transmission process of TDT, SDCT and NDCT and their probabilities.

1

0

Finally we get the average hop of TDT is 3, SDCT is 1.728 and NDCT is 1.648. Thus the advantage on energy saving of our new data transmission mechanism has been proved.

Table 1. Probabilities of hops used to complete the transmission process

	Probabilities of hops used to complete the trans- mission process			
	1-Hop	2-Hop	3-Нор	
NDCT	0.4	0.562	0.048	
SDCT	0.32	0.632	0.048	
TDT	0	0	1	

6. Conclusions

In this paper, we give out some effective solutions to improve previous variable-range control such as improved minimum spanning tree algorithm, adaptive power control mechanism and the "Incremental" of transmission power. We prove the improved variable-range control is an optimization of traffic capacity and energy saving. The variable-range control method improves the performance of the mobile ad hoc networks efficiently at the cost of a little more transmission power. Furthermore we describe two new data transmission mechanisms based on our cooperative communication mechanism in detail and show their advantage on energy saving with numerical simulation. By these improvements, we provide the possibility of using variable-range transmission power control in wireless ad hoc networks.

7. Acknowledgements

This work is partially supported by Shanghai Key Basic Research Project under grant number 09JC1400700; NSFC granted number 70271001, China Postdoctoral Fund granted number 2002032191, Shanghai Fund of Science and Technology granted number 00JG05047, the Shanghai Key Scientific Research Project under grant number 05dz05036, and the fund of Engineering Research Center of Digitized Textile & Fashion Technology, Ministry of Education.

REFERENCES

- J. Gomez and T. Campbell, "Variable-range transmission power control in wireless ad hoc networks," IEEE Transactions on Mobile Computing, Vol. 6, No. 1, pp. 87–99, 2007.
- [2] P. Santi, D. M. Blough, and F. Vainstein, "A probabilistic analysius for the range assignment problem in ad hoc networks," Proceeding ACM Mobihoc, pp. 212–220, August 2000.
- [3] MANET, IETF Mobile Ad-Hoc Network Working Group, 2006, http://www.ietf.org/html.charters/manet-charter.html.
- [4] R. Ramanathan and R. Rosales-Hain, "Topology control of multihop wireless network using transmit power adjustment," Proceeding IEEE INFOCOM, pp. 404–413, 2000.
- [5] M. Cardei, J. Wu, and S. Yang, "Topology control in ad hoc wirelsee networks using cooperative communication," IEEE Transactions on Mobile Computing, Vol. 5, No. 6, 2006.
- [6] B. Liang and Z. J. Haas, "Predictive distance-based mobility management for multidimensional PCS networks," IEEE/ACM Transactioons on Networking, Vol. 11, No. 5, pp. 718–732, October 2003.
- [7] M. Agarwal, J. H. Cho, and J. Wu, "Energy efficient broadcast in wireless ad hoc networks with hitch-hiking," Proceeding IEEE INFOCOM, 2004.
- [8] A. K. Sadek, W. F. Su, and K. J. R. Liu, "Multinode cooperative communications in wireless networks," IEEE Transactions on Signal Processing, pp. 341–355, 2007.



A Comparison Study of Input ESD Protection Schemes Utilizing NMOS, Thyristor, and Diode Devices

Jin Young Choi

Electronic and Electrical Engineering Departmet, Hongik University, Jochiwon, Korea E-mail: jychoi@hongik.ac.kr Received October 20, 2009; accepted December 22, 2009

Abstract: For three fundamental input-protection schemes suitable for high-frequency CMOS ICs, which utilize protection devices such as NMOS transistors, thyristors, and diodes, we attempt an in-depth comparison on HBM ESD robustness in terms of lattice heating inside protection devices and peak voltages developed across gate oxides in input buffers, based on DC, mixed-mode transient, and AC analyses utilizing a 2-dimensional device simulator. For this purpose, we construct an equivalent circuit model of input HBM test environments for CMOS chips equipped with input ESD protection circuits, which allows mixed-mode transient simulations for various HBM test modes. By executing mixed-mode simulations including up to six active protection devices in a circuit, we attempt a detailed analysis on the problems, which can occur in real tests. In the procedure, we suggest to a recipe to ease the bipolar trigger in the protection devices and figure out that oxide failure in internal circuits is determined by the peak voltage developed in the later stage of discharge, which corresponds to the junction breakdown voltage of the NMOS structure residing in the protection circuit for high-frequency ICs, and suggest valuable guidelines relating design of the protection devices and circuits.

Keywords: ESD protection, HBM, NMOS, thyristor, diode, mixed-mode

1. Introduction

CMOS chips are more vulnerable to electrostatic discharge (ESD) due to the thin gate oxides used, and therefore protection devices such as NMOS transistors are required at input pads. A large size for the protection devices is needed to reduce discharge current density and thereby to protect them against thermal-related problems. However, using the large devices adds parasitic capacitances to the input nodes to generate other problems such as gain reduction and poor noise characteristics in highfrequency ICs [1].

To reduce the added parasitics, various techniques have been suggested [1-3]. However, basic approaches should be to reduce the size of protection devices by utilizing, for example, thyristors or forward-biased diodes [4,5].

In this paper, we introduce three fundamental ESD protection schemes utilizing NMOS transistors, thyristors, and diodes, which can be implemented into input pad structures of high-frequency CMOS ICs, assuming usage of standard CMOS processes. While there can be many variants of the fundamental protection schemes, it is worthwhile to carefully examine the mechanisms leading

This work was supported by 2008 Hongik University Research Fund.

to device failures when using the fundamental protection schemes since it can provide valuable information in designing most of protection circuits. We analyze and compare in detail discharge characteristics of the three protection schemes for various discharge modes in input human-body model (HBM) tests. A 2-dimensional device simulator, together with a circuit simulator, is utilized as a tool for a comparative analysis. The analysis methodology utilizing a device simulator has been widely adopted with credibility [6,7] since it can provide valuable information relating the mechanisms leading to device failure, which may not be obtained by measurements.

In Section 2, we suggest three protection device structures, which will be utilized for the comparative analysis, and introduce device characteristics based on DC device simulations, which will be utilized to confirm the mixedmode simulation results analyzed in Section 4. In Section 3, we briefly explain discharge modes in HBM tests and introduce the input protection circuits utilizing each suggested protection device. In Section 4, we construct an equivalent circuit model of a CMOS chip equipped with input protection devices to simulate various input HBM test situations, and execute mixed-mode transient simulations on the circuits including up to six active protection devices. We figure out weak modes, and present in-depth

Table 1. Principal parameters of the NMOS device

analysis results on critical characteristics such as peak
voltages developed across gate oxides in input buffers,
locations of peak temperature inside protection devices,
and so on. In Section 5, we introduce AC device simula-
tion results to compare magnitudes of the added parasitics
when the suggested protection circuits are adopted. In
Section 6, considerations relating device design are dis-
cussed.

2. Protection Device Structures and DC Characteristics

Figure 1 shows the NMOS protection device structure assumed in this work. The scales of two axes are in micrometers. The structure represents a conventional protection device incorporating n⁺ source and drain ESD implants, which is implied by the relatively deep junctions. In order to alleviate drain-contact melting problems caused by lattice heating, the gate-drain contact spacing is chosen to be 3.5µm, which can be considered as ordinary. Table 1 summarizes the principal structure parameters. The n^+ and p^+ junctions shown in Figure 1 are assumed to have Gaussian doping profiles with about 10²⁰ cm⁻³ of peak concentration.

The p^+ junctions located at the upper left/right corners represent diffusions for substrate ground contacts. A series resistor of 1 M $\Omega\mu$ m, which is not shown in Figure 1, is connected at the bottom substrate node considering the distributed resistances leading to the substrate contacts located far away.

DC simulations were performed using a 2-dimensional device simulator ATLAS [8]. All necessary physical models including an impact ionization model were considered in the simulations. The latticeheating model included joule heat, generation-recombination heat, and Peltier-Thomson heat. The source, the gate, and the substrate were grounded, and the drain bias was varied for simulation.

Figure 2 shows the simulated drain current vs. voltage characteristics of the NMOS transistor in Figure 1 in a semi-log scale. We confirmed that a leakage current through the weakly inverted MOS channel dominates when the drain voltage is below 5V. Increasing the drain voltage, a leakage current through the reverse-biased n⁺-drain/p-sub junction starts to dominate, and the junction breakdowns by avalanche when the drain voltage is increased above 9.3V.

A generated hole current by avalanche flows to the substrate terminal to increase the body potential. With a sufficient hole current flowing, the body potential near the source junction gets high enough to forward-bias the n⁺-source/p-sub junction triggering a parasitic lateral npn (source/body/drain) bipolar transistor. The source, the body, and the drain act as an emitter, a base, and a collector, respectively. Generation of holes around the drain junction is augmented due to impact ionization caused by

Parameter	Values
Effective channel length	0.38µm
Gate oxide thickness	75µm
Substrate doping	10^{16}cm^{-3}
Channel peak doping	2.35×10 ¹⁷ cm ⁻³
Junction depth of n ⁺ diffusion	0.3µm
Junction depth of p^+ diffusion	0.1µm
Gate-drain contact spacing	3.5µm
Gate-source contact spacing	1.0µm

Table 2. Principal parameters of the lvtr thyristor device

Parameter	Values
p^+ & n^+ junction depth	0.1µm
n well depth	1.0µm
n^+ & p^+ anode contact spacing	2.7µm
NMOS effective channel length	0.38µm

the injected electrons from the source, and thereby the required drain-source voltage is reduced to show a snapback, as indicated as 'BJT trigger' in Figure 2. After the snapback at about 9.4V, the drain-source voltage drops to about 4.6V of a bipolar holding voltage.

In Figure 2, a 2nd breakdown [9] occurs when the drain current is about 1.3mA/µm, and the required drain-source voltage is further reduced to cause device failures relating drain-contact melting in real devices. It was confirmed that the 2nd breakdown in Figure 2 occurs when the peak lattice temperature inside the device exceeds about 1, 100°K.

Figure 3 shows the lvtr thyristor device structure assumed in this work. An lvtr thyristor device is a pnpntype device suggested to the lower snapback voltage by incorporating a NMOS transistor into it [4]. The device in Figure 3 can be easily fabricated in standard CMOS processes, and does not incorporate ESD implant steps, which is implied by the relatively shallow junctions.

Table 2 summarizes the principal structure parameters. The n well is assumed to have a Gaussian doping profile with 10^{17} cm⁻³ of peak concentration, and the doping profiles of the n^+ and p^+ junctions are similar to those of the p^+ junctions in Figure 1. A series resistor is also connected at the bottom substrate node as in the NMOS device in Figure 1. The n^+ and p^+ anodes in Figure 3 are tied together to serve as an anode. The cathode, the gate, and the substrate were grounded, and the anode bias was varied for simulation.

Figure 4 shows the simulated DC anode current vs. voltage characteristics of the lvtr thyristor device in Figure 3. As in the NMOS device, when the drain voltage is below 5V, a leakage current through the weakly inverted MOS channel between the n^+ well (the n^+ region at the



Figure 2. Drain I-V characteristics of the NMOS device

right-hand side of the n well) and the n^+ cathode dominates. When increasing the drain voltage, a leakage current through the reverse-biased n^+ -well/p-sub junction, where electric field intensity is highest, starts to dominate, and the junction breakdowns by avalanche when the drain voltage is increased above 8.8V. The n^+ well junction acts as a drain of the NMOS transistor, whose breakdown voltage is different from that of the NMOS



Figure 4. Anode I-V characteristics of the lvtr_thyristor device

device in Figure 1 due to the different junction-doping profile.

As the anode voltage increases, the p-sub/n⁺-cathode

junction is forward biased triggering a lateral npn $(n^+-cathode/p-sub/n^+-well)$ bipolar transistor. The n^+ cathode, the p substrate, and the n^+ well act as an emitter, a



Figure 5. Cross section of the diode device

base, and a collector, respectively. At this situation, a snapback is monitored as shown in Figure 4. The collector current from the n^+ anode flows through the n well to decrease the potential of the region under the p^+ anode by an ohmic drop. When the collector current is large enough, the p^+ -anode/n-well junction is forward biased to trigger a pnpn (p^+ -anode/n-well/p-sub/ n^+ -cathode) thyristor, which causes another decrease in the anode voltage, as indicated as 'pnpn trigger' in Figure 4. The resulting holding voltage drops to about 1V, which is much smaller compared to 4.6V of the NMOS transistor in Figure 2.

The 2nd breakdown in Figure 4 occurs when the anode current is about $12\text{mA}/\mu\text{m}$. The critical current for the 2nd breakdown is much larger than that in the NMOS device due to the reduced holding voltage, which implies superior ESD robustness of the lvtr_thyristor device in suppressing lattice heating.

Figure 5 shows the diode device structure assumed in this work, which is a p^+ -anode/n-well/n⁺-cathode junction. The doping profiles of the n well, the n⁺ and p⁺ junctions are similar to those in the lvtr_thyristor device, and the contact spacing between the n⁺ cathode and the p⁺ anode was chosen as 2.4 μ m. The reason for forming the p⁺n junction inside the n well is to use the same device as a protection device between V_{DD} and pad nodes as well as that between pad and V_{SS} nodes. A series resistor is also connected at the bottom substrate node as in the NMOS device in Figure 1.

The substrate and the p^+ anode were grounded and the n^+ cathode bias was biased positively or negatively to simulate DC reverse-bias or forward-bias characteristics, respectively. From the DC simulation results, it was confirmed that the forward diode drop is 0.95V when the diode current is 0.2mA/µm, and the reverse breakdown voltage is about 11.3V.

3. Input ESD Protection

Since parasitics added to an input pad by adopting ESD protection circuits should be minimized, it is desired to connect fewer number of protection devices to an input pad. An effective way to reduce the number is to use a V_{DD} - V_{SS} clamp device since it provides discharge paths without adding parasitics to an input pad. Figure 6, 7 and 8 show the fundamental ESD protection schemes utilizing the assumed three protection devices while minimizing the added parasitics. In the figure, a CMOS inverter was assumed as an input buffer.

The NMOS device shown in Figure 1 is used for M_1 and M_2 in Figure 6. M_1 is a protection device between the pad and V_{SS} nodes, and M_2 is a clamp device between the V_{DD} and V_{SS} nodes. It is important to locate all the protection devices close to the pad to minimize variation of the gate voltage in the input buffer when an ESD voltage is applied to the pad.

The lvtr_thyristor device shown in Figure 3 is used for T_1 in Figure 7. The NMOS device shown in Figure 1 is used for M_2 . In T_1 , the p⁺ and n⁺ anodes are connected to



Figure 6. Protection scheme utilizing the NMOS device



Figure 7. Protection scheme utilizing lvtr_thyristor device

the pad, and the p substrate and the n^+ cathode are connected to V_{SS} . Although it is not shown in Figure 7, the gate in T_1 is also connected to V_{SS} to maintain an off state in normal operations.

The diode device shown in Figure 5 is used for D_1 and D_2 in Figure 8. In D_1 , the n^+ cathode is connected to the pad, and the p^+ anode and the p substrate are connected to V_{SS} . In D_2 , the n^+ cathode, the p^+ anode, and the p substrate are connected to V_{DD} , the pad, and V_{SS} , respectively.

Since HBM tests for input pins should include all possible discharge modes, tests are performed for the five modes defined below.

1) PS mode: +V_{ESD} at an input pin with a $V_{SS}\xspace$ pin grounded

2) NS mode: -V $_{\text{ESD}}$ at an input pin with a V $_{\text{SS}}$ pin grounded

3) PD mode: $+V_{ESD}$ at an input pin with a V_{DD} pin grounded

4) ND mode: -V_{ESD} at an input pin with a V_{DD} pin grounded

5) PTP mode: $+V_{ESD}$ at one input pin with another input



Figure 8. Protection scheme utilizing the diode devices

pin grounded

Figure 9 shows main discharge paths in the protection scheme utilizing the NMOS device. In a PS mode, a parasitic npn bipolar transistor in M_1 provides a main discharge path, and in a NS mode, a forward-biased pn (p-sub/n⁺-drain) diode in M_1 provides it. In a PD mode, a parasitic npn bipolar transistor in M_1 and a forward-biased pn (p-sub/n⁺-drain) diode in M_2 in series provides a main discharge path, and in an ND mode, a parasitic npn bipolar transistor in M_2 and a forward-biased pn (p-sub/n⁺-drain) diode in M_1 in series provides it.

Local lattice heating is proportional to a product of current density and electric field intensity, and therefore temperature-related problems in the protection devices can occur in the parasitic npn bipolar transistor rather than in the forward-biased diode since the holding voltage of the bipolar transistor is much larger. Therefore we should assign sufficient device widths to M_1 considering PS and PD modes, and to M_2 considering an ND mode.

Main discharge paths in the protection scheme utilizing the lvtr_thyristor device are almost same as those shown in Figure 9 except that discharge paths inside T_1 replacethose in M_1 . That is, a pnpn thyristor in T_1 performs the role of the parasitic npn bipolar transistor in M_1 , and a forward-biased pn (p-sub/n⁺-anode) diode in T_1 performs the role of the forward-biased pn (p-sub/n⁺-drain) diode in M_1 .

Since lattice heating is not severe in pnpn thyristors by virtue of the smaller holding voltage, the width of the lvtr_thyristor device can be small. However, we should assign a sufficient device width to M₂ considering an ND mode.

Figure 10 shows main discharge paths in the protection scheme utilizing the diode devices. In a PS mode, forward-biased D_2 and an npn bipolar transistor in M_2 in series provides a main discharge path, and in an NS mode, forward-biased D_1 provides it. In a PD mode, forward-biased D_2 provides a main discharge path, and in an ND mode, an npn bipolar transistor in M_2 and forward-biased D_1 in series provides it.



Figure 9. Main discharge paths in the protection scheme utilizing the NMOS device



Figure 10. Main discharge paths in the protection scheme utilizing the diode devices

Since lattice heating is not severe in forward-biased diodes by virtue of the smaller voltage drop, the width of the diode devices can be small. However, we should assign a sufficient device width to M_2 considering PS and ND modes.

Figure 11 shows main discharge paths for a PTP mode. As shown on Figure 11(a), an npn bipolar transistor in M_1 and a forward-biased pn (p-sub/n⁺-drain) diode in M_3 in series provides a main discharge path in the protection scheme utilizing the NMOS device. It can be easily inferred that a pnpn thyristor in T_1 and a forward-biased pn (p-sub/n⁺-anode) diode in T_3 in series provides a main discharge path in the protection scheme utilizing the lvtr_thyristor device. As shown on Figure 11(b), two forward-biased pn (p⁺-anode/n⁺-cathode) diodes D₂, D₃ and an npn bipolar transistor in M_4 in series provides a main discharge path in the protection scheme utilizing the diode devices.



Figure 11. Main discharge paths for the PTP mode in the protection scheme utilizing (a) the NMOS device and (b) the diode devices

4. Mixed-Mode Transient Simulations

Figure 12 shows an equivalent circuit of an input HBM test situation assuming a PS mode. The portion indicated as 'Test environment' is an equivalent circuit for the test equipment connection. C_{ESD} and R_{ESD} represent a human capacitance and a human contact resistance, respectively, and 100pF and 1.5k Ω were assigned according to the international standard, respectively. C_s , C_t , and L_s , represent small parasitic elements present between test



Figure 12. Equivalent circuit of an input-pin HBM test situation



Figure 13. Variation of the drain current of M₁ in a PS mode in case of using the NMOS protection circuit



Figure 14. Variations of the voltages developed on C_{ngate} and C_{pgate} in a PS mode in case of using the NMOS protection circuit

equipment and an input pad, and typical values of 1pF, 10pF, and 5μ H [10] were assigned, respectively. V_{ESD} is a HBM test voltage, and the switch S₁ charges C_{ESD} and then the switch S₂ initiates discharge. By utilizing time-

varying resistors for the switches, the switching times of S_1 and S_2 were set short as 0.15ns.

In Figure 12, a V_{DD} - V_{SS} clamp NMOS device M_2 , protection devices P_1 and P_2 form a representative protection

circuit at an input pad. A CMOS inverter is assumed as an input buffer inside a chip, which is modeled by a capacitive network. C_{ngate} and C_{pgate} represent gate oxide capacitances of an NMOS transistor and a PMOS transistor, respectively. C_{ds} represents an n-well/p-sub junction capacitance. The reason for choosing this simple model for the inverter is based on the intension to minimize complexity of the equivalent circuit in this study focusing on the voltages developed across the gate oxides.

 R_{line} , R_{mdd} , and R_{mss} represent metal-line resistances, whose values were assigned relatively small as 5Ω assuming an input buffer located close to an input pad to simulate a more critical situation.

Using ATLAS, we performed mixed-mode transient simulations utilizing the equivalent circuit in Figure 12 equipped with one of the three input protection circuits shown in Figures 6–8. When a mixed-mode simulation is performed, active protection devices are solved by device and circuit simulations simultaneously. Notice that the number of the active protection devices included in a mixed-mode simulation in this work varies from two to six, which correspond to the PS, NS, PD, and ND mode simulations for the protection schemes in Figure 9, and the PTP mode simulation for the protection scheme util izing the diode devices in Figure 11(b), respectively.

For all the mixed-mode simulations performed for each test mode, $V_{ESD}=\pm 2000V$ was assumed. To make fair comparison on ESD robustness of the different protection schemes, the widths of the protection devices were adjusted to have utmost peak lattice temperature inside them below 500°K in all the mixed-mode simulations, resulting 250µm, 20µm, and 15µm for the NMOS device, the lvtr_thyristor device, and the diode device, respectively.

As an example of the mixed-mode simulation results, Figure 13 shows the variation of the M_1 drain current as a function of time in a PS mode in case of using the NMOS protection circuit in Figure 6. Notice that M_1 lies in the main discharge path in this case. The drain current reaches up to 2.2A, and shows a discharging characteristics with a time constant of roughly $R_{ESD}C_{ESD}=1.5k\Omega$ ×100pF=0.15µs.

Figure 14 shows the variations of the voltages developed on the capacitors C_{ngate} and C_{pgate} in the input buffer from the same simulation result. In Figure 14, the pad voltage is not shown since it is almost same with the voltage developed on C_{ngate} .

In Figure 13 and 14, we can see that the parasitic bipolar transistor in M_1 is triggered when the pad voltage in the early stage of discharge increases to about 11V, which is 0.68ns after S_2 in Figure 12 is closed. Main discharge through the parasitic bipolar transistor proceeds as the pad voltage, which is equal to the drain-source voltage of M_1 , drops to the holding voltage of about 5V.

We can also see that the pad voltage increases again

and reaches up to 9.5V at about 0.5 μ s, when the drain current is reduced below the holding current for the bipolar transistor action, and decreases very slowly thereafter. The peak voltage of 9.5V corresponds to the breakdown voltage of the NMOS device, which was explained relating Figure 2. The discharge thereafter continues a long time by the drain-junction leakage current in a breakdown mode. Up to 9.5V is developed on C_{ngate}, and in overall a lower voltage by about 1V is developed on C_{pgate} since the V_{DD} node does not lie in the main discharge path.

Since the discharge current decreases with time, the time for discharge to end is very long. We confirmed from additional simulations that it takes 7.4ms and 18.5 ms for the pad voltage to decrease down to 5V and 3V, respectively. If the pad voltage in the later stage of discharge is high, the NMOS gate oxide in the input buffer can be damaged since a large voltage is applied across it for a long time.

Figure 15 shows the variations of the voltages developed on C_{ngate} and C_{pgate} in a PS mode in case of using the lvtr_thyristor protection circuit in Figure 7. We confirmed that the pad voltage is almost same with the voltage developed on C_{ngate} again, and that the variation of the current through the anode of the lvtr_thyristor device T_1 is similar to that in Figure 13.

In Figure 15, the parasitic bipolar transistor in T_1 is triggered when the pad voltage in the early stage of discharge increases to about 12.8V, which is 0.77ns after S_2 in Figure 12 is closed. Main discharge through the pnpn thyristor proceeds as the pad voltage, which is equal to the anode-cathode voltage of T_1 , drops to the holding voltage of about 2V.

We can also see that the pad voltage increases again and reaches to 6.5V at about 0.9μ s, when the anode current is reduced below the holding current for the pnpn thyristor action, and decreases very slowly thereafter. We confirmed that main components of the anode current in this later stage of discharge are the leakage current through the n-well/p-sub junction and the weak-inversion MOS current. The developed voltage is smaller than the breakdown voltage (8.8V) of the lvtr_thyristor device shown in Figure 4. This seems to be caused by the longer duration (0.9μ s) of the main discharge through the pnpn thyristor, compared to that (0.5μ s) when using the NMOS protection scheme. The resulting discharge current in the later stage of discharge is too low for the lvtr_thyristor device to conduct in a breakdown mode.

We confirmed from additional simulations that it takes 165ms and 510ms for the pad voltage to decrease down to 5V and 3V, respectively.

In this case also, in overall a lower voltage by about 1V is developed on C_{pgate} since the V_{DD} node does not lie in the main discharge path.

Figure 16 shows the variations of the voltages developed on C_{ngate} and C_{pgate} as a function of time in a PS



Figure 15. Variations of the voltages developed on C_{ngate} and C_{pgate} in a PS mode in case of using the lvtr_thyristor protection circuit



Figure 16. Variations of the voltages developed on C_{ngate} and C_{pgate} in a PS mode in case of using the diode protection circuit

mode in case of using the diode protection circuit in Figure 8. We confirmed that variation of the anode current of D_2 , which lies in the main discharge path, is similar to that in Figure 13.

A forward-biased diode in D_2 is triggered when the pad voltage in the early stage of discharge increases to about

13.4V, which is 0.82ns after S_2 is closed. At this point the voltage developed across D_2 in Figure 8 corresponds to about 7.6V. Main discharge through the forward-biased diode in D_2 and the parasitic bipolar transistor in M_2 in series proceeds when the pad voltage drops to a sum of the forward diode drop and the holding voltage, which

We can see that the pad voltage increases again and reaches to 10.7V at about 0.5 μ s when the drain current of the clamp NMOS device is reduced below the holding current for the bipolar transistor action. The peak voltage 10.7V corresponds to a sum of the forward diode drop in D₂ (1.2V) and the breakdown voltage of M₂ (9.5V). Therefore, the maximum voltage developed on C_{ngate} is larger by an amount of the forward diode drop in D₂ when compared with that in case of using the NMOS protection circuit.

In Figure 16, the voltage developed on C_{pgate} is maintained low all the time since it is almost equal to the forward diode drop in D_2 .

4.1 Voltages across the Gate Oxides in the Early Stage of Discharge

In case of the PS mode analyzed up to this point, the trigger times for the parasitic bipolar transistor in M_1 , the parasitic bipolar transistor in T_1 , and the forward diode in D_2 are 0.68ns, 0.77ns, and 0.82ns, respectively, which are relatively short without big differences. However, due to these times, voltages larger than the snapback voltage or the ordinary forward diode drop appear across the devices right after S_2 is closed, resulting the high voltages developed on C_{ngate} in the early stage of discharge in Figures 14, 15 and 16.

Depending on test modes, larger peak voltages across the gate oxides appear at C_{ngate} or C_{pgate} in the early stage of discharge. If we define the test modes, which produce larger peak voltages in the mixed-mode transient simulations performed for 5 test modes, as weak modes, the results can be summarized as shown in Table 3.

The peak voltages in Table 3 could be regarded as excessive; however, the durations of the peak voltages are very short. We confirmed that, for example, the durations for which the voltages exceed 10V are at most 0.3ns. Therefore it may be inferred that the gate oxides won't be damaged in the early stage of discharge [11].

Notice that the peak voltages can be suppressed by reducing the bipolar trigger voltage of the NMOS protection device. To make the bipolar trigger voltages even lower than the off-state DC breakdown voltages, the gate-coupled NMOS (gcNMOS) structure [12] can be adopted. It is based on the technique to adding a RC network to the gate node, which is composed of a coupling capacitor (C_C) connected between the gate and the drain nodes to turn on the NMOS transistor immediately after a positive ESD pulse is applied to the drain, and a resistor (R_G) connected between the gate and V_{SS} nodes to discharge the gate node thereafter. The on duration is defined by R_GC_C . Turning on the NMOS transistor reduces the bipolar trigger voltage with enhanced hole generation at a lower drain-source voltage.

It seems possible to obtain a similar result by simply

Protection scheme	Week mode	Peak ve	oltage (V)	Time
Trottetion seneme	weak moue	C _{ngate}	C _{pgate}	(ns)
NMOS	PS	11.0		0.68
	PD		11.9	0.62
	ND		11.8	0.62
Lvtr_thyristor	PS	12.8		0.77
	PD		13.3	0.66
	ND		13.5	0.83
Diode	PS	13.4		0.82
	ND		13.3	0.82

adding a series resistor between the gate and V_{SS} nodes since the gate-drain overlap capacitance (C_{gd}) already exists in the NMOS structure, avoiding an increase of added parasitic to the input node. For the lvtr_thyristor device, the same technique may be applied since it includes the same NMOS structure in it.

We performed addition simulations to confirm that the early peaking can be suppressed by adding the series resistor to the gate node. For the 250µm NMOS device, adding a 10k Ω resistor between the gate and V_{SS} nodes is enough to turn on the NMOS for about 5ns duration suppressing the peak voltage on C_{ngate} down to 8.7V in case of the PS mode when using the NMOS protection scheme. The gate voltage peaks around 1.45V at 0.5ns, and the bipolar trigger time is also reduced to 0.5ns in accordance. Adding higher than a 50k Ω resistor turns on the NMOS transistor for an excessive duration more than 30ns, and tends to exacerbate lattice heating by confining the main discharge current towards the surface for a longer time.

When using the diode protection scheme, the same recipe on the NMOS clamp device suppresses the peak voltage down to 10.1V in the PS mode, which is still high due to the needed trigger voltage for the diode but lower than that (10.65V at 0.5μ s) in the later stage of discharge.

In case of the lvtr_thyristor device, there exist an n-well resistance (R_{NW}) between the n⁺ anode and the n⁺well junction, which tends to reduce the peak voltage developed at the gate node by the resistive division. For the 20µm lvtr_thyristor device, a 125k Ω resistor connected between the gate and V_{SS} nodes, which is larger by the same ratio (12.5) of the device sizes, suppresses the peak voltage down to 8.7V. The gate voltage peaks only 1.1V at 0.65ns, however it is still enough to utilize the recipe.

4.2 Voltages across the Gate Oxides in the Later Stage of Discharge

Depending on test modes, larger peak voltages across the gate oxides also appears at C_{ngate} or C_{pgate} in the later

Protection scheme	Week mode	Peak voltage (V)		
1 rotection scheme	weak mode -	C _{ngate}	C _{pgate}	
NMOS	PD	9.6	10.4	
	ND		10.4	
Lvtr_thyristor	ND		10.7	
Diode	PS	10.7		
	ND		10.7	
	PTP	10.8	10.8	

 Table 4. Peak voltage developed across the gate oxides in the later stage of discharge

stage of discharge. If we define the test modes, which produce larger peak voltages, as weak modes, the results can be summarized as shown in Table 4.

As explained relating the results in the PS mode, the high pad voltages in the later stage of discharge can damage the gate oxides since they last for long time.

In case of using the NMOS protection scheme in Table 4, the developed voltage on C_{pgate} in a PD mode is larger than that on C_{ngate} since the forward diode drop in M_2 is added to the breakdown voltage of M_1 , which can be easily explained from Figure 6 and 9. Due to the same reason, 10.4V is developed on C_{pgate} in an ND mode. In case of using the lvtr_thyristor protection scheme, 10.7V on C_{pgate} in an ND mode corresponds to a sum of the breakdown voltage of M2 and the forward diode drop in T_1 . The voltage is somewhat larger than that in the NMOS protection scheme since the diode drop in the lvtr thyristor device is larger due to the smaller device width adopted. In case of using the diode protection scheme, the same voltage (10.7V) is developed on C_{ngate} in a PS mode and on C_{pgate} in an ND mode, and this voltage corresponds to a sum of the breakdown voltage of M_2 and the forward diode drop in D_1 or D_2 . In a PTP mode, 10.8V is developed both on Cngate and Cpgate, which corresponds to a sum of the breakdown voltage of M₄ and the forward diode drop in D_2 or D_3 , which can be easily explained from Figures 8 and 11(b).

When judging from the peak voltages developed across the gate oxides in the later stage of discharge in Table 4, the weakest modes in case of using the NMOS protection scheme are PD and ND modes, and the PMOS gate oxide is more vulnerable to HBM ESD damages if the gate-oxide thicknesses of the NMOS and the PMOS are same. In case of using the lvtr_thyristor protection scheme, the weakest mode is an ND mode and the PMOS gate oxide is more vulnerable. In case of using the diode protection scheme, the weakest mode is a PTP mode and the NMOS and PMOS gate oxides are vulnerable in the same extent.

In Table 4, we can see that there is no big difference in the peak voltages developed across the gate oxides in the input buffers in each protection scheme. This is because the peak voltages in the later stage of discharge are determined mainly by the junction breakdown voltage of the NMOS structure in the NMOS devices or the lvtr_thyristor device. Since the breakdown voltage cannot be lowered with the gate coupling technique, junction engineering is essential to reduce it and to avoid possible oxide failures. We note that any junction engineering to lower the breakdown voltage was not tried in this work.

4.3 Location of Peak Temperature

As explained in Section 3, depending on test modes, utmost peak temperature resulting from lattice heating appears at the protection device connected to the input pad or at the V_{DD} - V_{SS} clamp NMOS device. In case of using the NMOS protection circuit, we confirmed that the utmost peak temperature in a PS mode appears at M₁, which lies in the main discharge path, and Figure 17 shows the variation of the peak temperature inside M₁. The peak temperature increases up to 495°K at about 30ns, when the bipolar transistor current still dominates the discharge, and decreases slowly as the discharge current decreases. By examining 2-dimensional temperature appears at the gate-side n⁺ drain junction.

In case of using the lvtr thyristor protection circuit, the utmost peak temperature in a PS mode appears at T_1 , which lies in the main discharge path. The peak temperature inside T_1 , whose device width is set to 20 μ m, increases sharply up to 473°K at 0.9ns and decreases down to 330°K as the pnpn thyristor in T_1 is triggered, and increases again up to 421°K at about 47ns, when the pnpn thyristor current still dominates the discharge, and decreases slowly. We confirmed that the peak temperature at 0.9ns appears at the n^+ well junction, where the electric field is high in a breakdown mode, and that at 47ns appears at the n⁺ cathode junction, where the current density is high. Notice that the peaking at 0.9ns can be avoided by adopting the gate coupling technique to reduce the bipolar trigger voltage. With a 125k Ω resistor connected at the gate, the early peaking is reduced down to 370°K providing a room for reducing the device size by a small amount.

In case of using the diode protection circuit, the forward-biased D_2 and the npn bipolar transistor in M_2 in series provides a main discharge path in a PS mode, and the utmost peak temperature appears at M_2 . The peak temperature inside M_2 increases up to 495°K at about 30ns, when the bipolar transistor current in M_2 still dominates the discharge, and decreases slowly. We confirmed that the peak temperature appears again at the gate-side n⁺ drain junction. We also confirmed that the peak temperature inside D_2 , whose device width is set to $15\mu m$, increases up to 485°K at about 45ns, and appears at the n⁺ cathode junction.

If we define the test modes, which produce larger tem-



Figure 17. Peak temperature variation inside the NMOS device (M₁) in a PS mode in case of using the NMOS protection circuit

Table 5. Peak temperature locations and times

Protection Weak		Peak	Peak temperature		
scheme	mode	temp. (°K)	Location	Time (ns)	
NMOS	PS, PD, PTP	495	Gate-side drain junction in M ₁	32	
	ND	495	Gate-side drain junction in M_2	31	
Lvtr_thyristor	PS	473	n^+ well junction in T_1	0.9	
		421	$n^{\scriptscriptstyle +}$ cathode junction in T_1	47	
	ND	495	Gate-side drain junction in M_2	33	
Diode	All	485	$\boldsymbol{n}^{\scriptscriptstyle +}$ cathode junction in \boldsymbol{D}_1 or \boldsymbol{D}_2	43~48	
	PS, ND	495	Gate-side drain junction in M_2	32, 34	

perature increase inside any protection device, as weak modes, the results can be summarized as shown in Table 5.

The peak temperature in the NMOS device, which is commonly used in all of the three protection circuits, appears at the gate-side n^+ drain junction. This is the reason for assigning a large spacing between the gate and the drain contact in Figure 1 to avoid drain contact melting.

In case of using the lvtr_thyristor protection scheme, the peak temperature in T_1 appears at the n^+ well junction even though it can be avoided with the gate coupling technique. However, a problem with contact melting will not occur in this junction since there is no contact on it. The 2nd peak temperature in T_1 appears at the n^+ cathode junction, and junction engineering such as increasing the junction area or adopting ESD ion implantation may be required to restrain temperature increase. However, it will not add parasitics to the input pad since the junction is not connected to it.

In case of using the diode protection scheme, the peak temperature in D_1 or D_2 appears at the n⁺ cathode junction, and similar junction engineering may be required to restrain temperature increase. However, it will not add parasitics to the input pad unless the n-well size is increased since the junction stays inside the n-well.

5. AC Device Simulations

We performed AC device simulations using ATLAS to compare magnitudes of the parasitics added to an input pad when using three different protection schemes in Figures 6–8.

Since only the drain is connected to an input pad when using the NMOS device in Figure 1, all nodes except the drain were grounded and an AC voltage was applied to the drain for a simulation to get admittances of the device as a function of frequency. In case of the lvtr_thyristor device in Figure 3, all nodes except the n^+ anode and the p^+ anode were grounded and an AC voltage was applied to the anode. To get admittances of the diode device (D₁) connected between the pad and the ground, all nodes except the n^+ cathode were grounded and an AC voltage was applied to the n^+ cathode. In case of the diode device (D₂) connected between V_{DD} and the pad, all nodes except the p^+ anode were grounded and an AC voltage was applied to the p^+ anode. The DC voltages for all nodes were assumed to be zero to simplify the analysis based on comparison.

Table 6. Series R and C parasitics of the protection dev
--

Protection device	C [F/µm]	R [Ω·μm]
NMOS	4.45×10 ⁻¹⁵	1.0×10 ³
lvtr_thyristor	3.10×10 ⁻¹⁵	1.4×10^{4}
Diode (D ₁)	2.35×10 ⁻¹⁵	3.5×10 ⁵
Diode (D ₂)	0.97×10 ⁻¹⁵	4.0×10 ³
Diode (total)	3.32×10 ⁻¹⁵	1.7×10 ⁵

 Table 7. Parasitics added to the input node in each protection scheme

Protection scheme	C [F]	R [Ω]
NMOS (250µm)	1.11×10 ⁻¹²	4
lvtr_thyristor (20µm)	6.20×10 ⁻¹⁴	700
Diode (15µm)	4.98×10 ⁻¹⁴	1.1×10^{4}

Simple series RC circuits seem adequate as the AC equivalent circuits for the protection devices to roughly compare magnitudes of the added parasitics [13], and Table 6 summarizes the R and C values extracted by fitting the modeled admittances assuming series RC equivalent circuits to those by the AC device simulations. In Table 6, the diode (total) device denotes the parallel combination of D_1 and D_2 .

Let's focus on the capacitance values in Table 6. Main portion of the capacitance in the NMOS device is the n⁺-drain/p-sub junction capacitance, whose value is relatively large since the n⁺ junction is large and deep as shown in Figure 1. The main portion of the capacitance in the lvtr_thyristor device is a parallel sum of the n-well/psub junction capacitance and the n⁺ well/p-sub junction capacitance. While main portion of the capacitance in D₁ is a parallel sum of the n-well/p-sub junction capacitance and the junction capacitance relating the p⁺ anode, that of the capacitance in D₂ is the p⁺-anode/n-well junction capacitance alone. Fore this reason, the capacitance in D₁ is larger than that in D₂.

Table 7 summarizes the parasitics added to an input pad, computed from the simulated parasitics in Table 6 by considering the device widths in each protection scheme, which are shown in the parentheses. From Table 7, we can see that when using the protection circuit utilizing the lvtr_thyristor device or the diode device, the added parasitic capacitance to an input pad can be reduced to 1/18 or 1/22 of that when using the NMOS protection circuit, respectively, while providing a similar level of ESD robustness in terms of lattice heating. Therefore we can confirm that the lvtr_thyristor protection scheme and the diode protection scheme are much superior to the NMOS protection scheme if they are adopted as an input protection scheme in high-frequency ICs, for example, in RF ICs.

6. Design Considerations

6.1 Considerations in Designing the NMOS Device

By performing additional simulations, we figured out that a serious problem could occur when p-type substrate contacts are not located close to the NMOS device for the reason explained below.

When the p-sub/n⁺-drain forward diode in M_1 or M_2 in Figure 9 gets on in the early stage of discharge in PD, ND, and NS modes, the n⁺-source/p-sub junction is excessively reverse biased due to an ohmic drop inside the p-type substrate if substrate contacts are not located close. In that case, a parasitic npn (n⁺-drain/p-sub/n⁺-source) bipolar transistor inside the NMOS deice can be triggered to increase temperature around the n⁺ source junction a lot, where electric field intensity is high. By the same mechanism, in PD and ND modes, a sum of the bipolar holding voltages of M_1 and M_2 , which is about 12V, is developed on C_{pgate} in a significant duration, which may damage the gate oxide. Therefore it is very important to locate the p-sub contacts close as shown in Figure 1.

6.2 Considerations in Designing the Lvtr_Thyristor Device

When p-type substrate contacts are not located close to the lvtr_thyristor device, the same problem with that in the NMOS device can occur in the lvtr_thyristor device. When the p-sub/n⁺-anode forward diode in T₁ in Figure 7 gets on in the early stage of discharge in NS and ND modes, a parasitic npn (n⁺-anode/p-sub/n⁺-cathode) bipolar transistor inside the lvtr_thyristor deice can be triggered to increase temperature around the n⁺ cathode junction a lot, where electric field intensity is high. By the same mechanism, in an ND mode, a sum of the bipolar holding voltages of T₁ and M₂ is developed on C_{pgate} in a significant duration, which may damage the gate oxide. Therefore it is also very important to locate the p-sub contacts close as shown in Figure 3.

6.3 Considerations in Designing the Diode Device

The diode device in Figure 5 does not have p-type substrate contacts close to it. By performing additional simulations, we figured out that a serious problem can occur in a PS mode if p-type substrate contacts are located close.

Let's assume that we attempt a PS mode test with an additional grounded p^+ -sub contact located at the upper right-hand side corner of the diode device in Figure 5. When the p^+ -anode/ n^+ -cathode diode in D_2 in Figure 10 gets on, a lateral parasitic pnp (p^+ - anode/n-well/ right-hand side p^+ -sub) bipolar transistor inside D_2 can be triggered to allow a large current and to increase temperature around the additional p^+ -sub contact a lot. We

confirmed that even with the proposed diode structure in Figure 5, a vertical pnp (p^+ -anode/n-well/p-sub) in D₂ is triggered. However, due to a resistance leading to p-sub contacts, amount of the bipolar current is restrained not to cause a temperature-related problem. Therefore it is very important to locate p-sub contacts as far away as possible.

6.4 Location of the Clamp NMOS Device

Since the clamp devices M_2 in Figures 6–8 are large and consume a large area if they are located in every input pad, we may consider locating them between V_{DD} and V_{SS} buses in V_{DD} and/or V_{SS} pad structures. Although a clamp device M_2 in that case can provide the same discharge paths explained, the ohmic voltage drops in the V_{DD} and/or V_{SS} bus with a very large discharge current flowing will increase the developed voltages across the gate oxides in input buffers, especially in case of adopting the diode protection circuit in Figure 8 since the ohmic voltage drops occur in both buses. Therefore it is recommended to locate M_2 in each input pad structure unless the chip size is not a critical issue.

7. Summary

For three fundamental input-protection schemes suitable for high-frequency CMOS ICs, which utilize protection devices such as NMOS transistors, thyristors, and diodes, we attempted an in-depth comparison on HBM ESD characteristics based on DC, mixed-mode transient, and AC analyses utilizing a 2-dimensional device simulator.

For this purpose, we construct an equivalent circuit model of input HBM test environments for CMOS chips equipped with input ESD protection circuits, which allows mixed-mode transient simulations for various HBM test modes. By executing mixed-mode simulations including up to six active protection devices in a circuit and analyzing the results, we attempted a detailed analysis on the problems, which can occur in real tests. Contributions of this work can be summarized as follows.

1) We demonstrated a simulation-based method to analyze problems occurring in all possible input HBM ESD test modes.

2) We figured out weak modes in terms of the peak voltages developed across gate oxides in input buffers in each protection scheme. We showed that the voltage peaking in the early stage of discharge can be suppressed by simply adding a series resistor to the NMOS gate, and figured out that oxide failure is determined by the peak voltage developed in the later stage of discharge, which corresponds to the junction breakdown voltage of the NMOS structure residing in the protection devices.

3) We figured out weak modes in terms of temperature increase inside the protection devices in each protection scheme, and also figured out the locations of peak temperature inside the protection devices.

4) We compared magnitudes of the added parasitics to an input pad in each protection scheme to confirm that the lvtr_thyristor and the diode protection schemes are more suitable for highfrequency ICs.

5) We also suggested the valuable design guidelines to minimize temperature increase inside the protection devices and to minimize the voltages developed across the gate oxides in input buffers.

REFERENCES

- [1] P. Leroux and M. Steyaert, "High-performance 5.2GHz LNA with on-chip inductor to provide ESD protection," Electronics Letters, Vol. 37, pp. 467–469, March 2001.
- [2] B. Kleveland, T. J. Maloney, I. Morgan, L. Madden, T. H. Lee, and S. S. Wong, "Distributed ESD protection for high-speed integrated circuits," IEEE Transactions on Electron Devices, Vol. 21, pp. 390–392, August 2000.
- [3] S. Hyvonen, S. Joshi, and E. Rosenbaum, "Cancellation technique to provide ESD protection for multi-GHz RF inputs," Electronic Letters, Vol. 39, No. 3, pp. 284–286, February 2003.
- [4] A. Chatterjee and T. Polgreen, "A low-voltage triggering SCR for on-chip ESD protection at output and input pads," IEEE Electron Device Letters, Vol. 12, pp. 21–22, August 1991.
- [5] E. R. Worley, R. Gupta, B. Jones, R. Kjar, C. Nguyen, and M. Tennyson, "Sub-micron chip ESD protection schemes which avoid avalanching junctions," in Processing, EOS/ ESD Symposium, pp. 13–20, 1995.
- [6] H. Feng, G. Chen, R. Zhan, Q. Wu, X. Guan, H. Xie, and A. Z. H. Wang, "A mixed-mode ESD protection circuit simulation-design methodology," IEEE Journal Soilid-State Circuits, Vol. 38, pp. 995–1006, June 2003.
- [7] B. Fankhauser and B. Deutschmann, "Using device simulations to optimize ESD protection circuits", in Processing, IEEE EMC Symposium, pp. 963–968, 2004.
- [8] ATLAS II Framework, Version 5.10.2.R, Silvaco International, 2005.
- [9] A. Amerasekera, L. van Roozendaal, J. Bruines, and F. Kuper, "Characterization and modeling of second breakdown in nMOST's for extraction and ESD-related process and design parameters," IEEE Transactions on Electron Devices, Vol. 38, pp. 2161–2168, September 1991.
- [10] C. H. Diaz, S. M. Kang, and C. Duvvury, "Modeling of electrical overstress in integrated circuit," Kluwer Academic Publishers, 1995.
- [11] Z. H. Liu, E. Rosenbaum, P. K. Ko, C. Hu, Y. C. Cheng, C. G. Sodini, B. J. Gross, and T. P. Ma, "A comparative study of the effect of dynamic stressing on high-field endurance and stability of reoxidized-nitrided, fluorinated and conventional oxides," in IEDM Technology Digest, pp. 723–726, 1991.
- [12] G. Chen, H. Fang, and A. Wang, "A systematic study of ESD protection structures for RF ICs," in Processing, IEEE Radio Frequency Integrated Circuit Symposium, Vol. 46, pp. 347–350, 2003.
- [13] J. Y. Choi, "AC modeling of the ggNMOS ESD protection device," ETRI Journal, Vol. 27, No. 5, pp. 628–634, October 2005.



On Possible A-Priori "Imprinting" of General Relativity Itself on the Performed Lense-Thirring Tests with LAGEOS Satellites

Lorenzo Iorio

INFN-Sezione di Pisa Permanent address for correspondence, Viale Unità di Italia, Bari (BA), Italy E-mail: lorenzo.iorio@libero.it Received October 23, 2009; accepted December 22, 2009

Abstract: The impact of possible a-priori "imprinting" effects of general relativity itself on recent attempts to measure the general relativistic Lense-Thirring effect with the LAGEOS satellites orbiting the Earth and the terrestrial geopotential models from the dedicated mission GRACE is investigated. It is analytically shown that general relativity, not explicitly solved for in the GRACE-based models, may "imprint" their even zonal harmonic coefficients of low degrees J_{ℓ} at a non-negligible level, given the present-day accuracy in recovering them. This translates into a bias of the LAGEOS-based relativistic tests as large as the Lense-Thirring effect itself. Further analyses should include general relativity itself in the GRACE data processing by explicitly solving for it.

Keywords: experimental studies of gravity, satellite orbits, harmonics of the gravity potential field

1. Introduction

The term "gravitomagnetism" [1–3] (GM) denotes those gravitational phenomena concerning orbiting test particles, precessing gyroscopes, moving clocks and atoms and propagating electromagnetic waves [4,5] which, in the framework of the Einstein's General Theory of Relativity (GTR), arise from non-static distributions of matter and energy. In the weak-field and slow motion approximation. the Einstein field equations of GTR, which is a highly non-linear Lorentz-covariant tensor theory of gravitation, get linearized [6], thus looking like the Maxwellian equations of electromagntism. As a consequence, a "gravitomagnetic" field \vec{B}_{g} , induced by the off-diagonal components g_{0i} , i=1,2,3 of the space-time metric tensor related to mass-energy currents, arises. In particular, far from a localized slowly rotating body with angular momentum Sthe gravitomagnetic field can be written as [7]

$$\vec{B}_{g}\left(\vec{r}\right) = \frac{G}{cr^{3}} \left[\vec{S} - 3\left(\vec{S} \cdot \hat{r}\right)\hat{r}\right],\tag{1}$$

where G is the Newtonian gravitational constant and c is the speed of light in vacuum. It affects, e.g., a test particle moving with velocity v with a non-central acceleration [7]

$$\vec{A}_{\rm GM} = \left(\frac{\vec{v}}{c}\right) \times \vec{B}_g \ . \tag{2}$$

It is the cause of the so-called Lense-Thirring¹ effect [9], which is one of the most famous and empirically inves-

tigated GM features; another one is the gyroscope precession [10,11], goal of the Gravity Probe B (GP-B) mission [12] whose data analysis is still ongoing [13].

The Lense-Thirring effect consists of small secular precessions of the longitude of the ascending node Ω and the argument of pericenter ω of the orbit of a test particle in geodesic motion around a slowly rotating body with angular momentum \vec{S} ; they are

$$\dot{\Omega}_{\rm LT} = \frac{2GS}{c^2 a^3 \left(1 - e^2\right)^{3/2}}, \quad \dot{\omega}_{\rm LT} = -\frac{6GS \cos I}{c^2 a^3 \left(1 - e^2\right)^{3/2}} \quad (3)$$

where *a* is the semimajor axis of the satellite's orbit, *e* is its eccentricity and *I* is the inclination of the orbital plane to the equatorial plane of the central body.

Concerning the possibilities of measuring it in the terrestrial gravitational field, soon after the dawn of the space age with the launch of Sputnik in 1957 it was proposed by Soviet scientists to directly test the Lense-Thirring effect with artificial satellites orbiting the Earth. In particular, V. L. Ginzburg [14–16] proposed to use the perigee of a terrestrial spacecraft in highly elliptic orbit, while A. F. Bogorodskii [17] considered also the node. In 1977-1978 Cugusi and Proverbio [18,19] suggested to use the passive geodetic satellite LAGEOS, in orbit around the Earth since 1976 and tracked with the Satellite Laser Ranging (SLR) technique, along with the other existing laser-ranged targets to measure the Lense-Thirring node precession. Since such earlier studies it was known that a major source of systematic error is represented by the fact that the even ($\ell = 2, 4, 6, ...$) zonal (m =

¹According to a recent historical analysis, it should be more correct to speak about an Einstein-Thirring-Lense effect [8]

0) harmonic coefficients J_{ℓ} , $\ell = 2,4,6$ of the multipolar expansion of the classical part of the terrestrial gravitational potential, accounting for its departures from spherical symmetry due to the Earth's diurnal rotation, induce competing secular precessions of the node and the perigee of satellites [20] whose nominal sizes are several orders of magnitude larger than the Lense-Thirring ones. In the case of the node, the largest precession is due to the first even zonal harmonic J_2

$$\dot{\Omega}_{J_2} = -\frac{3}{2}n\left(\frac{R_{\oplus}}{a}\right)\frac{\cos IJ_2}{\left(1-e^2\right)^2},\tag{4}$$

where R_{\oplus} is the Earth's mean equatorial radius and $n = \sqrt{GM_{\oplus}/a^3}$ is the satellite's Keplerian mean motion. For the other higher degrees the analytical expressions are more involved; since they have already been published in e.g., [21], we will not show them here.

Tests have started to be effectively performed about 15 years ago by Ciufolini and coworkers [22] with the LAGEOS and LAGEOS II satellites², according to a strategy by Ciufolini [23] involving the use of a suitable linear combination of the nodes Ω of both satellites and the perigee ω of LAGEOS II in order to remove the impact of the first two multipoles of the non-spherical gravitational potential of the Earth. Latest tests have been reported by Ciufolini and Pavlis [24,25], Lucchesi [26] and Ries and coworkers [27] with only the nodes of both the satellites according to a combination of them explicitly proposed by Iorio³ [28]. The total uncertainty reached is still matter of debate [32-38] because of the lingering uncertainties in the Earth's multipoles and in how to evaluate their biasing impact; it may be as large as ~20- 30% according to conservative evaluations [32,35–38], while more optimistic views [24,25,27] point towards ~10-15%.

To be more specific, the node-only combination used in the latest tests is

$$\dot{\Omega}^{\text{LAGEOS}} + k_1 \dot{\Omega}^{\text{LAGEOS II}}, \ k_1 = 0.554.$$
 (5)

It was designed to remove the effects of the static and time-varying components of J_2 , so that (5) is affected by the remaining even zonals of higher degree J_4 , J_6 ,... The gravitomagnetic trend predicted by (5) amounts to 47.8 milliarcseconds year⁻¹ (mas yr⁻¹ in the following) since the Lense-Thirring node precessions for the LAGEOS satellites are 30.7 mas yr⁻¹ (LAGEOS) and 31.5 mas yr⁻¹ (LAGEOS II). The Lense-Thirring signal is usually extracted from long time series of computed⁴ "residuals" of the nodes of LAGEOS and LAGEOS II obtained by processing their data with a suite of dynamical force

27

zonals is accounted for by using global solutions for the Earth's gravity field, in which general relativity has never been explicitly solved for⁵, produced by several institutions around the world from data of dedicated satel-lite-based missions like GRACE⁶ [42].

GRACE recovers the spherical harmonic coefficients of the geopotential from the tracking of both satellites by GPS and from the observed intersatellite distance variations [43]. The possible "memory" effect of the gravitomagnetic force in the satellite-to-satellite tracking was preliminarily addressed in [32]. Here we will focus on the "imprint" which may come from the GRACE orbits which is important for us because it mainly resides in the low degree even zonals.

2. A-Priori "Imprinting" of General Relativity on the GRACE-Based Models

Concerning that issue, Ciufolini and Pavlis write in [33] that such a kind of leakage of the Lense-Thirring signal itself into the even zonals retrieved by GRACE is completely negligible because the GRACE satellites move along (almost) polar orbits. Indeed, for perfectly polar (I = 90 deg) trajectories, the gravitomagnetic force is entirely out-of-plane, while the perturbing action of the even zonals is confined to the orbital plane itself. According to Ciufolini and Pavlis [33], the deviations of the orbit of GRACE from the ideal polar orbital configuration would have negligible consequences on the "imprint" issue. In particular, they write: "the values of the even zonal harmonics determined by the GRACE orbital perturbations are substantially independent on the a priori value of the Lense-Thirring effect [...]. The small deviation from a polar orbit of the GRACE satellite, that is 1.7×10^{-2} rad, gives only rise, *at most*, to a very small correlation with a factor 1.7×10^{-2} ". The meaning of such a statement is unclear; anyway, we will show below that such a conclusion is incorrect.

The relevant orbital parameters of GRACE are quoted in Table 1; the orbital plane of GRACE is, in fact, shifted by 0.98 deg from the ideal polar configuration, and contrary to what claimed in [33], this does matter because its classical secular node precessions are far from being negligible with respect to our issue. The impact of the Earth's gravitomagnetic force on the even zonals retrieved by GRACE can be quantitatively evaluated by computing the "effective" value⁷ $\overline{C}_{\ell 0}^{LT}$ of the normalized even zonal gravity coefficients which would induce classical secular node precessions for GRACE as large as those due to its Lense-Thirring effect, which is independent of the inclination *I*. To be more precise, $\overline{C}_{\ell 0}^{LT}$ come from solving the following equation which connects the classical even zonal

²LAGEOS II was launched in 1992.

See also [29–31].

⁴Actually, the nodes are not directly measurable quantities, so that speaking of "residuals" is somewhat improper.

For a critical discussion of such an issue, see [41].

⁶See on the WEB http://icgem.gfz-potsdam.de/ICGEM/ICGEM.html.

⁷It must be recalled that $J_{\ell} = -\sqrt{2\ell + 1}\overline{C}_{\ell 0}$ where $\overline{C}_{\ell 0}$ are the normalized gravity coefficients.

precession of degree $\ell \dot{\Omega}_{J_{\ell}} = \dot{\Omega}_{\ell} J_{\ell}$ to the Lense-Thirring node precession $\dot{\Omega}_{LT}$

 $\dot{\Omega}_{\ell}J_{\ell} = \dot{\Omega}_{\rm LT} \tag{6}$

In it

$$\Omega_{\ell} = f(a, e, I; R_{\oplus}, GM_{\oplus}) \tag{7}$$

are the coefficients of the classical node precessions depending on the satellite's orbital parameters and on the Earth's radius and mass. Table 2 lists $\overline{C}_{\ell 0}^{\text{LT}}$ for degrees ℓ =4,6, which are the most effective in affecting the combination (5). Thus, the gravitomagnetic field of the Earth contributes to the value of the second even zonal of the geopotential retrieved from the orbital motions of GRACE by an amount of the order of 2×10^{-10} , while for $\ell = 6$ the imprint is one order of magnitude smaller. Given the present-day level of accuracy of the latest GRACEbased solutions, which is of the order of 10^{-12} (Table 3), effects as large as those of Table 2 cannot be neglected. Thus, we conclude that the influence of the Earth's gravitomagnetic field on the low-degree even zonal harmonics of the global gravity solutions from GRACE may exist. falling well within the present-day level of measurability.

3. The Impact of the "Imprint" on the LAGEOS-LAGEOS II Tests

A further, crucial step consists of evaluating the impact of such an a-priori "imprint" on the test conducted with the LAGEOS satellites and the combination of Equation (5): if the LAGEOS-LAGEOS II uncancelled combined classical geopotential precession computed with the GRACEbased a-priori "imprinted" even zonals of Table 2 is a relevant part of, or if it is even larger than the combined Lense-Thirring precession, it will be demonstrated that the doubts concerning the a-priori gravitomagnetic "memory" effect are founded. It turns out that this is just the case because Equation (5) and Table 2 yield a combined geopotential precession whose magnitude is 77.8 mas yr⁻¹ (-82.9 mas yr⁻¹ for $\ell = 4$ and 5.1 mas yr⁻¹ for ℓ = 6), i.e. just 1.6 times the Lense-Thirring signal itself. This means that the part of the LAGEOS-LAGEOS II uncancelled classical combined node precessions which is affected by the "imprinting" by the Lense-Thirring force through the GRACE-based geopotential's spherical harmonics is as large as the LAGEOS-LAGEOS II combined gravitomagnetic signal itself.

We, now, comment on how Ciufolini and Pavlis reach a different conclusion. They write in [33]: "However, the Lense-Thirring effect depends on the third power of the inverse of the distance from the central body, i.e., $(1/r)^3$, and the J_2 , J_4 , J_6 ... effects depend on the powers $(1/r)^{3.5}$, $(1/r)^{5.5}$, $(1/r)^{7.5}$... of the distance; then, since the ratio of the semimajor axes of the GRACE satellites to the

Table 1. Orbital parameters of GRACE and its Lense-Thirring node precession. Variations of the orders of about 10 km in the semimajor axis *a* and 0.001 deg in the inclination *I* may occur, but it turns out that they are irrelevant in our discussion. (http://www.csr.utexas.edu/grace/ground/ globe.html)

<i>a</i> (km)	е	I (deg)	$\dot{\Omega}_{_{LT}}$ (mas yr ⁻¹)
6835	0.001	89.02	177.4

Table 2. Effective "gravitomagnetic" normalized gravity coefficients for GRACE ($\ell = 4,6; m=0$). They have been obtained by comparing the GRACE classical node precessions to the Lense-Thirring rate. Thus, they may be viewed as a quantitative measure of the leakage of the Lense-Thirring effect itself into the second and third even zonal harmonics of the global gravity solutions from GRACE. Compare them with the much smaller calibrated errors in \bar{C}_{40} and \bar{C}_{60} of the GGM03S model [44] of Table 3

$\overline{C}_{40}^{ ext{LT}}$	$ar{C}_{60}^{ ext{LT}}$
2.23×10 ⁻¹⁰	-2.3×10 ⁻¹¹

Table 3. Calibrated errors in the solved-for normalized gravity coefficients \overline{C}_{40} and \overline{C}_{60} according to the GGM03S global gravity solution by CSR [44]. They can be publicly retrieved at http:// icgem.gfz-potsdam.de/ ICGEM/ ICGEM. html. Compare them with the much larger "gravitomagnetic" imprinted coefficients of Table 2

$\sigma \overline{C}_{_{40}}$	$\sigma ar{C}_{_{60}}$
4×10 ⁻¹²	2×10 ⁻¹²

LAGEOS' satellites is $\sim \frac{6780}{12270} \approx 1.8$, any conceivable "Lense-Thirring Imprint" on the spherical harmonics at the GRACE altitude becomes quickly, with increasing distance, a negligible effect, especially for higher harmonics of degree $\ell > 4$. Therefore, any conceivable "Lense-Thirring imprint" is negligible at the LAGEOS' satellites altitude." From such statements it seems that they compare the classical GRACE precessions to the gravitomagnetic LAGEOS' ones. This is meaningless since, as we have shown, one has, first, to compare the classical and relativistic precessions of GRACE itself, with which the Earth's gravity field is solved for, and, then, compute the impact of the relativistically "imprinted" part of the GRACE-based even zonals on the combined LAGEOS nodes. These two stages have to be kept separate, with the first one which is fundamental; if different satellite(s) Y were to be used to measure the gravitomagnetic field of the Earth, the impact of the Lense-Thirring effect itself on them should be evaluated by using the "imprinted" even zonals evaluated in the first stage. Finally, in their latest statement Ciufolini and Pavlis write in [33]: "In addition, in (Ciufolini et al. 1997), it was proved with several

simulations that by far the largest part of this "imprint" effect is absorbed in the by far largest coefficient J_2 ." Also such a statement, in the present context, has no validity since the cited work refers to a pre-GRACE era. Moreover, no quantitative details at all were explicitly released concerning the quoted simulations, so that it is not possible to judge by.

4. Conclusions

We have analytically investigated the impact of possible a-priori "imprinting" effects of GTR itself on the ongoing Lense-Thirring tests with the LAGEOS satellites in the gravitational field of the Earth modeled from the dedicated GRACE mission.

The classical part of the terrestrial gravitational potential, acting as a source of major systematic error because of its even zonal harmonic coefficients $\overline{C}_{\ell 0}$, is retrieved from the data of the dedicated satellite-based GRACE mission. GTR, not explicitly solved for so far in GRACE data analyses, may impact the retrieved even zonals of the GRACE models at a non-negligible level ($\approx 10^{-10}-10^{-11}$ for $\ell = 4,6$), given the present-day level of accuracy (for $\ell =$ 4,6). It turns out that the resulting a-priori "imprint" of the Lense-Thirring effect itself on the LAGEOS-LAGEOS II data analysis performed to test it is of the same order of magnitude of the general relativistic signal itself.

Further, more robust tests should rely upon Earth gravity models in which GTR is explicitly solved for.

REFERENCES

- K. S. Thorne, "Gravitomagnetism, jets in quasars, and the stanford gyroscope experiment," in Near Zero: New Frontiers of Physics, J. D. Fairbank, B. S. Deaver, C. W. F. Everitt, and P. F. Michelson, Eds., W. H. Freeman and Company, New York, pp. 573–586, 1988.
- [2] W. Rindler, Relativity. Special, General and Cosmological, Oxford University Press, Oxford, pp. 195–198, 2001.
- [3] B. Mashhoon, "Gravitoelectromagnetism: a brief review", in The Measurement of Gravitomagnetism: A Challenging Enterprise, L. Iorio, Ed., Nova, Hauppauge, pp. 29–39, 2007.
- [4] M. L. Ruggiero and A. Tartaglia, "Gravitomagnetic effects," II Nuovo Cimento B, Vol. 117, No. 7, pp. 743– 768, July 2002.
- [5] G. Schäfer, "Gravitomagnetic effects," General Relativity and Gravitation, Vol. 36, No. 10, pp. 2223–2235, October 2004.
- [6] H. C. Ohanian and R. J. Ruffini, Gravitation and Spacetime, 2nd Edition, W. W. Norton and Company, New York, pp. 130–240, 1994.
- [7] B. Mashhoon, L. Iorio, and H. I. M. Lichtenegger, "On the gravitomagnetic clock effect," Physics Letters A, Vol.

292, No. 1-2, pp. 49-57, December 2001.

- [8] H. Pfister, "On the history of the socalled Lense-Thirring effect," General Relativity and Gravitation, Vol. 39, No. 11, pp. 1735–1748, November 2007.
- [9] J. Lense and H. Thirring, "Über den Einfluß der Eigenrotation der Zentralkörper auf die Bewegung der Planeten und Monde nach der Einsteinschen Gravitationstheorie," Physikalische Zeitschrift, Vol. 19, pp. 156–163, 1918.
- [10] G. E. Pugh, Proposal for a satellite test of the Coriolis prediction of general relativity WSEG Research Memorandum, The Pentagon, Washington DC, No. 11, November 1959.
- [11] L. I. Schiff, "Possible new experimental test of general relativity theory," Physical Review Letters, Vol. 4, No. 5, pp. 215–217, March 1960.
- [12] C. W. F. Everitt, S. Buchman, D. B. DeBra, G. M. Keiser, J. M. Lockhart, B. Muhlfelder, B. W. Parkinson, J. P. Turneaure and other members of the Gravity Probe B team, "Gravity Probe B: Countdown to Launch," in Gyros, Clocks, Interferometers: Testing Relativistic Gravity in Space, C. Lämmerzahl, C. W. F. Everitt and F. W. Hehl, Eds., Springer, Berlin, pp. 52–82, 2001.
- [13] C. W. F. Everitt, M. Adams, W. Bencze, S. Buchman, B. Clarke, J. W. Conklin, D. B. DeBra, M. Dolphin, M. Heifetz, D. Hipkins, T. Holmes, G. M. Keiser, J. Kolodziejczak, J. Li, J. Lipa, J. M. Lockhart, J. C. Mester, B. Muhlfelder, Y. Ohshima, B. W. Parkinson, M. Salomon, A. Silbergleit, V. Solomonik, K. Stahl, M. Taber, J. P. Turneaure, S. Wang, and P. W. Worden, "Gravity Probe B data analysis," Space Science Reviews, Vol. 148, No. 1–4, pp. 53–69, December 2009.
- [14] V. L. Ginzburg, "The use of artificial earth satellites for verifying the general theory of relativity," Advances in Physical Science (Uspekhi Fizicheskikh Nauk), Vol. 63, No. 1, pp. 119–122, 1957.
- [15] V. L. Ginzburg, "Artificial satellites and the theory of relativity," Scientific American, Vol. 200, No. 5, pp. 149– 160, May 1959.
- [16] V. L. Ginzburg, "Experimental verifications of the general theory of relativity," in Recent Developments in General Relativity, Pergamon press, London, pp. 57–71, 1962.
- [17] A. F. Bogorodskii, "Relativistic effects in the motion of an artificial earth satellite," Soviet Astronomy, Vol. 3, No. 5, pp. 857–862, October 1959.
- [18] L. Cugusi and E. Proverbio, "Relativistic effects on the motion of the earth's satellites," Journal of Geodesy, Vol. 51, pp. 249–252, 1977.
- [19] L. Cugusi and E. Proverbio, "Relativistic effects on the motion of earth's artificial satellites," Astronomy and Astrophysics, Vol. 69, pp. 321–325, October 1978.
- [20] W. M. Kaula, Theory of Satellite Geodesy, Blaisdell, Waltham, 1966.
- [21] L. Iorio, "The impact of the static part of the Earth's gravity field on some tests of General Relativity with satellite laser ranging," Celestial Mechanics and Dynamical

30

Astronomy, Vol. 86, No. 3, pp. 277-294, July 2003.

- [22] I. Ciufolini, D. M. Lucchesi, F. Vespe, and A. Mandiello, "Measurement of dragging of inertial frames and gravitomagnetic field using laser-ranged satellites," II Nuovo Cimento A, Vol. 109, No. 5, pp. 575–590, May 1996.
- [23] I. Ciufolini, "On a new method to measure the gravitomagnetic field using two orbiting satellites," II Nuovo Cimento A, Vol. 109, No. 12, pp. 1709–1720, December 1996.
- [24] I. Ciufolini and E. C. Pavlis, "A confirmation of the general relativistic prediction of the Lense-Thirring effect," Nature, Vol. 431, No. 7011, pp. 958–960, October 2004.
- [25] I. Ciufolini, E. C. Pavlis and R. Peron, "Determination of frame-dragging using Earth gravity models from CHAMP and GRACE," New Astronomy, Vol. 11, No. 8, pp. 527– 550, July 2006.
- [26] D. M. Lucchesi, "The lense thirring effect measurement and LAGEOS satellites orbit analysis with the new gravity field model from the CHAMP mission," Advances in Space Research, Vol. 39, No. 2, pp. 324–332, 2007.
- [27] J. C. Ries, R. J. Eanes and M. M. Watkins, "Confirming the frame-dragging effect with satellite laser ranging," in Proceedings of The 16th International Laser Ranging Workshop, "SLR-The Next Generation", Poznań (PL), 13–17 October 2008, S. Schillak, Ed. Available from: http://cddis.gsfc.nasa.gov/lw16/.
- [28] L. Iorio, "The new Earth gravity models and the measurement of the lense-thirring effect," in The Tenth Marcel Grossmann Meeting On Recent Developments in Theoretical and Experimental General Relativity, Gravitation and Relativistic Field Theories. Proceedings of the MG10 Meeting, Rio de Janeiro, Brazil 20–26 July 2003, M. Novello, S. P. Bergliaffa, and R. J. Ruffini, Eds. Singapore: World Scientific, 2006, pp. 1011–1020.
- [29] E. C. Pavlis, "Geodetic contributions to gravitational experi- ments in space," in Recent Developments in General Relativity: Proceedings of the 14th SIGRAV Conference on General Relativity and Gravitational Physics (Genova, IT, 18–22 September 2000), R. Cianci, R. Collina, M. Francaviglia, and P. Fré P., Eds. Milan: Springer, 2002, pp. 217–233.
- [30] J. C. Ries, R. J. Eanes and B. D. Tapley, "Lense-thirring precession determination from laser ranging to artificial satellites," in Nonlinear Gravitodynamics, The Lense-Thirring Effect, R. J. Ruffini and C. Sigismondi, Eds., World Scientific, Singapore, pp. 201–211, 2003.
- [31] J. C. Ries, R. J. Eanes, B. D. Tapley, and G. E. Peterson, "Prospects for an improved lense-thirring test with slr and the GRACE gravity mission," in Proceedings of The 13th International Laser Ranging Workshop, NASA CP (2003-212248), R. Noomen, S. Klosko, C. Noll and M. Pearlman, Eds., NASA Goddard, Greenbelt, 2003, http://cddisa.gsfc.nasa.gov/lw13/lw proceedings.html#science.
- [32] L. Iorio, "On the reliability of the so-far performed tests for measuring the Lense-Thirring effect with the LAG

EOS satellites," New Astronomy, Vol. 10, No. 8, pp. 603–615, August 2005.

- [33] I. Ciufolini and E. C. Pavlis, "On the measurement of the Lense-Thirring effect using the nodes of the LAGEOS satellites, in reply to 'On the reliability of the so-far performed tests for measuring the Lense-Thirring effect with the LAGEOS satellites' by L. Iorio," New Astronomy, Vol. 10, No. 8, pp. 636–651, August 2005.
- [34] D. M. Lucchesi, "The impact of the even zonal harmonics secular variations on the lense-thirring effect measurement with the two lageos satellites," International Journal of Modern Physics D, Vol. 14, No. 12, pp. 1989–2023, 2005.
- [35] L. Iorio, "A critical analysis of a recent test of the lense-thirring effect with the LAGEOS satellites," Journal of Geodesy, Vol. 80, No. 3, pp. 128–136, June 2006.
- [36] L. Iorio, "An assessment of the measurement of the lense-thirring effect in the earth gravity field, in reply to: "On the measurement of the lense-thirring effect using the nodes of the LAGEOS satellites, in reply to "On the reliability of the sofar performed tests for measuring the lense-thirring effect with the LAGEOS satellites" by L. Iorio," by I. Ciufolini and E. Pavlis," Planetary and Space Science, Vol. 55, No. 4, pp. 503–511, March 2007.
- [37] L. Iorio, "An assessment of the systematic uncertainty in present and future tests of the lense-thirring effect with satellite laser ranging," Space Science Reviews, Vol. 148, No. 1–4, pp. 363–381, December 2009.
- [38] L. Iorio, "Conservative evaluation of the uncertainty in the LAGEOS-LAGEOS II lense-thirring test," Central European Journal of Physics, Vol. 8, No. 1, pp. 25–32, February 2010.
- [39] D. M. Lucchesi and G. Balmino, "The LAGEOS satellites orbital residuals determination and the lense thirring effect measurement," Planetary and Space Science, Vol. 54, No. 6, pp. 581–593, May 2006.
- [40] D. M. Lucchesi, "The LAGEOS satellites orbital residuals determination and the way to extract gravitational and non-gravitational unmodeled perturbing effects," Advances in Space Research, Vol. 39, No. 10, pp. 1559– 1575, 2007.
- [41] K. Nordtvedt jr., "Slr contributions to fundamental physics," Surveys in Geophysics, Vol. 22, No. 5–6, pp. 597– 602, September 2001.
- [42] B. D. Tapley and Ch. Reigber, "The GRACE mission: status and future plans," EOS Transactions AGU 2001; 82: Fall Meeting Supplement G41, C–02.
- [43] C. Reigber, R. Schmidt, F. Flechtner, R. König, U. Meyer, K. H. Neumayer, P. Schwintzer and S. Y. Zhu, "An earth gravity field model complete to degree and order 150 from GRACE: EIGEN-GRACE02S," Journal of Geodynamics, Vol. 39, No. 1, pp. 1–10, January 2005.
- [44] B. D. Tapley, J. C. Ries, S. Bettadpur, D. Chambers, M. Cheng, F. Condi and S. Poole, "The GGM03 mean earth gravity model from GRACE," American Geophysical Union, Fall Meeting, abstract # G42A-03, 2007.



Neural Network Performance for Complex Minimization Problem

Tadeusz Wibig

Physics Department, University of Lodz; Cosmic Ray Laboratory, The A. Soltan Institute for Nuclear Studies, Uniwersytecka 5, 90-950 Lodz, Poland E-mail: wibig@zpk.u.lodz.pl Received October 19, 2009; accepted November 11, 2009

Abstract: We have analyzed the important problem of contemporary high-energy physics concerning the estimation of some parameters of the observed complex phenomenon. The standard statistical method of the data analysis and minimization was confronted with the Neural Network approaches. For the Natural Neural Networks we have used brains of high school students involved in our Roland Maze Project. The excitement of active participation in real scientific work produced their astonishing performance what is described in the present work. Some preliminary results are given and discussed.

Keywords: artificial neural network, natural neural network, curve fitting, minimization, interpolation, optimization

1. Introduction

The analysis of the surrounding physical reality, for last at least three thousand years, as we know it, follow the line of building simplified models and solving problems using specific tools developed with applied approximations making them easy or relatively easy to maintain. The unquestioned successes of such, scientific, way of thinking allow us to create so-called civilization.

However, this method can have limitations. Some problems can not be treated this way, at least at present. There is a belief that, e.g., mathematical tools needed to solve some problems in quantum field theory or hydro or thermodynamic will be developed in the future. However there are much common problems where the usual methods of standard analysis sometimes fail. The general problem of 'pattern' recognition is the perfect example.

We would like to discuss here a particular problem of the describing of the data registered by some cosmic ray physics experimental device. This is a problem of the general class of minimization or curve fitting. It is a good example to illustrate our general statement.

The statement is that the complex problems can be solved not only qualitatively but also quantitatively on the level of the standard statistical method precision not only by Artificial Neural Network (ANN) trained on the problem but with the over-sized, redundant Natural Neural Network (NNN) using their 'natural' abilities gathered in the past not obviously (obviously not) related to the particular problem. The method of the analyzing the NNN performance is described and some first results are given in this paper.

2. The Problem

The ultra high-energy cosmic ray particles, its origin and nature are one of the most intriguing questions on general interest among the physicists. The phenomenon of arriving form the cosmos of the elementary particle with energy of about 50 J is very rare and thus hard to investigate experimentally. Fortunately during the passage through the earth atmosphere the cascade of smaller energy secondary particles is created and eventually the surface is momentarily bombarded by billions of particles spread over the area of squared kilometers. The experimental setups for registrations of such events consists of several to several thousands detectors separated by hundreds of meters to few kilometers equipped with the triggering and recording devices.

Such arrays sample the mentioned showers of particles in not very big number of points and this is the only information we have about the event. (We do not discuss here the experiments recording the fluorescent light which is the distinct and complementary technique of study such phenomena). Each detector of the surface array registers actual number of particles passing the detector giving the information about particle density at the detector position. It is additionally smoothed by the physics of the detection process and electronic noises of different kinds. The transition from recorded digits to the physics in question is to estimate the shape of the distribution of cascade particles on the ground. The limited information allows us only to get the precise enough estimates of normalization constant total number of particles, first or at last second moment of the distribution (or any other, more suitable, parameters of this distribution). For doing this one has to use some prior assumptions, e.g. about the radial symmetry or the expected analytically approximated functional form of the distribution. After that one has to go through the procedure of making the estimate.

The standard is to make a χ^2 or likelihood measure of the goodness of the fit and than to use known textbook methods to minimize the respective distance between the 'theoretical line' and 'measured points'.

In general this is all we need, but practically there are classes of multiparameter problems and very noisy data when the minimization is not very straightforward. This is of course also the problem of the function to be minimized and its many local minima, distant and not very much different in depth. The problem of Extensive Air Shower (EAS) parameter estimation is a good example. The number of parameters is not very large. From physical point of view they are mainly: position of the shower axis and total number of particles and a parameter of the slope of theirs radial distribution. For simplicity and using the prior knowledge on the shower physics we use only one shape parameter [1]. The large spacing between detectors makes this simplification justified. We neglect also, for the purposes of this paper, the two parameters describing shower inclination. So we have only four parameter space with well defined physical meaning. It also provides a kind of independence of them all which is very helpful for minimization.

The problem arises because of the sharpness of the distribution of particles when one tests the distances close to the shower axis. When the detector gets close to the axis number of registered particles goes into thousands while a little far it goes to tens being on the edge of 1 and below in most other detection points. From the point of view of minimization procedure when the axis position is tested close to the one detection point it is the only one which controls the χ^2 or likelihood or whatsoever. This situation is caused by the physics of the process and one can't avoid it. The exclusion of close detectors is a remedy, but it is rather costly. The detector close to the shower core registered highest number of particles and thus the statistical importance of this point is the highest and in case of small number of detection points in general we can't afford to lose the most important one. We have to play with it.

Many methods and tricks were invented to get the minimization going with a number of problems as less as possible, but, as it will be shown, it is a hard task. The parameter which we will study comparing different methods of estimation is the number of lost events. We can define here the lost events as that having the χ^2 (or other studied measure) above some critical value. But to be comparable with others we defined them as the events

for which the minimization moves the values of the parameter of the axis position: x and y and shower size exceeding some limits (which can be treated as defining the divergence of the minimization procedure).

3. Artificial Neural Network Approach

The process of estimation EAS parameters with the help of Artificial Neural Network, as it is shown in [5] in the case of the hard shower component registered by the experiment KASCADE [4], can be quite successful. In the present work we used very similar network architecture which schematic view is shown in Figure 1. The input nodes are seeded with the registered particle densities, and the signal processing eventually gives the total number (its logarithm, to be precise) of shower particles. The particular network was build to work with the array of the Roland Maze Project being realized in Lodz [2]. The array is based on detectors placed on the roofs of city high school. In the final phase about 30 schools will be equipped with 4 one squared meter scintillator detector each. The sum of numbers of particles registered in each school carry the same information as the four numbers from all four detectors due to the Poissonian character of the cascade. The distance of about 10 meters within each school are negligible with the kilometers between the schools and the scale of the changes of the average particle density. Thus we need the ANN with 30 input nodes. The geometry of the network we used here was not optimized for the number of neurons and its final structure analyzed here is highly redundant. We used eventually two hidden layers with 20 and 10 neurons, respectively.

During tests we check that the input values could be logarithms of the value of the input signal surface enlarged by 1.0 to avoid the zeros from detectors with no muon registered. We do not add the electronic noise here as not very big. Each input is connected with each of the first hidden level neurons. The last hidden level neurons are connected to the single output unit. Tests with different number of hidden neurons shows that there is no effect on the network performance when we keep this numbers with reasonable limits. The number of the network parameters to be trained was from about 5000 to 25000. As the neuron response function the common sigmoid has been used. The network was trained with the standard back-propagation algorithm using the simple 'EAS generator'.

The generator works assuming particular shape of the particle distribution adopted from the measurements made by the one of the biggest arrays (in particular AGASA). The shower profile shape parameter known as "age parameter" defining the slope of the radial distribution was then smeared within physically reasonable limits. The number of particles is roughly proportional to the total energy of the primary cosmic ray particle. The normalization, total number of particles, was generated ac-
cording to the flat distribution in logarithm of particle energy scale. The cosmic ray shower spectrum is known to be of the power-law form and it is very steep with the index of about -3.0. This affects the estimated values. The generation used allows for the systematic bias, but on the other hand, too steep distribution in the training sample leads to the over-training of network with small size events while the events of energies, e.g., 100 times bigger, which are million times less abundant are practically not used for the training purposes. The uniform in log(*E*) is the compromise prior.

The steep spectrum makes the registered events consist mostly of events on the lower primary particle energy threshold which is defined by the trigger requirements. The 'artificial trigger' was applied to the training sample. We assumed that at least three 'schools' has to register some particles at least in two detectors each. Such trigger can be realized in the original Roland Maze Project array.

With the information limited so much, it is expected that standard minimization should fail quite often. The comparison of effectiveness of the standard and the ANN approach is one of the questions we want to answer here.

The network was trained first to estimate the most important shower parameter: the shower size (the total number of particles in the shower at the observation level). But we tested also the possibility of using network to estimate other parameters, and it was found that there is possible to train the network to estimate as well the *x*- and *y*-coordinate of the shower axis. The attempt to get the age parameter was not very successful.

In the Figure 2 the convergence of the training procedure is shown for the network trained with shower size a) and the axis position b). The dependence of the width of the distribution of the deviation of the estimated value from the *true* one is shown. The learning is guite a long lasting process. The first rational answers appear however already after the number of training events comparable with the number of internal weights. Then we observe the continuous improvement. An interesting feature appeared below 1 million events on Figure 2b. There are abrupt decreases of efficiency and then further and deeper improvements. The effect is seen for all networks we tested for both x- and y- axis position adjustment. It is seen always at the roughly some point and we suggests that it means the internal change of the network strategy of estimation. Something similar is seen also for EAS size estimation networks but at different length of training sample (around 10^4 at Figure 2a). The closer look at this phenomenon could put some light to the process of network learning but it is beyond of the scope of this paper.

We ended the learning process at the 10^8 event sample. The further improvement is interesting, but of no practical importance. The final state of the trained network allows us to use it as a tool for shower size and axis position determination. Such trained network was then applied



Figure 1. Schematic layout of the Artificial Neural Network used for the evaluation of the total number of the EAS particles

to the serial of 10000 events produced by the particles of energy generated by our event generator which build the library of the showers to by analyzed also with different methods. The ANN, when any event from the library is taken as an input, always give some answer. The accuracy of the ANN answer was studied in few 'modes'.

In the Figure 3 the illustration of an accuracy is presented as histograms showing the spread of the ANN guess errors. To get it easier to compare with other method some numbers should be given. Some measures of the 'goodness of the fit' are given below.

 σ_N - the accuracy of size determination measured as a dispersion of the difference between decimal logarithms of the *true* and *ANN reconstructed* shower size (total number of particles).

 $\Delta_{\rm N}$ - the bias of the shower size measured as a difference of the average *true* and ANN reconstructed shower size.

 $\xi_{\rm N}$ - the fraction of *perfect* reconstructions, which we defined as these within 10% around the *true* value of the shower size (logarithm).

 σ_R - the error in the localization measured as a average distance between *true* and ANN *reconstructed* shower axis position.

 $\xi_{\rm R}$ - the fraction of *perfect* localizations, by which we mean the *ANN reconstructed* axis closer than 100 m from the *true* one.

The values of all these five parameters obtained for our trained network are given in the Table 1 in the second column, labeled ANN.

4. Standard Minimization Approach

The data generated in the 10000 event shower library



Figure 2. The accuracy of the ANN answer as a function of the number of events using for the training process shown as a widths of the distribution of the difference between the decimal logarithm of the estimated shower size and the 'true' value a), and the difference of the respective spatial distance (in one x direction) difference b). Different lines represents different shower size samples. The solid one is for showers of the "true" size between 3 and 5 10^9 particles (the medium sizes), the dashed is for smaller showers, and the dotted one for really big showers

were also analyzed with the help of standard numerical minimization algorithms. We have used the CERN MIN-UIT package described in [3]. The straightforward application for such, not perfectly well determined, problem as shower parameter minimization works rather bad, so some slightly improved, thus much time consuming programs reaching the minimum of the likelihood in few steps, have to be developed. After careful adjustments of the proper divergence between 'the data' and predictions the program runs better. We have to mention here that some oversimplification was made here, because the radial distribution of shower particles used for minimization was exactly the one which was used inside the generator to calculate the averages. In the real case the particle distribution is rather unknown. This fact favours the minimization technique and the results given in this work should be treated with care, as the optimistic limit.

After applying to the library showers, the same as ana-

lyzed with the help of ANN, the results are as they are shown in the Table 1 in the columns labeled 'MINUIT'. There are two values in each case. The first on gives the average over all studied showers (we limited parameter ranges to reasonable values). Some showers due to the fluctuations can not give the minimum within the assumed ranges of parameters of the minimum found gives the value of χ^2 too high to be accepted. If we excluded them from the averaging procedures, the results get better



Figure 3. The spread of the trained ANN answer for the shower size a) and the spatial distance b) between guess and the true size and position of the shower axis. The result is obtained for the library showers

Table 1. Comparison of the performance of ANN, standard minimization and the three best NNNs found in our experiment

	ANN	MINUIT		NNN					
$\sigma_{\rm N}$	0.217	0.478	0.280	0.270	0.279	0.294	0.349		
Δ_{N}	0.29	-0.14	0.14	0.31	0.97	1.29	0.01		
ξN	6%	17%	20%	16%	17%	12%	13%		
σ_{R}	442	567	312	375	613	528	662		
ξ _R	5%	29%	34%	33%	36%	23%	18%		



Figure 4. The layout of the graphical interface to estimate the EAS parameters using NNN (e.g., 'by eye')

as they are given in the second MINUIT column. It is important to note that such 'bad' showers of of about 1/3 of all in the 5 x 10^{18} eV library.

5. Natural Neural Network Approach

The comparison of ANN and MINUIT methods of shower reconstruction given in the Table 1 shows that the Neural Networks can perform the competitive solutions of the complex minimization problem under study. However, the training procedure of the redundant network takes long and some systematic biases are seen. These are not very strong objections when taking into account the pros.

It would be interesting to test the performance of the extremely complex neural network one can imagine, which is, to some extend, the human brain. The training of the brain on a compound problem can be done in principle in two ways. The first one is similar to the ANN training process when the supervising teacher shows the examples of input data and told what the right answer should be. The more effective way is to use the ability of the brain collected in the whole 'common' life of the neural network as it is.

One can bet that in most cases the 'usual people' can react properly seeing the elephant running in their direction (whatever this reaction should be) even if they never met an elephant before. The brain (specially human) is so redundant that it can easy adapt the past solution to the new problem. The only requirement needed is that that the new situation must be similar to something seen before. If the similarity is closer than the probability of reaction is expected to be the right one is higher.

We want to use Natural Neural Networks (NNN) to

perform the estimations of the EAS parameters. The problem has to be transformed first to the form which can be understood by the 'common people'. The knowledge on the elementary particle passing through the atmosphere is helpful in principle, but is useless, in fact, in our case. We would like the NNNs to use their natural abilities.

We have built the graphical interface shown in Figure 4 which contains all what we know and all we should get.

On the left big panel the map of the city is shown in the scale, but without any unimportant information. The map shows position of the detectors (schools) in the Roland Maze Project shower array (the crosses). Next to some schools are the vertical lines (lighter and darker blue and red). The red shows the particle density registered by the detectors (its logarithm, but it is not relevant to the NNNs, and it hasn't been even told to them). The blue lines show 'the just proposed' solution.

The right panel shows, what can be told, the radial distribution of the particle density; horizontally: the distant to the shower axis is given, vertically: the particle density. Points (in red) show the same values as they are on the left plot but 'the just proposed' solution is given now by the line (blue).

The interface allows one to manipulate the shower parameter. The axis can be dragged usually using a computer mouse, or moves slightly with the batons with respective arrows next to the map side. The shower size (normalization of densities) can be changed on the left plot dragging the blue line vertically with the mouse. The horizontal movement of the clicked mouse increases the age parameter making the radial density curve wider or narrower. With this interface all the parameters can be adjusted comparing by eye the sizes of the bars of the right plot or/and, what is equivalent, controlling the positions of the points with respect to the curve on the left plot. The 'user manual' describing the interface is rather short and simple. The package was supplied with the set of examples showing how the *true* (known from the simulation) line should looks like. This explains additionally the task.

As the NNN donors we used pupils from the schools collaborating with the Roland Maze Project. They were a little familiar with the problem of EAS, but it was not a requirement, some of them were not. We assumed that each one of our volunteers perform the minimization of 100 showers. The practice shows that after the initial phase one shower takes about 2-5 minutes to be fitted. This gives few hours of the hard work. All NNNs were working on their free time, so we could not motivate them very strongly and there were a number of people which started and never finished the whole task. To avoid boring students we transferred the data to them in packages of 10 and the next 10 can be sent only after receiving the adjusted previous set. So the whole examining takes usually weeks of work.

The initial position of the shower axis, normalization and age parameters were taken randomly for each new event on display to avoid any unphysical guesses. Results were sent back by e-mail, but the program coded them. It was not possible to see what the numerical value was obtained and to correct them 'by hand'. The results once sent were put to the database and they couldn't be changed later on. The error once made remain what sometimes gives strong contribution to overall performance of the particular NNN.

6. Discussion

Anyway, we get some pupils completed their work. We (TW) did it also to be comparing with high school students' performance results. In the Table 1 results of TW followed by three (the best) student's NNNs are given in last four columns.

As it is seen the NNN accuracy is comparable with ANN concerning the shower size estimation (width and the bias), there are also no big differences concerning the axis position. Taking into account that NNN as well as ANN get an answer for each shower it should be compared with the 'all MINUIT' (third column in the Table 1).

It is interesting to compare results obtained by high school students and TW who can be called a specialist in the field, if not in the shower parameter estimation in general, then at least some specialist, because of building and testing the system of graphic interface *etc*. In fact there is no big difference (the sample of only 100 events was used to get the numbers). One can conclude that there is not experience needed.

Insights that the statistics education community badly needs to have, even though it may not know it yet.

Table 2. The improvement (or dis improvement) of the NNN performances concerning the parameter of σ_N in the course of the EAS analysis

			nun	iber of	analy	zed ev	ents			
	10	20	30	40	50	60	70	80	90	100
TW	0.142	0.203	0.199	0.210	0.220	0.239	0.250	0.264	0.261	0.270
NNN 1	0.257	0.237	0.244	0.233	0.235	0.312	0.302	0.291	0.286	0.279
NNN 2	0.379	0.374	0.324	0.307	0.287	0.322	0.305	0.308	0.295	0.294
NNN 3	0.551	0.427	0.390	0.389	0.368	0.363	0.358	0.349	0.335	0.349

It is not obvious, however, when we look how the individual NNN was improving its performance during the process. After analyzing the set of 10 events the results of the accuracy of their estimations were published in the web, and each participant can check how it has gone and what kind of error he made (specially the biases were easy to identify). The ability of the work with the interface could also getting better during the process of using it. Table 2 shows details in the case of the parameter describing the spread of the estimated shower size with respect to the true one. In some cases the improvement (NNN 3) is seen clearly, while for others (TW) the accuracy is diminishing with number of analyzed showers. This last is understood, because, in spite of the students, TW has not been limited to analyze only 10 events per day and the last 50 was taken just one by one continuously. The result is surprisingly big. If the constant care could be achieved during all the analysis process it is possible that the result of σ_N around 0.2. This value is exactly what has been achieved by the trained artificial neural network and significantly better than the standard statistical analysis.

7. Summary

With the help of number of enthusiastic young people we have shown that the redundant Neural Network, Artificial or Natural may work well and in fact in some cases even better than classical statistical tools of minimization. There is the evidence that NNN analyzed in the present work gone even better than the trained ANN. This suggests that the further studies of the over-sized networks and their performance are important and the minimization of the network size should not be taken on too early steps of the network arrangement, at least in some cases.

On the other hand the participation of young people, high school students on each level of the present work gives them a possibility to learn and understand the subject of statistics and data analysis on the level which is far beyond the standards even for the university students. This encourages us to propose further to the next groups of pupils the ambitious, extensive program for further studies of their brain performance and abilities. Interesting results are expected in the future.

REFERENCES

- [1] D. Barnhill, et al., [Pierre Auger Collaboration], Measurement of the lateral distribution function of UHECR air showers with the Pierre Auger observatory, Proceedings of the 29th International Cosmic Ray Conference, Pune, India, pp. 101–104; arXiv:astro-ph/0507590, 2005.
- [2] J. Feder, *et al.*, "The roland maze project: school-based extensive air shower network," Nuclear Physics Proceedings Supplements, No. 151, pp. 430–433, 2006.
- [3] F. James and M. Roos, "Minuit: A system for function minimization and analysis of the parameter errors and

correlations," Computer Physics Communications, Vol. 10, pp. 343–367, 1975.

- [4] H. O. Klages, *et al.*, "The KASCADE experiment," Nuclear Physics Proceedings Supplements, No. 52B, pp. 92–102, 1997.
- [5] T. Wibig, The artificial neural networks in cosmic ray physics experiment; I. Total muon number estimation. In A. P. del Pobil and J. Mira (Eds.) Lecture notes in computer science; Vol. 1416: Lecture Notes in Artificial Intelligence; Vol. 2 Tasks and methods in applied artificial intelligence, Springer-Verlag, Berlin, Heidelberg, New York, pp. 867, 1998.

Joint Power Control and Spectrum Allocation for Cognitive Radio with QoS Constraint

Zhijin Zhao¹, Zhen Peng¹, Zhidong Zhao¹, Shilian Zheng²

¹Telecommunication School, Hangzhou Dianzi University, Hangzhou, China ²State Key Lab. of Information Controlling Technology in Communications System of No. 36 Research Institute, China Electronic Technology Corporation, Jiaxing, China E-mail: {Zhaozj03, Zhaozd}@hdu.edu.cn, ipengzhen@163.com, lianshizhen@126.com Received October 11, 2009; accepted November 11, 2009

Abstract: Spectrum sharing with quality of service (QoS) requirement and power constraint on cognitive users is studied. The objective is to maximize the system throughput. This problem is modeled as a mixed integer nonlinear programming problem and then transformed to a continuous nonlinear programming problem through eliminating integer variables. We propose the joint power control and spectrum allocation algorithm based on particle swarm optimization to obtain the global optimal solution. Simulation results show that the proposed method can achieve higher system throughput and spectrum utilization under the constraints of transmit power and QoS requirement.

Keywords: spectrum sharing, power constraint, QoS, particle swarm optimization

1. Introduction

The spectrum of the wireless networks is generally regulated by governments via a fixed spectrum assignment policy. However, in recent years, the demand for wireless spectrum use has been growing dramatically with the rapid development of the telecommunication industry, which has caused scarcity in the available spectrum bands. Furthermore, the underutilization of the licensed spectrum bands makes the situation even worse [1]. Cognitive radio [2], with the ability to sense unused bands and adjust transmission parameters accordingly, is an excellent candidate for improving spectrum utilization. In cognitive radio networks, the cognitive (unlicensed) user needs to detect the presence of the primary (licensed) users as quickly as possible and dynamically changes the system parameters, such as transmit power level, so as to best utilize the valuable spectrum [3].

There are two kinds of spectrum sharing method: spectrum overlay and spectrum underlay. The researches of underlay spectrum limit the transmit power of the cognitive users and make sure that the interference temperature does not exceed certain threshold [4]. The related works on spectrum sharing schemes under interference temperature mainly include [5–7]. [5] regards the capacity of one cognitive link as an optimization problem with constraints in interference temperature and studies the optimal power allocation strategies. [6] studies the problem of channel selection in multi-hop cognitive mesh networks, but power allocation is not considered. With the assump-

tion that the primary users will always occupy the spectrum, these approaches can sufficiently increase the spectrum efficiency. [7] studies the joint of power control and random access under interference temperature, the optimization problem is transformed to a convex optimization problem. However, each cognitive user should be aware of the interference with the primary users and requires some kind of communications between the cognitive users and the primary users.

Previous works (such as [7]) on conventional OFDM systems are based on an implicit assumption that all the OFDM sub-carriers are fixed and always available. But in practice, the under-utilized spectrum which can be utilized by the secondary users varies over time, this is because the primary users can access to their spectrum unrestricted.

In this paper, we consider an overlay cognitive system, where multiple cognitive users coexist with multiple primary users and the availability of spectrum might not be contiguous because it is used by primary users. The multi-carrier system which dynamically operates in non-contiguous frequency bands and enabled by cognition technology is referred to as NC-OFDM [8]. The flexibility offered by NC-OFDM based CR can be employed to devise spectrum sharing schemes and provides QoS requirements by jointly considering variations in spectrum availability. We integrate the transmit power constraint and fairness of spectrum allocation in this paper. The optimization objective is to maximize the system throughput subject to maximize peak power constraints



and minimum QoS requirement on individual cognitive user. The QoS constraint is characterized by the minimum transmission rate requirement. To balance the power and QoS constraints, and further to efficiently and fairly utilize spectrum, transmit power and spectrum allocation must be determined by coordination among cognitive users.

The rest of this paper is organized as follows. The problem formulation and transformation are presented in Section 2. In Section 3, we propose the power and spectrum allocation algorithm based on particle swarm optimization. Section 4 includes simulation results and analysis. Conclusions are drawn in Section 5.

2. Mathmatics Model

In this paper, we consider a cognitive base station to multi cognitive users in wireless networks with rapid changes of spectrum opportunities. When the spectrum opportunities vary quickly, the cognitive users should frequently update the spectrum availably to avoid interference with the primary user. The cognitive base station balances the cognitive users' transmit power and spectrum to efficiently and fairly utilize spectrum.

Consider OFDM based CR system with a total bandwidth of B Hz and M primary users, each primary user with a bandwidth $B_m = B / M$ (m = 1, 2, ..., M), assume that B_m is less than the coherent bandwidth of the wireless channel, so that the channel response on each is flat. At the same time, there are N cognitive users in this system. At different location and time t, cognitive users have different available spectrum resource information because of the primary users' transmission activities. Define this available spectrum resource information as $L_{n,m}(t)$, where $L_{n,m}(t) \in \{0,1\}$. $L_{n,m}(t) = 1$ means the *m*th primary user use its own channel and the nth cognitive user can not use this channel, otherwise $L_{n,m}(t) = 0$. Let $G_{n,m}(t)$ denote the channel gain of the *n*th cognitive user on the *m*th channel and $p_{n,m}(t)$ be the transmit power of the *n*th cognitive user on the *m*th channel at time *t*, p_{nmax} is the maximum peak transmit power constraint of user n. We assume that the time variation of the wireless channel is stationary and slow enough, so that the cognitive users are able to perfectly estimate their local channels state information (CSI) on each channel and the cognitive base station knows all the CSI. Based on this CSI, cognitive base station balances the power and spectrum allocation to maximize the system throughput. Let $x_{n,m}(t) \in \{0,1\}$ indicate whether the spectrum is allocated to the cognitive user at time t. If $x_{n,m}(t) = 1$, the mth channel is assigned to the *n*th cognitive user, otherwise $x_{n,m}(t) = 0$. Each channel can be used by one cognitive user at any given time t, that is interpreted as $\sum_{n=1}^{N} x_{n,m}(t) \le 1$.

We assume that the network is under additive white Gaussian noise. We use M-ary quadrature amplitude modulation (MQAM) and then the maximum transmit rate of cognitive user n in channel m is given by:

$$r_{n,m}(t) = B_m \cdot \log_2(1 + \frac{-1.5 p_{n,m}(t) \cdot G_{n,m}(t)}{\log(5BER_{reg})\delta^2(t)})$$
(1)

where BER_{req} is an SNR gap parameter which indicates how far the system is operating from capacity, $\delta^2(t)$ is the interference power.

The objective is to maximize the cognitive system throughput

$$\sum_{n=1}^{N} \sum_{m=1}^{M} x_{n,m}(t) \cdot r_{n,m}(t)$$
 (2)

since power per user is finite in this system, every cognitive user has its own peak power constraint

$$\sum_{m=1}^{M} x_{n,m}(t) \cdot p_{n,m}(t) \le P_{n\max}$$
(3)

In practice, cognitive user transmission rate requirement is required no less than a certain threshold r_{n0} . It is defined as the QoS constraint of each cognitive user and expressed as

$$\sum_{m=1}^{M} r_{n,m}(t) \ge r_{n0}$$
 (4)

Note that if the available spectrum information $L_{n,m}(t) = 0$ or the *m*th channel has not been allocated to the *n*th cognitive user $x_{n,m}(t) = 0$, the transmit power must be zero. The base station should optimize the spectrum allocation matrix X and power matrix P. In this problem, $x_{n,m}(t)$ is an element of X and $x_{n,m}(t) \in \{0,1\}$,

 $p_{n,m}(t)$ is an element of P and $p_{n,m}(t) \in (0, p_{n \max})$.

Due to the discrete nature of channel and continuous nature of power, this optimization problem is a mixed integer nonlinear programming problem (MINLP). The difficulties in solving this MINLP problem come from the conflicting constraint sets, and coupled control variables. In [7] and many other works, they relax the binary valued constraint on the integer variable and replace it by a continuous variable. While this method causes inaccuracy of the algorithms and it can not find the optimal solution. In this paper, we first transform the MINLP problem to a continuous nonlinear programming (NLP) problem by introducing variable transformation, then we solve this problem by particle swarm optimization algorithm.

We substitute the variable $x_{n,m}(t)$ and $p_{n,m}(t)$ by

$$p_{n,m}(t) = x_{n,m}(t) \cdot p_{n,m}(t)$$
 (5)

so the variable $r_{n,m}(t)$ is transformed to

$$r_{n,m}(t) = x_{n,m}(t) \cdot r_{n,m}(t)$$
 (6)

Then the optimization problem is transformed as the following problem P1: P1.

$$Max \sum_{n=1}^{N} \sum_{m=1}^{M} r'_{n,m}(t)$$
 (7)

s.t.

$$r'_{n,m}(t) = B_m \cdot \log_2(1 + \frac{-1.5 p'_{n,m}(t) \cdot G_{n,m}(t)}{\log(5BER_{req})\delta^2(t)}) \quad (\forall n \in N)$$
(8)

$$\sum_{m=1}^{M} r'_{n,m}(t) \ge r_{n0} \quad (\forall n \in N)$$
(9)

$$\sum_{m=1}^{M} p'_{n,m}(t) \le P_{n\max} \quad (\forall n \in N)$$
(10)

$$p'_{n,m}(t) \cdot (\sum_{i \neq n}^{N} p'_{i,m}) = 0, \quad p'_{n,m}(t) \ge 0$$

 $(\forall n \in N, m \in M, i \in N)$ (11)

In P1, one continuous variable $p'_{n,m}(t)$ replaces the integer variable and the continuous variable, this substitution reduces the solution space dramatically. In addition, the new model is suitable for heuristic and search algorithms.

3. Power and Spectrum Allocation Algorithm Based on PSO

The particle swarm algorithm (PSO) is a swarm intelligence optimization algorithm modeled on the flight characteristics of birds [9,10]. In PSO, each solution is a 'bird' in the flock and is referred to as a 'particle', each particle has a position vector and velocity vector. The location of particles is the solution of optimization problem, the performance of each particle depends on the value of optimization objective's fitness function. Velocity vector used to determine particle velocity.

The following notation is needed in PSO. The number of particles in the population is denoted as Q. Let $y_i^k = [y_{i1}^k, y_{i2}^k, ..., y_{iD}^k]$ be the position of particle *i* $(1 \le i \le Q)$ at iteration *k*, where *D* is the number of dimensions to represent a particle and y_{id}^k is the *d*th $(1 \le d \le D)$ dimension of the position of particle *i*. Note that y_i^k is treated as a potential solution of the optimization problem. The velocity of particle *i* at iteration *k* is denoted as $v_i^k = [v_{i1}^k, v_{i2}^k, ..., v_{iD}^k]$, $v_{id}^k \in [-V_{\max}, +V_{\max}]$. Each particle in the swarm is assigned a fitness value indicating the merit of this particle such that the swarm evolution is navigated by best solutions. Let $s_i^k = [s_{i1}^k, s_{i2}^k, ..., s_{iD}^k]$ be the best solution that particle *i* has obtained until iteration *k*, and $s_b^k = [s_{b1}^k, s_{b2}^k, ..., s_{bD}^k]$ be the global best solution obtained from the population at iteration *k*.

The evoluationary process of the PSO is as follows:

$$v_{id}^{k+1} = \varpi v_{id}^{k} + c_1 u_1 (s_{id}^{k} - y_{id}^{k}) + c_2 u_2 (s_{bd}^{k} - y_{id}^{k})$$
(12)

$$y_{id}^{k+1} = y_{id}^{k} + v_{id}^{k+1}$$
(13)

where c_1 and c_2 are two positive constants named learning factors or acceleration coefficients, u_1 and u_2 are uniform random numbers distributed in the range [0, 1], and ϖ is an inertia weight employed to control the impact of the previous history of velocities on the current velocity. Note that Equation 12) specifies that the velocity of a particle at iteration k is determined by the previous velocity of the particle, the cognition part, and the social part.

In the PSO-based spectrum and power allocation algorithm, each particle's position vector specifies a possible spectrum and power allocation scheme. The penalty function is used to solve the constrained optimization problem. Ordinary penalty function only calculates the total violation of individuals, but does not make full use of the violation information of the infeasible solutions. We use the penalty function which is not only depends on the number of constraint violations but also on the degree of constraint violations. The performance of this method is better than that using the ordinary penalty function [11]. As a result of the different scales in constraints, it is possible that some certain constraints play a dominant role in the total constraints and other constraints may not reflect their degree of constraint violations. In addition, the objective function and the violations of constraint functions may be in different scales, so we normalize the objective function and constraint functions to solve this problem.

We use the following fitness function to evaluate the particle:

$$F(t) = \sum_{n=1}^{N} \sum_{m=1}^{M} r'_{n,m}(t) - \min(\sum_{n=1}^{N} \sum_{m=1}^{M} r'_{n,m}(t)) \times (w1 \cdot sum_viol + w2 \cdot num_viol)$$
(14)

where sum_viol represents the total amount of the constraint violations and num_viol represents the number of the constraint violations. If any user *n* in particle **p**' violates the transmission rate constraint (9) or power constraint (10), the num_viol of **p**' will plus



Figure 1. Coding scheme of particle

one.

$$sum_viol = \sum_{n=1}^{N} \frac{J_n(\mathbf{p}')}{\max_i (J_n(\mathbf{p}'))} + \sum_{n=1}^{N} \frac{H_n(\mathbf{p}')}{\max_i (H_n(\mathbf{p}'))} \quad (15)$$

$$J_n(\mathbf{p}') = \max(0, j_n(\mathbf{p}')) \tag{16}$$

$$H_n(\mathbf{p}') = \max(0, h_n(\mathbf{p}')) \tag{17}$$

$$j_{n}(\mathbf{p}') = r_{n0} - \sum_{m=1}^{M} r'_{n,m}(t)$$
(18)

$$h_n(\mathbf{p}') = \sum_{m=1}^{M} p'_{n,m} - P_{n\max}$$
(19)

where $\sum_{n=1}^{N} \sum_{m=1}^{M} r'_{n,m}(t)$ is the system throughput of the particle **p**' in time slot *t*. **p**' is an *N* by *M* matrix representing the power and channel allocation, if $p'_{n,m} > 0$, the channel *m* is assigned to the cognitive user *n*, otherwise $p'_{n,m} = 0$.

In the proposed PSO-based power and spectrum allocation algorithm, a particle specifies a possible power and spectrum allocation assignment. As $p_{n,m} = 0$ when $L_{n,m}(t) = 0$, if we use one bit to encode every element in **p**', there will be a lot of redundancy in the particle. We encode only those elements which may take the value 1, i.e., $p_{n,m}$ where (n,m) satisfies $L_{n,m}(t) = 1$. As a consequence, the length of the coding string is equal to the number of elements equal to 1 in *L*. Figure 1 illustrates the structure of an example particle, where N = 5, M = 6. Note that encoding all the elements needs 30 bits, while encoding only the elements with underline only needs 9 bits. In order to evaluate the fitness of the particle, we need to map the particle to the assignment matrix **p**', as the arrows show in Figure 1.

The value of every bit in the particle is randomly generated at the initial population and this coding scheme reduces the searching space of the optimization problem efficiently.

The proposed PSO-based power and spectrum as-

signment algorithm proceeds as follows:

Step 1: cognitive user gets the available spectrum resource information matrix L and channel information matrix $G_{n,m}$, then transmits these information to the cognitive base station.

Step 2: set k = 0, and randomly generate y_{id}^k and v_{id}^k , where $v_{id}^t \in [-V_{\max}, +V_{\max}], 1 \le d \le D$, thus obtaining $y_i^k = [y_{i1}^k, y_{i2}^k, ..., y_{iD}^k], 1 \le i \le Q$.

Step 3: map y_{id}^k $(1 \le i \le Q)$ to $p_{n,m}$, where (n,m) is the *d*th element with $L_{n,m}(t) = 1$.

Step 4: compute the fitness value of each particle in the population according to Equation 14), set $s_i^k = [y_{i1}^k, y_{i2}^k, ..., y_{iD}^k]$ and $s_b^k = [y_{b1}^k, y_{b2}^k, ..., y_{bD}^k]$, where *b* is the index of the particle which has the highest fitness value.

Step 5: set k = k + 1, and update v_{id}^k according to Equation 12). If $v_{id}^k > V_{\max}$, then set $v_{id}^k = V_{\max}$; if $v_{id}^k < -V_{\max}$, set $v_{id}^k = -V_{\max}$.

Step 6: update y_{id}^k according to Equation 13) and map y_{id}^k to $p'_{n,m}$.

Step 7: compute the fitness value of each particle in the population. For particle *i*, if it's fitness value is greater than the fitness value of s_i^{k-1} , then set $s_i^k = [y_{i1}^k, y_{i2}^k, ..., y_{iD}^k]$. If particle *i*'s fitness value is greater than the fitness value of s_b^{k-1} , then set $s_b^k = [y_{i1}^k, y_{i2}^k, ..., y_{iD}^k]$.

Step 8: if *k* equals to the predefined maximum iteration, then the algorithm is terminated, map $s_b^k = [s_{b1}^k, s_{b2}^k, ..., s_{bD}^k]$ to **p**'; else, go to Step 5.

4. Simulation Result and Analysis

To evaluate the proposed algorithm, simulations were performed for the OFDM based CR system. The bandwidth of the OFDM system is B = 6 MHz, which is li-



Figure 2. Convergence of proposed algorithm under different QoS constraints($P_{nmax} = 1.5$)



Figure 3. Convergence of proposed algorithm under different QoS constraints($P_{n \max} = 2.5$)

censed to M = 12 primary users, every primary user's transmission uses one channel and the available spectrum resource information matrix *L* is generated randomly. The number of cognitive users is N = 10. The required bit error rate of each transmission is supposed to be $BER_{req} = 10^{-6}$. For simplicity, each cognitive link's average channel gain is chosen randomly within (0,0.01) and the interference power is 0.5 mW.

The parameters for the PSO are Q = 20, $c_1 = c_2 = 2$, and $V_{\text{max}} = 4$, and PSO would be terminated after 3000 iterations.

Figure 2 and Figure 3 illustrate the convergence proc-



Figure 4. Convergence of proposed algorithm under different transmit power constraints



Figure 5. System throughput at different time

ess of the proposed PSO-based power and spectrum allocation algorithm with different transmission rate requirements $r_{n0} = r_0$ under the peak transmit power constraint. The peak transmit power constraints are $P_{n\max}$ = 1.5 and $P_{n\max}$ = 2.5 respectively. The QoS requirement of each user is set to 1500bps, 2500bps and 3500bps respectively. As can be observed in Figure 2 and Figure 3, after about 2500 iterations, the proposed algorithm achieves the optimal solution. Further more, the system throughput doesn't increase with the transmission rate requirement increase, this is because the system throughput is also constrained by users' transmit power. In addition, the peak transmit power provides allocation fairness. In principle cognitive users with high channel gains are



Figure 6. Convergence of proposed algorithm and lagrange algorithm

taking more channels, but the more channels they take, the more power will be consumed. Then other users with weaker channel gains and more available power can take the rest channels to transmit and further increase the system throughput.

Figure 4 shows the convergence process of the proposed algorithm with same transmission rate requirements $r_0 = 2000$ bps under different peak transmit power constraints $P_{n\max}$. The peak transmit power is set to 2W, 2.5W and 3W respectively. We can clearly see that the system throughput is increasing as the peak transmits power increases.

Figure 5 shows that the system throughput is fluctuating at different time. The peak transmits power constraints are $P_{n\text{max}} = 2$ and QoS requirement of each user is set to 2500bps. At different time *t*, the CR system has different available spectrum information and channel state information because of the activities of primary users. Sometimes the primary users are not active, so the cognitive users have more available spectrum resource and the channel gains are better. As a result, the system through is higher than some situations which primary users are active.

Figure 6 shows the convergence processes of the proposed algorithm's performance and the Lagrange algorithm in [7]. The peak transmits power constraints are $P_{n\max} = 2$ and QoS requirement of each user is set to 2500bps. We can see that the proposed algorithm has higher system throughput and faster convergence speed

than the lagrange algorithm.

5. Conclusions

We model the power control and spectrum allocation problem as a mixed integer nonlinear optimization problem. This MINLP problem is difficult to find the optimal solution, so we transform the MINLP problem to an NLP problem. Then we use a coding scheme and PSO-based power control and spectrum allocation algorithm to solve the NLP problem. Simulations show that the proposed model provides the fairness of the assignment and the proposed algorithm performs better than the Lagrange algorithm.

REFERENCES

- J. Mitola, "The software radio architecture," IEEE Communications Magazine, 1995.
- [2] J. Mitola and G. Q. Maguire, "Cognitive radio: making software radios more personal," IEEE Pers. Commn, Vol. 6, pp. 13–18, 1999.
- [3] J. Mitola, "Cognitive raido: an integrated agent architecture for softoware defined radio," Doctor of Technology, Royal Institute Technology Stockholm, Sweden 2000.
- [4] Federal Communications Commission, Spectrum policy task force, Report ET Docket, 2002.
- [5] A. Ghasemi and E. S. Sousa, "Fundamental limits of spectrum sharing in fading environments," IEEE Trans Wirel. Commun, Vol. 6, pp. 649–658, 2007.
- [6] M. Sharma, A. Sahoo, and K. Nayak, "Channel selection under interference temperature model in multi-hop cognitive mesh networks," Proceeding of IEEE DySPAN, 2007.
- [7] B. Yang, Y. Y. Shen, and G. Feng, "Distributed power control and random access for spectrum sharing with QoS constraint," Computer Communications, Vol. 31, pp. 4089–4097, 2008.
- [8] M. Wylie-Green, "Dynamic spectrum sensing by multiband OFDM radio for interference mitigation," Proceeding of IEEE DySPAN, pp. 619–625, November 2005.
- [9] J. Kennedy and R. Eberhart, "Particle swarm optimization [A]," Proceeding of IEEE International Conference on Neural Networks, 1942–1948, 1995.
- [10] Y. Shi and R. Eberhart, "A modified particle swarm optimizer," Proceeding of the IEEE International Conference on Evolutionary Computation, pp. 69–73, 1998.
- [11] Z. Michalewicz and N. F. Attia, "Evolutionary optimization of constrained problems," Proceeding of CEP, pp. 98–108, 1994.



A Historical Narrative of Study of Fiber Grating Solitons

Xiaolu Li¹, Yuesong Jiang², Lijun Xu¹

¹School of Instrumentation Science and Opto-electronics Engineering, Beihang University, Beijing, China ²School of Electronic and Information Engineering, Beihang University, Beijing, China *E-mail: xiaolu5253@126.com Received September 18, 2009; accepted November 10, 2009*

Abstract: A brief historical narrative of the study of grating solitons in fiber Bragg grating is presented from the late 1970's up to now. The formation of photogeneration gratings in optical fiber by sustained exposure of the core to the interference pattern produced by oppositely propagating modes of argon-ion laser radiation was first reported in 1978. One important nonlinear application of fiber Bragg grating is grating solitons, including gap soliton and Bragg soliton. This paper summarily introduces the numerous theoretical and experimental results on this field, each indicating the potential these solitons have in all-optical switching, pulse compression, limiting, and logic operations, and especially important for the optical communication systems.

Keywords: nonlinear optics, periodic structure, fiber Bragg grating, kerr nonlinearity, dispersion, grating solitons, Bragg soliton, gap soliton

1. Introduction

After the invention of the laser, there has been much interest in propagating nonlinear pulses through the periodic medium such as a fiber Bragg grating (FBG), which is a periodic variation of the refractive index of the fiber core along the length of the fiber. Since the tirst demonstration of photo-induced optical fiber Bragg gratings by Hill and coworkers in 1978 [1], significant progress was made in the fabrication technology of fiber Bragg reflectors [2-5]. The concept of "photonic band structure" is introduced by Yablonovitch in the late 1980's [6]. A notable feature of this linear periodic structure is the presence of stop gap in the dispersion curve popularly known as photonic band gap (PBG) [7,8]. This PBG exists at frequencies for which the medium turns highly reflective and hence the light pulse will not be able to propagate through the periodic structure. Light interaction with nonlinear periodic media yields a diversity of fascinating phenomena, among which two solitonic phenomena have been studied most intensively, namely, discrete (or lattice) solitons [9–11] and gap (or Bragg) solitons [12–17]. While discrete solitons are spatial phenomena in twodimensional or three-dimensional arrays of coupled waveguides, gap solitons are usually considered as a temporal phenomenon in one-dimensional (1D) periodic media [18-20]. Perhaps the most fascinating feature of solitons is their particle like behavior. Survival of two such colliding solitons is even more remarkable if one notes that solitons interact strongly with each other during the collision. But for copropagating solitons, the interaction is either attractive or repulsive, depending on the relative phase between two solitons. In both cases the evolution of the soliton pair is well understood [21–24].

As first pointed out by Winful [25], because the dispersion is many orders of magnitude larger than the total dispersion due to the combined effects of material and waveguide dispersions that arise in the conventional fibers, the interactions lengths are reduced accordingly. Hence, the grating induced dispersion dominates over the total dispersion in the conventional fibers. When the entire spectral components of the input pulse lie within the PBG structure, the grating induced dispersion counterbalanced by the Kerr nonlinearity through the self-phase modulation (SPM) and cross-phase modulation (XPM) effects, forming solitons are referred to as gap solitons since their spectral components are within the PBG structure. Many research groups [3-10] theoretically predicted the existence of gap solitons and Bragg grating solitons in FBG and the investigations on these exciting entities are going on. However, it can be noticed that, in literatures, nowadays the distinction between gap solitons and Bragg solitons is hardly maintained and, in general, they are simply called grating solitons [26]. Ul [25], because the dispersion is many orders of magnitude larger than the total dispersion due to the combined effects of material and waveguide dispersions that arise in the conventional fibers, the interactions lengths are reduced accordingly. Hence, the grating induced dispersion dominates over the total dispersion in the conventional

fibers. When the entire spectral components of the input pulse lie within the PBG structure, the grating induced dispersion counterbalanced by the Kerr nonlinearity through the self-phase modulation (SPM) and crossphase modulation (XPM) effects, forming solitons are referred to as gap solitons since their spectral components are within the PBG structure. Many research groups [3-10] theoretically predicted the existence of gap solitons and Bragg grating solitons in FBG and the investigations on these exciting entities are going on. However, it can be noticed that, in literatures, nowadays the distinction between gap solitons and Bragg solitons is hardly maintained and, in general, they are simply called grating solitons [26].

2. Theory

The usual quantitative description of grating solitons employs coupled-mode theory, leading to the nonlinear coupled-mode equations. In addition, in the appropriate limit, the envelope of the electric field satisfies the nonlinear Schrödinger (NLS) equation. The pulse propagation through the FBG is described by the nonlinear-coupled mode (NLCM) equations which are nonintegrable in general. Therefore, the analytical solutions of the NLCM equations are not solitons but solitary waves that can propagate through FBG without changing their shape. These are obtained from the approximated nonlinear Schrödinger (NLS) equation that results from reducing the NLCM equations using the multiple scale analysis. The relation between the NLSE and the more general CME description, which was discussed earlier [28], is important. Gap solitons are obtained from the NLCM equations and their spectra lie within the photonic bandgap structure. There is another class of solitons called Bragg solitons obtained from the NLS equations whose frequencies fall close to, but outside, the band edge of the photonic bandgap. Generally speaking, the gap solitons are the special class of Bragg solitons.

For the first time, Chen and Mills [12] have analyzed the properties of these gap solitons in nonlinear periodic structure. Thereafter, Sipe and Winful published analyses showing that these "gap-solitons" are not only fundamental solutions in the weak-field regime but could be detected as propagating solutions in structures of finite length [14]. The general gap soliton solutions to the coupled mode equations were first obtained in a limiting case by Christodoulides and Joseph [16]. The solutions were first reported in their most general form by Aceves and Wabnitz [17]. Aceves and Wabnitz appoint parameters to form gap solitons in fiber Bragg grating, and the unique dispersion relation of the fiber grating, and the corresponding solitons, allows in theory all velocities from zero to the speed of light in the bare fiber. Their starting point is the massive Thirring model(MTM), and quantitative description of gap solitons employs coupled-mode theory, leading to the nonlinear coupled-mode equations [16,17]. At same time, Sipe and de Sterke examined, in further publications [27-29], the pulse transmission behavior as a function of both pulse energy and detuning from the Bragg resonance. Among the contributions of de Sterke. Sipe and others was a rigorous development of coupled-wave and multiple-scales approximations as well as the description of numerical methods [30] suitable for examining the regimes of instability of these structures. In a word, Sipe and Winful [14], Christodoulides and Joseph [16], Aceves and Wabnitz [17], and Winful et al. [31] have obtained the analytical solutions for the grating solitons. These solitons in FBGs have been extensively reviewed in [19,32]. Comprehensive analyses of Bragg solitons stability have also been reported [33,34]. Still other generalizations have been discussed by Feng and Kneubuhl [35] and by Feng [36]. In order to better simulate experimental conditions. Broderick, de Sterke and Jackson presented a method of numerically modeling periodic structures having optical nonlinearities [37]. Other important extensions and generalizations include a series of papers by Aceves and coworkers extending many of these principles to waveguide arrays [38].

Inverse scattering transform (IST) is currently the standard analytical technique for obtaining the soliton solution for the homogenous NLSE [39,40]. IST has been used to solve the two-dimensional space-time NLSE with initial-boundary conditions and coupled NLSE in the form of fundamental and higher-order solitons [39]. To our knowledge, no other analytical method has been published besides the IST for solving the NLSE systems. Another method can be described as effective particle pictures EPP's, since they represent the continuous field distribution as a point particle with a limited number of degrees of freedom. The key difference between the NLSE and NLCME's is that the NLSE is integrable, whereas NLCME's are not [37], hence that an EPP would be more accurate in that case [42-46]. However, previously, gap soliton propagation in the presence of uniform gain and loss was succesfully treated using an EPP [43,47] method, which was also used by Capobianco et al. to treat propagation between two quadratically nonlinear materials [48]. One method to analyze deep gratings is using Bloch wave solutions as the fundamental waves. Actually the modulation of a single Bloch wave is known to obey the nonlinear Schrödinger equation in Kerr optical media [13,49,50], and its fundamental soliton corresponds to gap solitons in this geometry. Note that the Bloch function formalism has the feature that the linear system needs to be solved first, and the nonlinearity is then considered as a perturbation which can be treated in a variety of approximations. A different formalism developed for linear gratings only to treat deep gratings was reported by Sipe et al. [51]. The linear

properties are therefore not obtained exactly, but in terms of an asymptotic series, only a few terms of which are retained. Nonetheless, the method leads naturally to low-order corrections to the coupled mode equations for shallow gratings. Then, one may expect that the model may give rise to two qualitatively different families of gap solitons: low-frequency ones, in which the self-focusing (cubic) nonlinearity is balanced by the dispersion branch with a sign corresponding to anomalous dispersion, and high-power solitons, supported by the balance between self- defocusing (quintic) nonlinearity and the normal branch of the dispersion. The simplest model of this type may be based on the cubic-quintic (CQ) nonlinearity that has recently attracted considerable attention, as the combination of the SF cubic and SDF quintic terms prevents collapse and makes it possible to anticipate the existence of stable solitons [52-60]. Atai and Malomed introduced the quintic nonlinearity into the NLCM equations and investigated two different families of zero-velocity solitons. One family was the usual Bragg grating solitons supported by the cubic nonlinearity. The other family was named as twotier solitons supported by the quintic nonlinearity [26]. In fact, in the cubic model, the final soliton retains only 11.6% of the initial energy, while the energy-retention share in the cubic-quintic model is 92.4% [59].

3. Experimentation and Applications

Recently conducted experiments have provided strong evidence for the existence of the grating solitons in FBGs [61-66]. To our knowledge, it was Larochelle, Hihino, Mizrahi and Stegeman [67] who were the first to report (in 1990) an experimental investigation of the optical response of nonlinear periodic structures. They employed an optical Kerr-effect cross-phase modulation in fiber gratings to achieve switching of a probe beam by a control beam. The first detailed experimental observation of all-optical switching dynamics in a nonlinear periodic structure was reported by Sankey, Prelewitz and Brown in 1992 [68]. Experimental observations of nonlinear grating behaviour are limited, principally by the difficulty in getting sufficiently high power densities within the core of a FBG in a suitable spectral and temporal range. In order to reduce the nonlinear threshold for gap soliton formation one can use the somewhat weaker dispersive properties of FBGs outside of the band gap. An investigation of nonlinear pulse propagation in uniform fiber gratings was published by Eggleton et al. in 1996 [61].In this report, the Bragg solitons are most easily generated in the laboratory travel at 60-80% of veocity of light in fiber absence of grating [61,64]. This was followed by further reports from the same group, which both refined the experimental technique and broadened the experimental understanding of the dynamics of pulse

propagation in periodic structures [65]. In their initial experimental observations of Bragg solitons [61.62.64]. the agreement between the experiments and the numerical calculations was qualitative. However, stationary (or nearly stationary) gap solitons have not been observed vet. Subsequently, the Southhampton group [69] first demonstrated switching at the important optical communication wavelength of 1550 nm, and in doing so have confirmed certain key aspects of the physics of pulse propagation in nonlinear periodic structures. We now understand that a Bragg soliton need not be centered near the Bragg resonance--indeed, some very interesting propagation effects occur rather far from the band edge. Experimental studies of BG solitons were further developed including, in particular, formation of multiple BG solitons inRefs [42]. Broderick et al. also report the first experimental demonstration of a novel type of all-optical pulse compression [71]. It is significant experimentation that Taverner et al. [42,70] reported the first observation of gap soliton generation in a Bragg grating at frequencies within the photonic bandgap. Furthermore the sets of experiments were performed in relatively short gratings. Thus, in these experiments, pure soliton propagation effects are difficult to distinguish from effects due to soliton formation. The occurrence of modulational instability (MI) in fibers had been first suggested by Hasegawa and Brickman [72] and experimentally verified by Tai et al. [73]. The effects of MI which occurs when a perturbed continuous wave experiences an instability that leads to an exponential growth of its amplitude or phase during the course of propagation in optical fibers due to an interplay between the nonlinearity and group velocity dispersion (GVD) act in opposition. THE studies on modulational instability (MI) have some impacts on solitons [8,74,75].

The researchers recently have realized the potential applications of these solitons in fiber Bragg grating for all-optical switching [67,76,77], pulse compression [69,71,78], limiting [80], and logic operations [81,84], also promising for the fiber-sensing technology [79], especially important for the optical communication systems [78,82]. One would hope to achieve zero velocity by a clever tailoring of the Bragg grating. This research goes beyond its intellectual value; all optical buffers and storing devices can be based on such fibers. About logic operations, for the first time to our knowledge, an alloptical 'AND' gate based on a configuration proposed by S. Lee and S.T. Ho [84]. The operation of the gate relies on the formation and propagation of coupled gap solitons by two orthogonally polarised high intensity input beams incident within the bandgap of a FBG [81]. Recently Nuran Dogru was pursueing for the hybrid soliton pulse source (HSPS) developed as a pulse source for the soliton transmission system [88-92]. In a Bragg grating SPM results in the transmission being bistable with one state (high power) having a transmission of unity while in the other (low power) the transmission is vanishingly small [31]. For strong optical pulses this behavior can result in all-optical switching. The all- optical switching of a fiber Bragg grating (FBG) was first seen by La-Rochelle et al. in 1990 [67] using a self-written grating centered at 514 nm. In their experiment the probe beam was centered on the grating, while the pump beam had a wavelength of 1064 nm. It was in this vein that Radic, George and Agrawal suggested the use of 1/4 phaseshifted gratings for use in optical switching [77]. Ju Han Lee [85–87] demonstrate the use of a superstructured fiber Bragg grating obtain more optimal operation of nonlinear all-optical switches [85], all-optical modulation and demultiplexing systems [86], tunable optical pulse source [87]. In long distance communications, that a third-order nonlinear effect is together with anomalous dispersion, can result in the formation of bright temporal optical solitons. Beacause of the shape-preserving property of the bright and dark solitons, they have received considerable attention from optical communication industries. Solitons are particularly desirable for dtra-long distance communication system and high-bit-rate fiber communications. A challenging possibility is to use fiber gratings for the creation of pulses of slow light, which is a topic of great current interest. A possible way to trap a zero-velocity soliton is to use an attractive finite-size or local defect [83] in BG. The interaction of the soliton with an attractive defect in the form of a local suppression of BG was studied recently in Refs [78,79].

4. Conclusions

My attempt on this article is to give a survey and update some of fiber Bragg grating solitons. There have been two papers for summarizing to Bragg solions [93] and gap solitons [20], gave readers insight into a series of working methods and results before these generalize. Clearly grating solitons have played an important role in past and ongoing nonlinear optical research in fiber Bragg grating, and we believe fiber Bragg grating solitons to have their greatest impact in the years to come.

5. Acknowledgements

This work was supported by the National Natural Science Foundation of China (No. 40571097) and the National Program on Key Basic Research Project (973 Program) (No. 2009CB724001)

REFERENCES

- K. O. Hill, Y. Fujii, D. C. Johnson, and B. S. Kawasaki, Appl. Phys. Lett., Vol. 32, pp. 647, 1978.
- [2] G. Meltz, W. W. Money, and W. H. Glenn, Optics Lett., Vol. 14, pp. 823, 1989.

- [3] R. Kashyap, J. R. Armitage, R. Wyatt, S. T. Davey, and D. L. Williams, Electron. Lett., Vol. 26, pp. 730, 1990.
- [4] K. O. Hill, B. Malo, F. Bilodeau, D. C. Johnson and J. Albert, Appl. Phys. lett., Vol. 62, pp. 1035, 1993.
- [5] L. Dong, J. L. Archambault, L. Reekie, *et al.* Electron. Lett., Vol. 29, pp. 1577, 1993.
- [6] E. Yablonovitch and T. J. Gmitter, Phys. Rev. Lett., Vol. 63, pp. 1950, 1989.
- [7] L. Brillouin, Wave propagation in periodic structures (New York: Dover), 1953.
- [8] Y. Kivshar and G. P. Agrawal, Optical solitons: from fibers to photonic crystals (New York: Academic Press), 2003.
- [9] D. N. Christodoulides and R. I. Joseph, Opt. Lett., Vol. 13, pp. 794, 1988.
- [10] D. N. Christodoulides, F. Lederer and Y. Silberberg Nature (London), Vol. 424, pp. 817, 2003.
- [11] A. A. Sukhorukov, Y. S. Kivshar, H. S. Eisenberg and Y. Silberberg, IEEE J. Quantum Electron, Vol. 39, pp. 31, 2003.
- [12] W. Chen and D. L. Mills, Phys. Rev. Lett., Vol. 58, pp. 160, 1987.
- [13] D. L. Mills and S. E. Trullinger, Phys. Rev. B, Vol. 36, pp. 947, 1987.
- [14] J. E. Sipe and H. G. Winful, Opt. Lett. Vol. 13, pp. 132, 1988.
- [15] C. M. de Sterke and J. E. Sipe, Phys. Rev. A, Vol. 38, pp. 5149, 1988.
- [16] D. N. Christodoulides and R. I. Joseph, Phys. Rev. Lett., Vol. 62, pp. 1746, 1989.
- [17] A. B. Aceves and S. Wabnitz, Phys. Lett. A, Vol. 141, pp. 37, 1989.
- [18] C. M. de Sterke, B. J. Eggleton and J. E. Sipe, Spatial Solitons, edited by S. Trillo and W. Torruellas (Berlin: Springer-Verlag), pp. 169, 2001.
- [19] C. M. de Sterke and J. E. Sipe, "Gap solitons" in Progress in Optics, E. Wolf, Ed. (Amsterdam, The Netherlands: Elsevier) Vol. XXXIII, ch. 3, pp. 203, 1994.
- [20] A. B. Aceves, Chaos 10, pp. 584, 2000.
- [21] J. P. Gordon, Opt. Lett., Vol. 8, pp. 596, 1983.
- [22] F. M. Mitschke and L. F. Mollenauer, Opt. Lett., Vol. 12, pp. 407, 1987.
- [23] J. S. Aitchison, A. M. Weiner, Y. Silberberg, D. E. Leaird, M. K. Oliver, J. L. Jackel and P. W. E. Smith, Opt. Lett., Vol. 16, pp. 15, 1991.
- [24] M. Shalaby, F. Reynaud and A. Barthelemy, Opt. Lett., Vol. 17, pp. 778, 1992.
- [25] H. G. Winful, Appl. Phys. Lett., Vol. 46, pp. 527, 1985.
- [26] K. Senthilnathan, P. Ramesh Babu, K. Porsezian, V. Santhanam and S. Gnanasekaran, Chaos: Solitons and Fractals, In Press, Corrected Proof, Available online, February

2006.

48

- [27] C. M. de Sterke and J. E. Sipe, Phys. Rev. A, Vol. 39, pp. 5163, 1989.
- [28] C. M. de Sterke and J. E. Sipe, Phys. Rev. A, Vol. 42, pp. 550, 1990.
- [29] C. M. de Sterke and J. E. Sipe, Phys. Rev. A, Vol. 42, pp. 2858, 1990.
- [30] C. M. de Sterke, K. R. Jackson and B. D. Robert, J. Opt. Soc. Am. B, Vol. 8, pp. 403, 1991,
- [31] H. G. Winful, J. H. Marburger and E. Gamire, Appl. Phys. Lett, Vol. 35, pp. 379, 1979.
- [32] G. Kurizki, A. E. Kozhenkin, T. Opatrny and B. A. Malomed, "Optical solitons in periodic media with resonant and off-resonant nonlinearities," in Progress in Optics, E. Wolf, Ed. (Amsterdam, The Netherlands: Elsevier) XXXXII, pp. 93, 2001.
- [33] I. V. Barashenkov, D. E. Pelinovsky and E.V. Zemlyanaya, Phys. Rev. Lett. Vol. 80, pp. 5117, 1998.
- [34] A. D. Rossi, C. Conti and S. Trillo, Phys. Rev. Lett., Vol. 81, pp. 85, 1998.
- [35] J. Feng and F. K. Kneubuhl, IEEE J. Quantum Electronics, Vol. 29, pp. 590, 1993.
- [36] J. Feng, Opt. Lett. Vol. 18, pp. 1302, 1993.
- [37] N. G. R. Broderick, C. M. d. Sterke and K. R. Jackson, Opt. Quantum Electron. Vol. 26, pp. S219, 1994.
- [38] A. B. Aceves, et al. Opt. Lett., Vol. 19, pp. 332, 1994.
- [39] V. E. Zakharov and A. B. Shabat, Soviet Phys. JETP, Vol. 34, pp. 62, 1972.
- [40] P. A. Bélanger and C. Paré, Phys. Rev. A, Vol. 41, pp. 5254. 1990.
- [41] D. J. Kaup and A. C. Newell, Lettere AL Nuovo Comento, Vol. 20, pp. 325, 1977.
- [42] D. Taverner, N. G. R. Broderick, D. J. Richardson, *et al.* Opt. Lett., Vol. 23, pp. 328, 1998.
- [43] C. Martijn de Sterke and J. E. Sipe, Phys. Rev. A, Vol. 43, pp. 2467, 1991.
- [44] A. B. Aceves, J. V. Moloney, and A. C. Newell, Phys Rev A, Vol. 39, pp. 1809, 1989.
- [45] N. G. R. Broderick and C. Martijn de Sterke, phys. rev. E, Vol. 51, pp. 4978, 1995.
- [46] N. G. R. Broderick and C. Martijn de Sterke, physical review E, Vol. 58, pp. 7941, 1998.
- [47] M. J. Steel and C. M. de Sterke, Phys. Rev. A Vol. 48, pp. 1625, 1993.
- [48] A. D. Capobianco, C. De Angelis, A. Laureti Palma and G. F. Nalesso, J. Opt. Soc. Am. B, Vol. 14, pp. 1956, 1997.
- [49] T. Iizuka and M. Wadati, J. Phys. Soc. Jpn. Vol. 66, pp. 2308, 1997.
- [50] T. Iizuka and C. M. de Sterke, Phys. Rev. E, Vol. 61, pp. 4491, 2000.
- [51] J. E. Sipe, L. Poladian and C. M. de Sterke, J. Opt. Soc.

Am. A, Vol. 11, pp.1307, 1994.

- [52] M. Quiroga Teixeiro and H. Michinel, J. Opt. Soc. Am. B, Vol. 14, pp. 2004, 1997.
- [53] M. L. Quiroga Teixeiro, A. Berntson and H. Michinel, J. Opt. Soc. Am. B, Vol. 16, pp. 1697, 1999.
- [54] A. Desyatnikov, A. Maimistov and B. Malomed, Phys. Rev. E, Vol. 61, pp. 3107, 2000.
- [55] D. Mihalache, D. Mazilu, L. C. Crasovan, B. A. Malomed and F. Lederer, Phys. Rev. E, Vol. 61, pp. 7142, 2000.
- [56] E. N. Tsoy, C. M. de Sterke, Phys. Rev. E, Vol. 62, pp. 2882, 2000.
- [57] E. N. Tsoy, C. M. de Sterke, Opt. Soc. Am. B, Vol. 18, pp. 1, 2001.
- [58] N. L. litchinitser, B. J. Eggleton and D. B. Patterson, J. Lightwave Technol, Vol. 15, pp. 1303, 1997.
- [59] J. Atai, B. A. Malomed, Phys. Lett. A, Vol. 284, pp. 247, 2001.
- [60] B. J. Eggleton, C. Martijn de Sterke, and R. E. Slusher, J. Opt. Soc. Am. B, Vol. 16, pp. 587. 1999.
- [61] B. J. Eggleton, R. E. Slusher, C. M. de Sterke, P. A. Krug, and J. E. Sipe, Phys. Rev. Lett., 76, pp. 1627, 1996.
- [62] C. M. de Sterke, B. J. Eggleton, and P. A. Krug, J. Lightwave Technol, Vol. 15, pp. 1494, 1997.
- [63] B. J. Eggleton, R. E. Slusher, T. A. Strasser, and C. M. deSterke, OSA Technical Digest Series (Washington, D.C.: Optical Society of America) 17 paper BMB1-1, 1997.
- [64] B. J. Eggleton, C. M. de Sterke, and R. E. Slusher, J. Opt. Soc. Am. B, Vol. 14, pp. 2980, 1997.
- [65] B. J. Eggleton, C. M. de Sterke, R. E. Slusher, A. Aceves, J. E. Sipe, and T. A. Strasser, Opt. Commun, Vol. 149, pp. 267, 1998.
- [66] G. Citation Lenz and B. J. Eggleton, J. Opt. Soc. Am. B, Vol. 15, pp. 2979, 1998.
- [67] S. Larochelle, V. Mizrahi, and G. Stegeman, Electron. Lett., Vol. 26, pp. 1459, 1990.
- [68] N. D. Sankey, D. F. Prelewitz and T. G. Brown, Appl. Phys. Lett. Vol. 60, pp. 1427. 1992,
- [69] N. G. R. Broderick, D. Taverner, D. J. Richardson, M. Ibsen and R. I. Laming, Opt. Lett. Vol. 22, pp. 1837, 1997.
- [70] D. Taverner, N. G. R. Broderick, D. J. Richardson, M. Ibsen, and R. I. Laming, Opt. Lett., Vol. 23, pp. 259, 1998.
- [71] N. G. R. Broderick, D. Taverner, D. J. Richardson, M. Ibsen and R. I. Laming, Phys. Rev. Lett., Vol. 79, pp. 4566, 1997.
- [72] A. Hasegawa and W. F. Brinkman, IEEE J. Quantum Electron. Vol. 16, pp. 694, 1980.
- [73] K. Tai, A. Hasegawa, and A. Tomita, Phys. Rev. Lett., Vol. 56, pp. 135, 1986.
- [74] G. P. Agrawal, Applications of Nonlinear Fiber Optics

(San Diego, CA: Academic), 2001.

- [75] H. He, A. Arraf, C. Martijn de Sterke, P. D. Drummond, and B. A. Malomed, Phys. Rev. E, Vol. 59, pp. 6064, 1999.
- [76] N. G. R. Broderick, D. Taverner, and D. J. Richardson Opt. Express, Vol. 3, pp. 447, 1998.
- [77] S. Radic, N. George, and G. P. Agrawal, Opt. Lett., Vol. 19, pp. 1789, 1994.
- [78] R. H. Goodman, R. E. Slusher, and M. I. Weinstein, J. Opt. Soc. Am. B, Vol. 19, pp. 1635, 2002.
- [79] W. C. K. Mak, B. A. Malomed, and P. L. Chu, J. Opt. Soc. Am. B, Vol. 20, pp. 725, 2003.
- [80] D. E. Pelinovsky, L. Brzozowski, and E. H. Sargent. Phys. Rev. E, Vol. 62, pp. 4536, 2000.
- [81] L. Brzozowski and E. H. Sargent, IEEE J. Quantum Electron, Vol. 36, pp. 550, 2000.
- [82] G. P. Agrawal, Nonlinear Fiber Optics. (New York: Academic), 1989.
- [83] K. T. Mc-Donald, Am. J. Phys., Vol. 68, pp. 293, 2000.
- [84] S. Lee and S. T. Ho, Opt. Lett., Vol. 18, pp. 962, 1993.
- [85] J. H. Lee, P. C. Teh, P. Petropoulos, M. Ibsen, and D. J. Richardson, IEEE Photon. Technol. Lett., Vol. 14, pp. 203,

2002.

- [86] J. H. Lee, L. Katsuo, K. S. Berg, A. T. Clausen, D. J. Richardson, and P. Jeppesen, J. Lightwave Technol, Vol. 21, pp. 2518, 2003.
- [87] J. H. Lee, Y. M. Chang, Y. G. Han, S. H. Kim, H. Chung, and S. B. Lee, IEEE Photon. Technol. Lett., Vol. 17, pp. 34, 2005.
- [88] N. Dogru and M. S. Ozyazici, International Workshops on Laser and Fiber-Optical Networks Modeling, (LFNM' 02: Kharkiv, Ukrania) june 2002, pp. 5.
- [89] N. Dogru and M. S. Ozyazici, International Workshops on Laser and Fiber-Optical Networks Modeling (LFNM' 04: Kharkiv, Ukrania) September 2004, pp. 119.
- [90] N. Dogru and M. S. Ozyazici, LFNM. September 2004, pp. 115.
- [91] N. Dogru and M. S. Ozyazici, International Conference on Indium Phosphide and Related Materials, 2005, May 2005, pp.291.
- [92] N. Dogru, 2005, IEEE NUSOD'05, September 2005, pp. 89.
- [93] T. G. Brown and B. J. Eggleton, Opt. Express Vol. 3, pp.385, 1998.



ADPF Algorithm for Target Tracking in WSN

Chunhe Song, Hai Zhao, Wei Jing, Dan Liu

Institute of Information and Technology, Northeastern University, Shenyang, China E-mail: songchunhe@tsinghua.org.cn, {zhhai, jingw, liud}@mail.neuera.com Received October 27, 2009; accepted December 29, 2009

Abstract: Particle filtering (PF) has been widely used in solving nonlinear/non Gaussian filtering problems. Inferring to the target tracking in a wireless sensor network (WSN), distributed PF (DPF) was used due to the limitation of nodes' computing capacity. In this paper, a novel filtering method—asynchronous DPF (ADPF) for target tracking in WSN is proposed. There are two keys in the proposed algorithm. Firstly, instead of transferring value and weight of particles, Gaussian mixture model (GMM) is used to approximate the posteriori distribution, and only GMM parameters need to be transferred which can reduce the bandwidth and power consumption. Secondly, in order to use sampling information effectively, when target moving to the next cluster head region, the GMM parameters are transfer to the next cluster head, and combine with the new local GMM parameters to compose the new GMM parameters incrementally. The ADPF can also deal with the situation of different number of nodes in different cluster when using the dynamic cluster structure. The proposed ADPF is compared to some other DPF for WSN target tracking, and the experimental results show that not only the precision is improved, but also the bandwidth and power is reduced.

Keywords: WSN, target tracking, asynchronous distributed particle filtering

1. Introduction

One of the major goals of WSN is to detect and track changes. The problem concerned is performing on-line state estimation for multi-dimensional signals that can be modeled using markovian state-space models that are nonlinear and non-Gaussian, Particle filter is one of the widely used tracking algorithms in non-linear/ Gaussian dynamic systems. When using such algorithm in sensor networks the energy cost related to computation in each sensor node and communication between sensor nodes is significant. Currently there are several distributed particle filters [1-3], in which the distributed nature is achieved by either transmitting local statistics of particles to a centralized unit or using the parameters passing method. Transmitting local statistics of particles to a centralized unit is not an efficient approach. Failure of the centralized unit is vital to the entire network. In the parameters passing method, the algorithms construct a path through the networks, which passes through all nodes. Global statistics of particles are accumulated by adding local statistics in each node through a forward pass. Then there needs a backward pass, which runs the important sampling and selection steps in each sensor node by using the accumulated global statistics.

In this paper, a novel filtering method – asynchronous DPF (ADPF) for target tracking in WSN is proposed. There are two keys in the proposed algorithm. Firstly, instead of transferring value and weight of particles,

Gaussian mixture model (GMM) is used to approximate the posteriori distribution, and only GMM parameters need to be transferred which can reduce the bandwidth and power consumption. Secondly, in order to use sampling information effectively, when target moving to the next cluster head region, the GMM parameters are transfer to the next cluster head, and combine with the new local GMM parameters to compose the new GMM parameters incrementally. Because of computing asynchronously, this process can be regarded as ADPF.

The remaining of the paper is organized as follows: a brief description of PF and DPF in WSN are presented in Section 2. The details of the new PF this paper proposed – ADPF is presents in Section 3. In Section 4 the proposed algorithm is compared to other DPFs and finally, we give some concluding remarks in Section 5.

2. Tracking in WSN

Issues considered in this paper is tracking a moving target based on position measurements from multiple distributed sensors.

2.1 Target Motion Model

The target motion model used in this paper is the same as [4]:

$$x_k = f(x_{k-1}) + Gv_k, \quad k = 1, 2, \dots$$
 (1)

where the target state vector $x_k = [x, \dot{x}, y, \dot{y}, w]^T$, consist of the position, velocity and the turn rate; $v_k \sim N(0, Q_k)$ is white process nosie: and

$$f(x) = \begin{bmatrix} 1 & \frac{\sin(\varepsilon_k T)}{\varepsilon_k} & 0 & -\frac{(1-\cos(\varepsilon_k T))}{\varepsilon_k} & 0\\ 0 & \cos(\varepsilon_k T) & 0 & \sin(\varepsilon_k T) & 0\\ 0 & \frac{(1-\cos(\varepsilon_k T))}{\varepsilon_k} & 1 & \frac{\sin(\varepsilon_k T)}{\varepsilon_k} & 0\\ 0 & \sin(\varepsilon_k T) & 0 & \cos(\varepsilon_k T) & 0\\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x\\ \dot{x}\\ y\\ \dot{y}\\ w \end{bmatrix}$$
(2)
$$Gk = \begin{bmatrix} T^2/2 & 0 & 0\\ T & 0 & 0\\ 0 & T^2/2 & 0\\ 0 & T & 0\\ 0 & 0 & T \end{bmatrix}$$
(3)

2.2 Target Motion Model

In this paper, measurements of range and bearing are given by [5]:

$$z_k = h(x_k) + w_k^i, \quad i = 1, 2, ..., N$$
 (4)

with

$$h(x_k) = \begin{bmatrix} ri\\bi \end{bmatrix} = \begin{bmatrix} \sqrt{x_k^2 + y_k^2}\\\tan^{-1}(\frac{y_k}{x_k}) \end{bmatrix}$$
(5)

And white measurement noise $w_k^i \sim N(0, R_k^i)$.

3. Basic PF and DPF in WSN

3.1 Basic Particle Filter

In the bayes filtering framework, the posterior distribution is updated recursively over the current state x_t given all observations $Z_t = \{z_i\}_{i=1}^t$ up to and including time *t* as follows[6]:

$$p(x_t | Y_{t-1}) = \int p(x_t | x_{t-1}) p(x_{t-1} | Y_{t-1}) dx_k$$
(6)

$$p(x_t | Y_t) = \frac{p(y_t | x_t) p(x_t | Y_{t-1})}{p(y_t | Y_{t-1})}$$
(7)

$$p(y_t | Y_{t-1}) = \int p(y_t | x_t) p(x_t | Y_{t-1}) dx_k$$
(8)

Using Monte Carlo sampling points, particle filter executes the filtering process by generating weighted sampling points of state variances recursively. Generic particle filter algorithm can be found in [4].

3.2 Distributed Particle Filter in WSN

Particle filter has been widely used in target tracking. In WSN, the information transferred between nodes and the computing process in nodes are limited due to the restriction of power, computing capacity and bandwidth, so some changes must be conducted in particle filter in order to use it in WSN target tracking. The main idea of distributed particle filter is to deal with the computation at different sensor rather than at a central unit.

One typical DPF is the forward-backward type of distributed particle filter [6], which may be the first DPF for WSN. Supposing K nodes in WSN, and N particles for each node, in this algorithm, firstly, it is assumed that measurements at sensor are independent with each others, and in the particle filtering process, the likelihood $p(y_i | x_i)$ can be approximated by a parametric model $p(y_t | x_t) = L(x_t; \theta_t^k |_{k=1}^K)$. Secondly, only a single communication chain exists from node 1 to node k, with any node i in the interior of the chain communicating only with nodes i-1 and i+1. In the initialization step, each node samples N particles from $p(x_0)$. At time t, Node i sample from its importance distribution to generate N particles. Node i calculates the value of its likelihood for each one of these particles for the current observ ation and then trains the model $\{(\tilde{x}_t^{(j)}, L(x_t; \{\theta_t^k\}_{k=1}^i)\}$ $p(y_t^k | \tilde{x}_t^{(j)}))_{i=1}^N$. The parameters $\{\theta_t^k\}_{k=1}^i$ are then appropriately quantized and transmitted to node i+1 in the chain. In the next phase of this algorithm, the estimated parameters $\{\theta_t^k\}_{k=1}^K$ are propagated back along the communication chain. And each node uses the parameters to calculate estimates of likelihood for its samples $\{\tilde{x}_t^{(j)}\}_{j=1}^N$, which will be used to calculate the importance weights of particles. In the mentioned process above, it is assumed that the sensors act synchronously and record their measurements at the same time. So it can be regarded as synchronously particle filter.

There is complicated training process in this algorithm, and all nodes compute synchronously during target tracking. One simple idea is to transfers value and weight of particles directly between nodes and represents the posterior distribution. But there are N particles in each node and the transferred bits will be very large, so the posterior distribution of particle filter is assumed to be a GMM with C mixture probabilities. [2] is such type of DPF. In this algorithm, firstly, the WSN is divided into a series of group misrelated, and a single particle filter runs in each group. Through the head of current group, parameters of filter are transferred to the next head and update the posterior distribution. On the last group of sensors target tracking was estimated. As need transfer number of value and weight of particle, in this algorithm, low dimension GMM is used to approximate the likelihood distribution of DPF. To implement a distributed particle filter (DPF), particles and weights are distributed over entire network. Each sensor should maintain N particles $x_{m,k}^{(n)}$ and weights $W_{m,k}^{(n)}$. The posterior distribution of particle filter is assumed to be a GMM with C mixture

probabilities. For each unobserved state y_c , observation $z_{m, k}$ follows a Gaussian distribution with mean μ_c and variance Σ_c :

$$p(z_{m,k} \mid \mu_{c}, \Sigma_{c}) = \frac{1}{\sqrt{2\pi \|\Sigma_{c}\|}} e^{-\frac{1}{2}(z_{m,k} - \mu_{c})^{T} \Sigma_{c}^{-1}(z_{m,k} - \mu_{c})}$$
(9)

The Gaussian mixture distribution for observation $z_{m,k}$ is:

$$p(z_{m,k} \mid \theta) = \sum \alpha_{m,c} p(z_{m,k} \mid \mu_c, \Sigma_c)$$
(10)

where θ is the set of the distribution parameters to be estimated, $\theta = \{\alpha_{m,c}, \mu_c, \Sigma_c; c = 1, ..., C, m = 1, ..., M\}$

Assume all observed data from all nodes are sent to a centralized unit where a standard EM algorithm is used to estimate the parameter set θ . The log-likelihood for the observed data satisfies:

$$L(\theta \mid z) = \log \prod_{m=1}^{M} \prod_{j=1}^{k} p(z_{m,j} \mid \theta) = \sum_{m=1}^{M} \sum_{j=1}^{k} p(z_{m,j} \mid \theta)$$
$$= \sum_{m=1}^{M} \sum_{j=1}^{k} \alpha_{m,c} p(z_{m,k} \mid \mu_{c}, \Sigma_{c})$$
(11)

4. Asynchronous Distributed Particle Filter

4.1 Dynamic Cluster Structure

In some former DPF, the cluster structure is fixed. In the ADPF, dynamic cluster structure is used. Supposing all nodes have the same detecting capacity, choose one cluster head, forming all nodes in the range of sing-hop of cluster head into a cluster, and this cluster head is used to deal with the sampling data and get local estimation. Supposing C is the max distance of sing-hop communication, R is the max range of detecting, and D is the distance of cluster head and other node. When target entering into the detecting region of WSN, and the number of node already detecting the target is over a predefined threshold, choose the nearest node as the cluster head. Forming the cluster and recall all nodes in this cluster. Following the movement of target, some nodes in the cluster will be out of the detecting region. If the number of nodes out of the detecting region is over a predefined threshold or D+C>R, choose a new cluster head and construct a new cluster. Otherwise predict the new location of target using filtering algorithms.

4.2 Gaussian Mixture Model Using EM

Using EM algorithm, the parameters of Equation 7) can be calculated. Given observation z and current parameter set θ^{t} where t is the time step between two consecutive sensor observations at k and k+1, the conditional expectation of joint distribution $p(z, y | \theta)$ is defined as:

$$Q(\theta, \theta^{t}) = \sum_{c=1}^{C} \sum_{m=1}^{M} \sum_{j=1}^{k} \log[p(z_{m,k}, y_{c} \mid \theta))] p(y_{c} \mid z_{m,j}, \theta^{t})$$
$$\sum_{c=1}^{C} \sum_{m=1}^{M} \sum_{j=1}^{k} \log[\alpha_{m,k} p(z_{m,k} \mid \mu_{c}, \Sigma_{c})] p(y_{c} \mid z_{m,j}, \theta^{t})$$

Detail information about the Gaussian mixture model using EM algorithm can be found in [6].

4.3 Asynchronous Updating GMM Parameters

When running DPF in WSN asynchronously using GMM approximate posterior distribution, after getting into the new cluster region, former sampling information will be lost. Here propose a new GMM incrementally updating algorithm using PCA concept.

Supposing m0 nodes in the first cluster, each node send its parameters $\theta = [\alpha_{m,c}, \mu_c, \Sigma_c]^T$ to the first cluster head, and compose the parameters matrix $P_0 = [\theta_1, ..., \theta_{n0}]$. It will be used to approximate the posterior distribution and transfer to the next cluster head.

Supposing the former cluster is the (i-1)th cluster with n_{i-1} sensors, parameters matrix is P_{i-1} , the current cluster is the ith cluster with n_i sensors, parameters matrix is $P_i \cdot \overline{P}_{i-1}$ and \overline{P}_i are the raw average vectors of P_{i-1} and P_i .

decomposing P_i with SVD:

$$P_i = U\Sigma V^T \tag{12}$$

and denote \overline{P}_i as the row average vector of P_i .

Using the P_i and P_{i-1} to compose a new matrix $P_i^* = [P_i, P_{i-1}]$, decomposing it with SVD:

$$P_i^* = [P_i, P_{i-1}] = \tilde{U}\tilde{\Sigma}\tilde{V}^T$$
(13)

where $\tilde{U} = [U, U^*], \tilde{\Sigma} = \begin{bmatrix} \Sigma & U^T P_{i-1} \\ 0 & U^{*T} P_{i-1} \end{bmatrix}$ and $\tilde{V}^T = \begin{bmatrix} V^T & 0 \\ 0 & I \end{bmatrix};$

$$P_{ix} = [P_i - \overline{P}_i | sqrt(\frac{n_i * n_{i-1}}{n_i + n_{i-1}}) * (\overline{P}_i - \overline{P}_{i-1})] \quad (14)$$

 $P_{iy} = P_{ix} * (I - U * U^T)$, calculate the QR decomposing of $P_{iy} : P_{iz} = QR(P_{iy})$. Using V^T , U, P_{ix} and P_{iz} to compose the matrix \tilde{P}_i :

$$\tilde{P}_{i} = \begin{bmatrix} ff * V & UT * P_{ix} \\ 0 & P_{iy}^{T} * P_{ix} (I - U * U^{T}) \end{bmatrix}$$
(15)

where ff is the forgotten factor with value in the range of [0, 1]. It indicates the weights of last parameters in the current computing time. In this experiment, the ff is set to 0.8.

Calculating SVD of \tilde{P}_i :



Figure 1. True states and observations



Figure 2. Results of mean errors

$$\tilde{P}_{i} = \tilde{U}\tilde{\Sigma}\tilde{V}^{T} \tag{16}$$

5. The Simulation Experiments

Ì

In order to test the proposed algorithm, the ADPF is compared to the DPF algorithms in [1] and [2]. Experimental results are shown in Figure 1–Figure 2 and Table 1. All results are the means of 100 runs.

As shown in the experimental results, it is clear that, the proposed ADPF has better performance than other two typical DPF algorithms. It can be explained as the proposed algorithms can use the sampling information as incremental updating GMM parameters more effectively.

6. Conclusions

Synchronous DPF and GMM parameters transferred DPF have their own disadvantages which limit their using range. In this paper, a novel filtering method – asynchronous DPF (ADPF) for target tracking in WSN is proposed. With incremental updating GMM parameters, ADPF can use the sampling information effectively. And ADPF can also deal with the situation of different number of nodes in different cluster when using the dynamic cluster structure. Simulation result shows that ADPF has better performance than other two typical DPF algorithms.

REFERENCES

- D. Guo and X. Wang, "Quasi-monte carlo filtering in nonlinear dynamic systems," IEEE Trans. Signal Process, Vol. 54, No. 6, pp. 2087–2098, 2006.
- [2] M. S. Arulampalam, S. Maskell, N. Gordon, *et al.* "A tutorial on particle filters for online nonlinear/non- gaussian bayesian tracking [J]," IEEE Trans on Signal Proceeding, Vol. 20(2), pp. 174–188, 2002.
- [3] D. Crisan and A. Doucet, "A survey of convergence results on particle filtering methods for practitioners [J]," IEEE Trans on Signal Proceeding, Vol. 50(3), pp. 736– 746, 2002.
- [4] B. D. Anderson and J. B. Moore, "Optimal filtering," Prentice-Hall, New Jersey, 1979.
- [5] X. R. Li and V. P. Jilkov, "Survey of maneuvering target tarcking part I: dynamci models," in IEEE Trans. Aerospace and Electronic System, Vol. 39, 2003.
- [6] X. R. Li and V. P. Jilkov, "A survey of maneuvering target tracking-part III: measurement models," in SPIE Conf. on Signal and Data Proceeding of Small Target, 2001.
- [7] M. Coates, "Distributed particle filters for sensor networks," in Proceeding of 3rd Intl sysmosium on Information Proceeding in sensor networks, Berkely, CA, USA.
- [8] S. J. Julier and J. K. Uhlmann, "A new extension of the Kalman filter to nonlinear systems," Proceeding of Aero-Sense: The 11th International Sysmpsium on Aerospace/ Defence Sensing, Simulation and Controls, Orlando, Florida, 1997. Vol. Muti Sensor Fusion, Tracking and Resource Mangement II pp.182–193.
- [9] Y. Shi and R. C. Eberhart "A modified particle swarm optimizer," In Proceedings of the IEEE International Conference on Evolutionary Computation. Piscataway, NJ: IEEE Press, pp. 69–73, 1998.
- [10] J. Riget, "A diversity-guided particle swarm optimizer," the ARPSO. EVALife Technical Report 2002–02, Dept. of Computer Science, University of Arhus, 2002.
- [11] D. Guo, X. Wang, and R. Chen, "New sequential monte carlo methods for nonlinear dynamic systems," Statistics and Computing, Vol. 15, No. 2, pp. 135–147, 2005.
- [12] Y. Bar-Shalom and X. R. Li, "Kirubarajan T. Estimation with applications to tracking and navigation: theory, algorithm and software [M]," New York: Wiley, 2001.

Designing Intrusion Detection System for Web Documents Using Neural Network

Hari Om, Tapas K. Sarkar

Department of Computer Science and Engineering, Indian School of Mines, Dhanbad, India E-mail: hariom63@rediffmail.com, aastitva@gmail.com Received November 17, 2009; accepted December 29, 2009

Abstract: Cryptographic systems are the most widely used techniques for information security. These systems however have their own pitfalls as they rely on prevention as their sole means of defense. That is why most of the organizations are attracted to the intrusion detection systems. The intrusion detection systems can be broadly categorized into two types, Anomaly and Misuse Detection systems. An anomaly-based system detects computer intrusions and misuse by monitoring system activity and classifying it as either normal or anomalous. Misuse detection systems can detect almost all known attack patterns; they however are hardly of any use to detect yet unknown attacks. In this paper, we use Neural Networks for detecting intrusive web documents available on Internet. For this purpose Back Propagation Neural (BPN) Network architecture is applied that is one of the most popular network architectures for supervised learning. Analysis is carried out on Internet Security and Acceleration (ISA) server 2000 log for finding out the web documents that should not be accessed by the unauthorized persons in an organization. There are lots of web documents available online on Internet that may be harmful for an organization. Most of these documents are blocked for use, but still users of the organization try to access these documents and may cause problem in the organization network.

Keywords: intrusion detection system, neural network, back propagation network, anomaly detection, misuse detection

1. Introduction

The information is the most important resource that must be managed efficiently. Besides management, its protection is also very important as it may lead to economic losses in today's electronic environment. For example, we can control our bank accounts from almost anywhere in the world using a suitable network, such as satellite and cellular phone networks to interact with the bank representatives, or the specialized wired ATM networks and the Internet for online banking services. The services supported by networks are very much useful and efficient, but these can be subverted by unscrupulous elements for their own benefits. So, suitable mechanism needs be employed to protect the information. In a survey of fraud against auto teller machines [1], it is reported that the patterns of fraud depends on those who were responsible for implementing and managing the systems. In USA, if a customer disputes a transaction, this is the responsibility of the bank to prove that the customer is mistaken or lying. This forced the US banks to protect their systems properly. But, in Britain, Norway and the Netherlands, the burden of proof lies on the customer. The bank is right if the

customer could not prove it wrong. That is why the banks in these countries became careless. Eventually, epidemics of fraud demolished their satisfaction and in the meanwhile the US banks suffered much less fraud. Though they spent less money on security than their European counterparts, yet they spent it more effectively [2]. A different kind of incentive failure was also seen in early 2000 with distributed denial of service attacks against a number of high profile websites. Those attacks exploited a number of weak machines to launch a large coordinated packet flood at a host. Since many of them flooded the victim at the same time, the traffic was more than the host could handle. Furthermore, because it came from many different sources, it could be very difficult to stop. Varian [3] discusses different kind attacks and their effects. The suggestions made in [3] are: the costs of distributed denial-of-service attacks should fall on the operators of the networks from which the flooding traffic originates. And assign legal liability to the parties that are best able to manage the risk as they will develop expertise for computer security and provide the required services to their clients. In next section we review the intrusion detection systems.



2. Early Intrusion Detection System

An intrusion occurs when an attacker gains unauthorized access to a valid user's account and performs disruptive behavior while masquerading as that user. The attacker may harm the user's account directly or can use it to launch attacks on other accounts or machines. In such scenario a useful method to detect it is to develop "patterns" of users of a computer system. The early intrusion detection efforts used to do manual review of a system audit trail that was inefficient approach as many systems did not collect enough data to provide an audit trail, or failed to protect the data against modification. Studies in [4] show that nearly all large corporations and most medium-sized organizations have installed some form of intrusion detection tool. In [5], the misuse detection methods using mobile agents are discussed. The methods to detecting intrusions can be anomaly detection or misuse detection. Misuse detection is mainly suitable for reliably detecting known patterns, but they are hardly of any use yet unknown attack methods. The mobile agents provide computational security by constantly moving around the Internet and propagating rules to solve misuse detection. The paper [6] discusses an Intrusion Detection System (IDS) architecture integrating both anomaly and misuse detection approaches. This architecture consists of three main modules: an anomaly detection module, a misuse detection module, and a decision support system module. The anomaly detection module uses a Self-Organizing Map (SOM) structure to model normal behavior and any deviation from the normal behavior is considered as an attack. The misuse detection module uses J.48 decision tree algorithm to classify different types of attacks. The decision support system analyzes and interprets the results for interpreting the results of both anomaly and misuse detection modules. In [7], strict anomaly detection method is discussed that uses the neural networks to a great effect. Now we review the important approaches used in the intrusion detection systems.

2.1 Rule Based Intrusion Detection Systems

The basic assumption in the rule-based intrusion detection systems is that the intrusion attempts can be characterized by sequences of user activities that lead to compromised system states and based on that they predict intrusion. These systems fire rules when audit records or system status information begins to indicate illegal activity. Two major approaches are followed in rule-based intrusion detection: state-based and model-based approach. In the former, the rule base is codified using the terminology found in the audit trails and Intrusion attempts are the sequences of system state as defined by audit trail information leading from an initial and limited access state to a final compromised state [8]. In the later, the known intrusion attempts are modeled as sequences of user behavior. The intrusion detection system itself is responsible for determining how an identified user behavior may manifest itself in an audit trail. These systems have many benefits, such as large data processing, more intuitive explanations of intrusion attempts, and prediction of future actions. The rule-based systems however have some limitations. They lack flexibility in the rule-to-audit record representation. Slight variations in an attack sequence can affect the activity-rule comparison up to that extent that the intrusion may not be detected. While increasing the level of abstraction of the rule-base does provide a partial solution to this weakness, it also reduces the granularity of the intrusion detection device. A number of non-expert system-based approaches to intrusion detection have been discussed in [9-12]. Most current approaches to detecting intrusions utilize some form of rule-based analysis. Expert systems are the most common form of rule-based intrusion detection approaches [13-16]. An Expert system consists of a set of rules that encode the knowledge of a human "expert". These rules are used by the system to make conclusions about the security-related data from the intrusion detection system. Unfortunately, the expert systems require frequent updates to remain current. While the expert systems offer an enhanced ability to review audit data, the required updates may be ignored or performed infrequently by the administrator. At a minimum, this leads to an expert system with reduced capabilities. At worst, this will degrade the security of the entire system by causing the system's users to be mislead into believing that the system is secure, even as one of the key components becomes increasingly ineffective over the time.

2.2 Network-Based and Host-Based Intrusion Detection Systems

A network-based intrusion detection system (NIDS) observes the traffic at specified points in the network and then checks that traffic packet by packet in real time to detect intrusion patterns. It can examine the activity at any layer of the network such as network layer, transport layer, and application layer protocol. The network-based systems are generally best at detecting the unauthorized outsider access and bandwidth theft/denial of service. When an unauthorized user logs in successfully, or attempts to log in, they are tracked with host-based IDS. However, detecting the unauthorized users before their logon attempt is best accomplished with network-based IDS. The packets that initiate bandwidth theft attacks can best be noticed with use of network-based IDS. Some of the network-based IDS are Shadow, Dragon, NFR, RealSecure, and NetProwler.

Host-based Intrusion Detection systems are first of IDSs developed and implemented. They collect and analyze the data originated on a computer that provides a

service, such as web server. After collecting the data from a given computer, it is analyzed. One example of the host-based system is programs that operate on a system and receive application or operating system audit logs. These programs are highly effective for detecting insider abuses. Residing on the trusted network systems themselves, they are close to the network's authenticated users. If one of these users attempts an unauthorized activity, the host-based systems usually detect and collect the most pertinent information in the quickest possible manner. In addition to detecting unauthorized insider activity, the host-based systems are also effective at detecting unauthorized file modification. The host-based IDSs are Windows NT/2000 Security Event Logs, RDMS audit sources, Enterprise Management systems audit data (such as Tivoli), and UNIX Syslog in their raw forms.

Graph-Based Intrusion Detection System (GrIDS) [17] uses a graphical representation to monitor the activity of entire network. EMERALD eXpert-BSM, a real-time forward-reasoning expert system, uses a knowledgebase to detect multiple forms of system misuse [18]. In [19], a technique is discussed for detecting intrusions at the level of privileged processes. It is reported that short sequences of system calls executed by running programs are a good discriminator between normal and abnormal operating characteristics of several common UNIX programs. Analyzing the system calls made by a program is a reasonable approach to detect intrusions based on program behavior profiles [20].

2.3 Neural Network Based Intrusion Detection Systems

The neural network based intrusion detection systems have the ability to be trained and learn patterns in a given environment, which can be used to detect intrusions by recognizing patterns of an intrusion. The Artificial Neural Network based methods for intrusion detection are quite popular. Recently an investigation on the unsupervised neural network models and choice for most appropriate one among them for evaluation and implementation is discussed in [21]. These can be used for both host-based and network based intrusion detection systems. For the success of IDS is the failure of firewalls to prevent many security intrusions. The intrusion detection systems can detect many of them that slip through firewalls. Many Anomalies based and Misuse based intrusion detection techniques have been designed to detect the abnormal behavior exhibited by the user in [22-27]. Artificial neural networks have been suggested as alternatives to the statistical analysis [28-30]. Statistical Analysis involves statistical comparison of current events to a predetermined set of baseline criteria. Neural networks are specifically discussed to identify the typical characteristics of system users and identify statistically

significant variations from the user's established behavior. Artificial neural networks have also been discussed for use in the detection of computer viruses. In [31], neural networks are discussed as statistical analysis approaches in the detection of viruses and malicious software in computer networks. The neural network intrusion detection (NNID) system [32] uses neural networks to predict the next command a user will enter based on previous commands. Now we discuss our neural network based intrusion detection system.

3. Audit Logs Analysis Using Neural Networks

In this work, we collect the data from the ISA 2000 Web Access Log to analyze for possible intrusion attacks using the neural networks and then use the back propagation neural (BPN) network model for analyzing the input data. Different numbers of hidden layers are considered in the PBN algorithm.

3.1 ISA 2000 Web Access Log Analysis

Internet bandwidth is consumed by a variety of internet application protocols. The most popular application layer protocol that accesses Internet resources is the HTTP protocol. It is used to access the resources on the World Wide Web. Although bandwidth cost per-kilobyte or per-megabyte has come down over the years, yet the amount of bandwidth consumed by users on the campus network increases year after year. HTTP connections to Internet resources not only lead to increase in bandwidth usage, they also reduce the amount of bandwidth available on the Internet link for other important protocols and applications, such as SMTP, POP3 and VPN. In order to provide the desired data resources to users, it is stored at different locations using some kind of servers. To further help the user in computer network environment, proxy servers are employed. A proxy server is a server (a computer system or an application program) which provides the services to user requests by making requests to other servers. A user connects to the proxy server, requesting a file, connection, web page, or other resource available from a different server. In an enterprise that uses the Internet, a proxy server is a server that acts as an intermediary between a workstation user and the Internet so that the enterprise can ensure security, administrative control, and caching service. It can receive a request for an Internet service (such as a Web page request) from a user. On clearing filtering requirements, the proxy server, assuming it is also a cacheserver, looks in its local cache of previously downloaded Web pages. If the desired pages are there, it returns them to the user without needing to forward the request to the Internet. In case the required pages are not in the cache, the proxy server, acting as a client on behalf of the user, uses one of its own IP addresses to request the pages

5	7
э	1

Field name	Description
c-ip	The Internet Protocol (IP) address of the requesting client.
cs-username	The account of the user making the request. If ISA Server access control is not being used, ISA Server uses Anony- mous.
c-agent	The name and version of the client application sent by the client in the <i>Hypertext Transfer Protocol (HTTP)</i> User-Agent header.
date	The date on which the logged event occurred.
time	The local time when the logged event occurred.
r-host	The domain name for the remote computer that provides service to the current connection.
r-ip	The network IP address of the remote computer that provides service to the current connection.
r-port	The reserved <i>port number</i> on the remote computer that provides service to the current connection.
time-taken	The total time, in milliseconds, that is needed by ISA Server to process the current connection
cs-bytes	The number of bytes sent from the remote computer and received by the client during the current connection.
sc-bytes	The number of bytes sent from the client to the remote computer during the current connection.
cs-protocol	The application protocol used for the connection. Common values are http for Hypertext Transfer Protocol, https for Secure HTTP, and ftp for <i>File Transfer Protocol</i> .
s-operation	The HTTP method used. Common values are GET, PUT, POST, and HEAD.
cs-uri	The URL requested.
s-object-source	The type of source that was used to retrieve the current object. A table of some possible values is provided in Object Source Values.
sc-status	A Windows (Win32) error code (for values less than 100), an HTTP status code (for values between 100 and 1,000), a Winsock error code (for values between 10,004 and 11,031), or an ISA Server error code.

Table 1. Attributes in ISA server 2000 log file

from the server out on the Internet. When the pages are received, the proxy server forwards them onto the user.

3.2 ISA Server 2000 Web Access Log

Internet Security and Acceleration (ISA) Server 2000 can help in reducing overall bandwidth usage and cost by caching Web contents on the ISA Server 2000. We use Microsoft ISA Server 2000 log to monitor and analyze the status of the Web proxy requests to find out the documents that are worthless in an organization. Table 1 shows the attributes used in ISA Server 2000 Log file.

3.3 Experiment

The input data is collected in terms of above mentioned attributes. Table 2 contains the values of the input data.

The data shown in Table 2 is not a valid input pattern

for BPN. Before providing the data for training to the BPN, it needs be converted in the valid pattern. We perform the following steps for making a valid input for BPN.

• Select the ip address part of the destination web server and convert it in the integer number without delimiter. For example, the ip 216.239.63.83 is converted into 2162396383. This is a long number which in itself is not a valid input pattern for BPN.

• Normalize the input pattern in real numbers. After normalization the input data pattern is shown in Table 3. First column shows the normalized ip addresses and the second column shows 1 as valid ip address and 0 as invalid ip addresses.

• Train the BPN for this input pattern by taking dif ferentnumber of hidden layers. We use 2, 5 and 10 hidden layers. The number of epochs is taken as 50,000. Results

Table 2. ISA server 2000 web access log

	CE HEAT							Time	05	60	CE	6.00		s objece	80
c-ip	name	c-agent	date	time	r-host	r-ip	r-port	-ta	bytes	bytes	protocol	eration	cs-uri	-source	stauts
10.0.4.36	anonymous	Mozilla/4	12/14/2006	7:01.41	Images3.0	72.14.209	80	797	796	3053	http	GET	http://image	VCache	30
10.0.4.46	Anonymous	Mozilla/5	12/14/2006	7:01.41	www.orku	72.14.209	80	797	981	253	http	GET	http://www	Inet	30
10.0.4.46	Anonymous	Mozilla/5	12/14/2006	7:01.41	Images3.0	72.14.209	80	813	995	253	http	GET	http://image	VCache	30
10.0.14.23	Anonymous	Mozilla/4	12/14/2006	7:01.41	Immail.re	210.161.32	80	640	1243	237	http	GET	http://image	Inet	30
10.0.7.221	Anonymous	Mozilla/4	12/14/2006	7:01.41	In.f89.mail	203.84.222	80	5844	2332	79644	http	POST	http://in.f89	Inet	20
10.0.4.123	Anonymous	Mozilla/5	12/14/2006	7:01.41	www.orku	72.14.209	80	3109	1135	7267	http	GET	http://www	Inet	20
10.0.4.165	Anonymous	Mozilla/4	12/14/2006	7:01.41	Jdelivery	210.161.32	80	593	1014	277	http	GET	http:// jesliv	Inet	30
10.0.98.43	Anonymous	Mozilla/4	12/14/2006	7:01.41	In.wrs.yal	216.252.12	80	1359	816	601	http	GET	http://in,wrn	Inet	30
10.0.4.185	Anonymous	Mozilla/4	12/14/2006	7:01.41	Mum.inte	220.226.20	80	4531	358	2312	http	GET	http:// mum	Inet	00
10.0.4.46	Anonymous	Mozilla/5	12/14/2006	7:01.41	Imagas3.0	72.14.209	80	797	995	253	http	GET	http://image	VCache	30
10.0.4.46	Anonymous	Mozilla/5	12/14/2006	7:01.41	Imagas3.0	72.14.209	80	781	1000	253	http	GET	http://image	VCache	30
10.0.4.36	Anonymous	Mozilla/4	12/14/2006	7:01.41	Imagas3.0	72.14.209	80	766	796	2257	http	GET	http://image	VCache	30
10.0.4.36	Anonymous	Mozilla/4	12/14/2006	7:01.42	Imagas3.0	72.14.209	80	781	796	2215	http	GET	http://image	VCache	30
10.0.4.179	Anonymous	Mozilla/4	12/14/2006	7:01.42	www.goo	72.14.235	80	859	969	1532	http	GET	http://www	Inet	20
10.0.4.163	Anonymous	Mozilla/4	12/14/2006	7:01.42	Images3.0	72.14.209	80	1563	747	1882	http	GET	http://image	Inet	20
10.0.4.36	Anonymous	Mozilla/4	12/14/2006	7:01.42	www.orku	72.14.209	80	5312	968	18229	http	GET	http://www	Inet	20
10.0.4.36	Anonymous	Mozilla/4	12/14/2006	7:01.42	Images3.0	72.14.209	80	797	994	2281	http	GET	http://image	VCache	30
10.0.4.54	Anonymous	Mozilla/4	12/14/2006	7:01.42	www.orku	72.14.209	80	3953	1013	18507	http	GET	http://www	Inet	20
10.0.4.46	Anonymous	Mozilla/5	12/14/2006	7:01.42	Images3.0	72.14.209	80	796	1014	201	http	GET	http://image	VCache	30
10.0.4.46	Anonymous	Mozilla/5	12/14/2006	7:01.42	Images3.0	72.14.209	80	812	1008	201	http	GET	http://image	VCache	30
10.0.4.165	Anonymous	Mozilla/4	12/14/2006	7:01.42	jdelivery	210.161.32	80	594	1030	276	http	GET	http://jdeliv	VCache	30
10.0.4.39	Anonymous	Mozilla/4	12/14/2006	7:01.42	www2.nu	69.25.142	80	133125	1683	1119	http	GOST	http://www	Inet	6
10.0.4.174	Anonymous	Mozilla/4	12/14/2006	7:01.42	Mail.goog	209.85.139	80	2563	1683	361	http	GET	http://mail	Inet	20
10.0.4.193	Anonymous	Mozilla/4	12/14/2006	7:01.42	www.go	72.14.235	80	703	340	234	http	GET	http://www	VCache	30
10.0.4.46	Anonymous	Mozilla/5	12/14/2006	7:01.42	Images3.0	72.14.209	80	781	1013	201	http	GET	http://image	VCache	30
10.0.4.46	Anonymous	Mozilla/5	12/14/2006	7:01.42	Images3.0	72.14.209	80	797	1014	201	http	GET	http://image	VCache	30
10.0.4.36	Anonymous	Mozilla/4	12/14/2006	7:01.42	Images3.0	72.14.209	80	766	796	2330	http	GET	http://image	VCache	30

Table 3. Normalized training patterns for BPN

Normalized IP addresses	Valid(0) / Invalid(1)	Normalized IP addresses	Valid(0) / Invalid(1)
0.549298	0	0.815306	0
0.57671	0	0.819334	0
0.588196	0	0.819424	0
0.753141	0	0.819514	0
0.760483	0	0.819537	0
0.780321	0	0.026241	1
0.780564	0	0.007997	1
0.791906	0	0.027761	1
0.795886	0	0.28298	1
0.803925	0	0.002624	1
0.803937	0	0.0819573	1
0.808023	0	0.000331	1
0.811643	0	0.081742	1
0.81187	0	0.076052	1
0.811933	0		

1.2

1

0.8 0.6

0.4

0.2

n

Predicted Output Values

• After training the BPN, it is tested with test patterns as shown in Table 4.

4. Results

The training of the neural networks has been conducted using the Back Propagation neural network algorithm for 50,000 iterations of the selected training data. After training the BPN, the following results are obtained.

The results obtained match very closely with the desired root mean square (RMS) error as shown in Table 5. Though this method is not designed to be used as a complete intrusion detection system, yet the results show the potential of neural networks to detect individual instances of possible misuse from a representative webbased data. Graphs in Figure 1 show the results for different number of hidden layers used in the BPN. It is evident from the graphs that the results are very close to desired output values, when we use 10 numbers of neurons for hidden layer.

5. Discussions

The above mentioned method can be used to find out the web documents that should not be allowed in the organization. Web Server log file is divided into two parts. One file contains only the destination ip addresses and the second file contains the corresponding source ip and date

Table 4. Normalize testing patterns for BPN

IP Patterns for Testing	Valid(0) / Invalid(1)	IP Patterns for Testing	Valid(0) / Invalid(1)
0.000771	0	0.00082	0
0.000776	0	0.000823	0
0.000788	0	0.000826	0
0.000793	0	0.000831	0
0.000794	0	0.819573	1
0.000795	0	0.000331	1
0.000796	0	0.081742	1
0.000798	0	0.002624	1
0.000799	0	0.076052	1
0.0008	0	0.259212	1
0.000802	0	0.008221	1
0.000813	0	0.027213	1
0.000818	0	0.000819	1

Table 5. RMS error corresponding to hidden layers

No of Hidden Layers	RMS Error (Training Data)
2	0.026315
5	0.024311
10	0.023302



Hidden Layer = 2









Figure 1. Predicted output for test patterns: in (a) 2, in (b) 5, and in (c) 10 hidden layers are used

and time of the site being accessed. Input of the first file having ip addresses of the sites being accessed is converted into normalized ip address. This is the input pattern to Neural Network for testing. For the ip addresses having errors (invalid websites) and no errors (valid websites) the Neural Network is already trained. When a user tries to access a website that is in the invalid website record, it is detected by the system. At the time there is a

Table 6. Web site address to be included in the invalid web site record

Address of the Web Site	ip address
www.bollyexpress.com	208.101.17.60
www.maxalbums.com	64.246.28.216

deviation in the log files under testing it will be figured out. Here in our case Normalized ip pattern 0.002624 is reported as invalid and its corresponding website is www.mp3fine.com. The corresponding source ip address, time, and date can be found from the second file.

We have manually analyzed Web Server log for duration of 15 minutes after the first detection is reported in the system. This is because there is a probability that the user on the system may try to access some similar sites that should be in the invalid web site record, but are not included in the invalid website record previously. This analysis gives us positive results and two sites have been included in the invalid website record as mentioned in Table 6.

There are lots of web documents which provide anonymous downloads of the files of larger size like movie and songs files. If a user is allowed to access these sites, then a large portion of the network bandwidth will be wasted. Many of the sites are already blocked by the Network Administrator, but some sites are still in use. When a user is stopped to access a web document he/she will try to access another web document with similar facility that is missed to block by the Network Administrator. The analysis discussed above can be used to block these types of Web documents.

6. Conclusions

Research and development of intrusion detection system has been ongoing last couple of decades and the challenges faced by designers have increased many fold. Misuse detection is particularly difficult problem because of the extensive number of vulnerabilities in computer systems and the creativity of the attackers. Neural networks provide a number of advantages in the direction of these attacks. The results of our tests for the Proxy Server (Microsoft ISA Server 2000) log show that this technique can be applied for detecting worthless web document access to save the network bandwidth.

REFERENCES

- R. J. Anderson, "Why cryptosystems fail," In Communications of the ACM, Vol. 37, No. 11, pp. 32–40, November 1994.
- [2] http://www.cert.org/reports/dsit_workshop-final.html.
- [3] H. Varian, "Managing online security risks," Economic Science Column, The New York Times, June 2000.

- [4] SANS Institute staff, "Intrusion detection and vulnerability testing tools: what works?" 101 Security Solutions E-Alert Newsletters, 2001.
- [5] T. K. Kim, D. Y. Lee, and T. M. Chung, "Mobile agentbased misuse intrusion detection rule propagation model for distributed system," Lecture Note in Computer Science, Vol. 2510, pp. 842–849, 2002.
- [6] O. Depren, M. Topallar, E. Anarim, and M. K. Ciliz, "An intelligent intrusion detection system (IDS) for anomaly and misuse detection in computer networks," Expert Systems with Applications, Vol. 29, No. 4, pp. 713–722, November 2005.
- [7] T. Konno and M. Tateoka, "Accuracy improvement of anomaly-based intrusion detection system using taguchi method," Proceeding of Symposium on Applications and the Internet Workshops (SAINT-W'05), 0-7695-2263-7/05, 2005.
- [8] K. Ilgun, "USTAT: A real-time intrusion detection system for UNIX," Proceeding of the 1993 Computer Society Symposium on Research in Security and Privacy, Oakland, California, Los Alamitos, pp. 16–28, May 1993.
- [9] K. Fox, R. Henning, J. Reed, and R. Simonian, "A neural network approach towards intrusion detection," Proceeding of 13th National Computer Security Conference, Baltimore, MD, pp. 125–134, 1990.
- [10] J. Frank, "Artificial intelligence and intrusion detection: current and future directions," Computers and Security, Vol. 14, No. 1, pp. 31–31(1), 1995.
- [11] L. Fu, "A neural network model for learning rule-based systems," Proceeding of the International Joint Conference on Neural Networks, pp. 343–348, 1992.
- [12] D. Hammerstrom, "Neural networks at work," IEEE Spectrum, pp. 26–53, June 1993.
- [13] J. Zimmermann, L. Mé, and C. Bidan, "An improved reference flow control model for policy-based intrusion detection," Proceeding of the 8th European Symposium on Research in Computer Security (ESORICS), pp. 291– 308, October 2003.
- [14] G. J. Nalepa, "Application of the XTT rule-based model for formal design and verification of internet security systems," Lecture Notes in Computer Science, Vol. 4680, pp. 81–86, 2007.
- [15] D. Dorothy, "An intrusion-detection model," IEEE Transactions on Software Engineering, Vol. 13, No. 2, pp. 222– 232, February 1987.
- [16] M. M. Sebring, E. Shellhouse, M. E. Hanna, and R. A. Whitehurst, "Expert systems in intrusion detection: a case study," Proceeding of the 11th National Computer Security Conference, Baltimore, MD, pp. 74–81, October 1988.
- [17] S. Staniford-Chen, S. Cheung, R. Crawford, M. Dilger, J. Frank, J. Hoagland, K. Levitt, C. Wee, R. Yip, and D. Zerkle. "GrIDS, a graph based intrusion detection system for large networks," Proceeding of the 20th National Information Systems Security Conference, Vol. 1, pp. 361– 370, October 1996.
- [18] P. A. Porras and P. G. Neumann, "Emerald: event moni-

toring enabling responses to anomalous live disturbances," Proceeding of the 20th National Information systems Security Conference, pp. 35–365, October 1997.

- [19] S. Freeman, "Host based intrusion detection using user signatures," Computer Science Master's project, May 2002.
- [20] A. K. Ghosh, A. Schwartzbard, and M. Schatz, "Learning program behavior profiles for intrusion detection," Proceeding of the 1st Workshop on Intrusion Detection and Network Monitoring, pp. 51–62, April 1999.
- [21] A. "Oks"uz, "Unsupervised intrusion detection system," Master Thesis, Technical University of Denmark, 2007.
- [22] A. Boukerche, K. R. Lemos Juc, J. B. Sobral, and M. Sechi Moretti Annoni Notare, "An artificial immune based intrusion detection model for computer and telecommunication systems," Parallel Computing, Vol. 30, No. 5–6, pp. 629–646, 2004.
- [23] R. Beghdad, "Modelling and solving the intrusion detection problem in computer networks," Computers and Security, Vol. 23, No. 8, pp. 687–696, 2004.
- [24] T. F. Lunt and R. Jagannathan, "A prototype real-time intrusion-detection system," Proceeding of the Symposium on Security and Privacy, New York, pp. 59–66, April 1988.
- [25] T. D. Garvey and T. F. Lunt, "Model based intrusion detection," Proceeding of the 14th National Computer Security Conference, pp. 372–385, October 1991.

- [26] K. Ilgun, "Ustat: A real-time intrusion detection system for UNIX," Master's thesis, Computer Science Dept, UCSB, July 1992.
- [27] S. Kumar and E. H. Spafford, "A pattern matching model for misuse intrusion detection," The COAST Project, Purdue University, 1996.
- [28] J. Ryan, M. Lin, and R. Miikkulainen, "Intrusion Detection with Neural Networks," AI Approaches to Fraud Detection and Risk Management: Papers from the 1997 AAAI Workshop (Providence, Rhode Island), pp. 72–79, 1997.
- [29] H. Debar and B. Dorizzi, "An application of a recurrent network to an intrusion detection system," Proceeding of the International Joint Conference on Neural Networks, pp. 478–483, 1992.
- [30] A. Abraham, C. Grosan, and C. Martin-Vide, "Evolutionary design of intrusion detection programs," International Journal of Network Security, Vol. 4, No. 3, pp. 328–339, March 2007.
- [31] M. Denault, D. Gritzalis, D. Karagiannis, and P. Spirakis, "Intrusion detection: approach and performance issues of the securenet system," Computers and Security, Vol. 13, No. 6, pp. 495–500, 1994.
- [32] S. E. Smaha, "Haystack: an intrusion detection system," Proceeding of the Fourth AeroSpace Computer Security Applications Conference, Orlando, FL, pp. 37–44, December 1988.

On Solvable Potentials, Supersymmetry, and the One-Dimensional Hydrogen Atom

R. P. Martínez-y-Romero¹, H. N. Núñez-Yépez², A. L. Salas-Brito^{3*}

¹Facultad de Ciencias, Universidad Nacional Autónoma de México, Apartado Postal, Coyoacán, México ²Departamento Física, Universidad Autónoma Metropolitana-Iztapalapa, Apartado Postal, Iztapalapa, México ³Laboratorio de Sistemas Dinámicos, Departamento de Ciencias Básicas, Universidad Autónoma Metropolitana-Azcapotzalco, Apartado Postal, Coyoacán, México E-mail: asb@correo.azc.uam.mx, nyhn@xanum.uam.mx

Received November 20, 2009; accepted December 15, 2009

Abstract: The ways for improving on techniques for finding new solvable potentials based on supersymmetry and shape invariance has been discussed by Morales *et al.* [1] In doing so they address the peculiar system known as the one-dimensional hydrogen atom. In this paper we show that their remarks on such problem are mistaken. We do this by explicitly constructing both the one-dimensional Coulomb potential and the superpotential associated with the problem, objects whose existence are denied in the mentioned paper.

Keywords: one-dimensional hydrogen atom, one-dimensional Coulomb potential, supersymmetric quantum mechanics.

A paper of Morales *et. al.* [1] has discussed the use of supersymmetric and shape invariance techniques and of Darboux and intertwining transformations, for building new solvable potentials.

To illustrate these ideas they apply them to hydrogen-like potentials and to radial and one-dimensional problems. They assert [page 23 of [1], in the paragraph after Equation (39)] that the potential corresponding to a one-dimensional hydrogen atom, *i. e.* a one-dimensional Coulomb potential, is nonexistent. They further claim that there is no superpotential associated with the -1/|x|potential energy term [page 22 of [1], in the paragraph just before their Equation (36)]. In this letter we want to challenge these two affirmations. Throughout this work we use atomic units $q_e = -m = 1$. In this paper we want to discuss their results concerning such 1D problem.

We recognize from the start that the potential deserving the name one-dimensional Coulomb potential is not the one usually alluded to in the literature—*i. e.* it is not -1/|x|. The true Coulomb potential in one dimension must be the solution of the corresponding Poisson equation

$$\nabla^2 \phi 1 D C = -4\pi \delta(x) \tag{1}$$

where $\delta(x)$ is a Dirac delta function which is really not a function but a distribution also termed a generalized function [2]. As it is very easy to realize, just solving Equation (3), the 1D Coulomb potential definitively exist and is given by

$$\phi 1DC = -2\pi \left| x \right| \tag{2}$$

so the potential energy function needed in the Schrödinger equation should be

$$V_{1DC}(x) = 2\pi |x| \tag{3}$$

In this sense the one-dimensional Coulomb potential does indeed exist. However, $V_{1DC}(x)$ is not the potential energy usually referred to as the one-dimensional hydrogen atom potential. But even if Morales *et. al.* are referring to this potential, namely $V_{1DH} = -1/|x|$, corresponding to a Hamiltonian

$$H_{1DC} = -\frac{1}{2} \frac{d^2}{dx^2} - \frac{1}{|x|}$$
(4)

the existence of a superpotential is beyond doubt, as we intend to exhibit in this work, see also [3,4]. The result [Equation (36) in [1]] they base their argument on the nonexistence of a superpotential for the one-dimensional hydrogen atom cannot be right since it does not have any explicit *r*-dependence. Even though this problem is surely just a misprint, the limit $l \rightarrow 0$ has no meaning for discrediting Hamiltonian (4) because the problem really comes from the need to describe Coulomb systems constrained to one-dimensional motions with no spherical symmetry and hence described by states with no well defined angular momentum.



Any system described by Hamiltonian (4) is one with baffling properties [5–7]. Its properties are so peculiar that people is prone to express erroneous concepts about it. For example, it has been claimed that the potntial energy term in Hamiltonian (4) is its own supersymmetric partner [8], or, as in [1], that the Hamiltonian itself cannot really be written since its potential energy function does not exist. On the other hand, it has been proven that it violates which the nondegeneracy theorem for onedimensional quantum problems [5], and it has been shown that a superselection rule, analogous to the one preventing the so-called paradox of optical isomers of quantum chemistry, operates in the system [6,9–13]; see also [14,15] for other similar points of view. The Hamiltonian (4) is not in general self-adjoint (in conventional physics parlance, is not Hermitian). Self-adjoint 4 parameter extensions have been derived in [16], such extension admits Hamiltonian (4) as one of its members [7,16–19]. Let us emphasize that Hamiltonian $H_{\rm D}$ together with the matching condition $\varphi(x)|_{x=0}=0$ is selfadjoint.

We think the misconception in the Morales et. al. paper could have arisen from their ideas on how the onedimensional hydrogen atom problem come to be. As they say that, according to certain authors [1], its equation arises from the radial Schrödinger equation of the (3D) hydrogen atom merely by substituting r by x and a vanishing angular momentum l = 0. Given such assertion, we assume that they think the 1D hydrogen atom is a purely formal problem with little or no relation to any actual systems. This, however, is not so. There are specific problems which lead to essentially one-dimensional quantum motions which may be described by Hamiltonian (4). Examples of such problems are an hydrogen atom placed in a constant but super-strong magnetic field **B** [20–22], or the problem of the motion of an electron sitting on a surface producing an image charge as happens to electrons over a pool of liquid helium. In this last case, given the charge and its image is hence clear that the electron is acted by a Coulomb interaction [23]. In the case of the hydrogen atom within a **B** field, any electron state may be expressed as a product of transverse Landau states times a state depending on a coordinate parallel to \mathbf{B} — states with no spherical symmetry [21]. The motion tranverse to the magnetic field is classically restrained to distances of the order of $\rho_c = (c/\mathbf{B})^{1/2}$. In the quantum case ρ_c may be called the mean size of the Landau states. So, as the intensity of the magnetic field is increased, $\rho_c \rightarrow 0$ leaving only the motion along **B** for a dynamical description [20]. When the (x-pointing) magnetic field is super-strong the potential felt by the electron can be approximated as

$$V(r) \approx \lim_{B \to \infty} -\frac{1}{\sqrt{\rho_c^2 + x^2}} - \frac{1}{|x|}$$
(5)

This is the potential used in Equation (4). Hence the name one-dimensional hydrogen atom is justified: it is just an hydrogen atom constrained to move in one direction and under the assumption that any transverse motions can be disregarded for field strenghts $B \sim 10^9$ Gauss typical of neutron stars [24] they are certainly very small. It is worth noting that an hydrogen atom in a magnetic field has two integrable cases: 1) when B=0, and, 2) when B= ∞ .

As we have shown previously [3–6], the two eigenfunctions describing the ground state of the one dimensional hydrogen atom are

$$\psi_0^+(x) = \begin{cases} 2xL_0^1(2x)\exp(-x) & \text{if } x \ge 0\\ 0 & \text{if } x < 0 \end{cases}$$
(6)

and

$$\psi_0^{-}(x) = \begin{cases} 0 & \text{if } x > 0\\ 2xL_0^1(-2x)\exp(x) & \text{if } x \le 0 \end{cases}$$
(7)

where the $L_0^1(x)$ are generalized Laguerre polynomials [17]. Notice the vanishing of the eigenfunctions at x = 0 and the explicit separation between the x > 0 and the x < 0 regions. This is one of the manifestations of the superselection rule which, among other things, prohibits any superposition of the right ψ_0^+ with the left ψ_0^- eigenstates. The energy eigenstates of the problem are given by a Balmer-like formula [4,13,25,26] $E_n = -1/2 n^2$, $n=1,2,3, \ldots$, so the ground state energy is $E_1 = -1/2$.

With the ground eigenstates given above, the superpotential can be easily calculated as [3,27,28]

$$W(x) = -\frac{\psi_0(x)}{\psi_0(x)} = \operatorname{sgn}(x) - \frac{1}{x}$$
(8)

where sgn(x) is the signum function and we have included in a single formula the consequences of both the right and the left eigenfunctions. Using the superpotential, the corresponding partner potentials are readily evaluated

$$V + (x) = \frac{1}{2} - \frac{1}{|x|} + \frac{1}{x^2}, \text{ and } V - (x) = -\frac{1}{|x|} + \frac{1}{2}$$
(9)

where, clearly, $V_{.}$ is the one-dimensional hydrogen atom potential, V_{1DH} , but shifted so that its ground state energy is zero, and V_{+} is the partner potential. Also, the raising and lowering operators are

$$A^{+} = -\frac{d}{dx} + W \tag{10}$$

and

$$A = \frac{d}{dx} - W \tag{11}$$

where, as it is easy to show,

$$[A, A^+] = 2\frac{dW}{dx} \tag{12}$$

and

$$V_{+} + V_{-} = 2W^2 \tag{13}$$

The results (8) to (13) establish that the one dimensional potential V_{1DH} can be regarded as stemming from the superpotential W(x) in Equation (8). In [4,29], we have discussed a complete supersymmetric extension of the one-dimensional hydrogen atom problem, with Hamiltonian

$$H_{susy} = -\frac{1}{2}\frac{\partial^2}{\partial x^2} + \frac{1}{2x^2} - \frac{1}{|x|} + \frac{1}{2} + \frac{1}{2x^2}\sigma_z \qquad (14)$$

where σ_z is a standard Pauli matrix which is needed to operate on both the the fermionic and bosonic sectors of the system. But, as the motivations of [29] were the similarities between light-cone singularities in quantum field theory with the singularity in (4), the results in [29] are not all related to the present discussion. Second, that Morales *et al.* have mistaken the paper they cite (reference [21] in their paper, reference [30] in this work) for other of our papers dealing with the one-dimensional hydrogen atom, since [30] has nothing to do with the problem at hand. It deals with a solvable model in relativistic quantum mechanics, the Dirac oscillator, which at the time was thought to have applications in QCD. They should have cited [3,5,29] instead.

Acknowledgements

We acknowledge with thanks the comments and suggestions of P. M. Schwartz, P. N. Zeus, P. M. Mec, and G. R. Maya.

REFERENCES

- J. Morales, J. J. Peña, J. L. López-Bonilla, and J. Mol. Struct. (Teochem), Vol. 621, pp. 19, 2003.
- [2] M. J. Lighthill, "Fourier analysis and generalized functions, cambridge university press," Cambridge, 1975.
- [3] R. P. Martínez-Romero, C. A. Vargas, A. L. Salas-Brito, and H. N. Núñez-Yépez, Rev. Mex. Fis., Vol. 35, pp. 617, 1989.
- [4] B. Jaramillo, R. P. Martínez-y-Romero, H. N. Núñez-Yépez, and A. L. Salas-Brito, Phys. Lett. A, Vol. 374, pp. 150, 2009.
- [5] H. N. Núñez-Yépez, C. A. Vargas, A. L. Salas-Brito, Eur. J. Phys., Vol. 8, pp. 189, 1987.
- [6] H. N. Núñez-Yépez, C. A. Vargas, A. L. Salas-Brito, J. Phys. A, Math. Gen., Vol. 21, pp. L651, 1988.

- [7] I. Tsutsui, T. Fülöp, and T. Cheon, J. Phys. A: Math. Gen. Vol. 36, pp. 275, 2003.
- [8] T. D. Imbo and U. P. Sukhatme, Phys. Rev. Lett., Vol. 54, pp. 2184, 1985.
- [9] H. N. Núñez-Yépez, C. A. Vargas, A. L. Salas-Brito, Phys. Rev. A, Vol. 39, pp. 4307, 1989.
- [10] P. Pfeifer, "Dissertation, Eidgenössiche Technisch Hochschule," Zürich, 1980.
- [11] P. Pfeifer, in J. Hinze ed. Energy Storage and Redistribution in Molecules, Plenum Press, New York, pp. 315, 1983.
- [12] R. P. Martínez-Romero, H. N. Núñez-Yépez, A. L. Salas-Brito, "A simple introduction to superselection rules in nonrelativistic quantum mechanics can be found in C. Cisneros," Eur. J. Phys. Vol. 19, pp. 237, 1998.
- [13] L. J. Boya, M. Kmiecik, A. Bohm, Phys. Rev. A., Vol. 37, pp. 3567, 1988.
- [14] R. G. Newton, J. Phys. A: Math. Gen., Vol. 27, pp. 4717, 1994.
- [15] U. Oseguera, M. de Llano, and J. Math. Phys. Vol. 43, pp. 4575, 1993.
- [16] W. Fischer, H. Leschke, P. Muller, and J. Math. Phys., Vol. 36, pp. 2313, 1995.
- [17] M. M. Nieto, Phys. Rev. A, Vol. 61, pp. 034901, 2000.
- [18] C. R. de Oliveira and A. A. Verria, Annals of Physics Vol. 324, pp. 251, 2009.
- [19] S. Nouri, Phys. Rev. A, Vol. 65, pp. 062108, 2002.
- [20] R. P. Martínez-Romero, H. N. Núñez-Yépez, A. L. Salas-Brito, and C. A. Vargas, "Actas de la 3a. Reunión Latinoamericana de Colisiones Atómicas, Moleculares y Electrónicas," CNEA, Bariloche, Argentina, 1989.
- [21] W. Rössner, G. Wunner, H. Herold, and H. Ruder, J. Phys. B: At. Mol. Opt. Phys., Vol. 17, pp. 29, 1984.
- [22] W. Edelstein and H. N. Spector, Phys. Rev. B, Vol. 39, pp. 7697, 1989.
- [23] M. W. Cole and M. H. Cohen, Phys. Rev. Lett. Vol. 23 pp. 1238, 1969. M. W. Cole, Phys. Rev. B, Vol. 2, pp. 4239. 1970.
- [24] H. Ruder, G. Wunner, H. Herold, and F. Geyer, "Atoms in strong magnetic fields," Springer Verlag, Berlin, 1994.
- [25] S. H. Patil, Phys. Rev. A, Vol. 64, pp. 064902, 2001.
- [26] G. Abramovici and Y. Avishai, J. Phys. A: Math. Theory, Vol. 42, pp. 28532, 2009.
- [27] F. Cooper and J. G. Ginocchio, A. Khare, Phys. Rev. D, Vol. 36, pp. 2458, 1987.
- [28] R. Montemayor and L. D. Salem, Phys. Rev. A, Vol. 40, pp. 2170, 1989.
- [29] R. P. Martínez-Romero and H. N. Núñez-Yépez, A. L. Salas-Brito, Phys. Lett. A, Vol. 142, pp. 318, 1989.
- [30] J. Benítez, R. P. Martínez-y-Romero, H. N. Núñez-Yépez, and A. L. Salas-Brito, Phys. Rev. Lett. Vol. 64, pp. 1643. 1990.



How to Measure in the Near Field and in the Far Field

Tomasz Dlugosz, Hubert Trzaska

Wroclaw University of Technology Institute of Telecommunications, Teleinformatics and Acoustics, Wyspianskiego, Wroclaw, Poland Email: Tomasz.Dlugosz@pwr.wroc.pl Received November 5, 2009; accepted December 24, 2009

Abstract: A background of the electromagnetic field (EMF) measurements is presented in the work. A special attention is given to the specificity of the measurements performed in the Near Field. Factors, that should be taken into consideration as during the measurements as well during their analysis, are discussed. Without their understanding and considering a comparison of the measurements' results, meters' calibration and EMF standards comparison between different centers is impossible.

Keywords: electromagnetic fields measurements, the near field, the far field

1. Introduction

Surfing on the World Wide Web, when in one of the most popular browsers we enter the words: "electromagnetic field" (EMF), we obtain over 1.5 million answers. In various libraries we also can find a few hundred thousand documents, publications and books pertaining to EMF measurements. It would seem that one more publication on this subject is superfluous, but experience shows something totally different. In reality, in many cases the manner in which EMF measurements are performed is an affront to any forms of correctness and has nothing to do with accuracy and engineering diligence. Even people familiar with this domain forget about some conditions which have to be met in order to carry out EMF measurement correctly, which means, with the required accuracy [4]. A good example that there is no understanding of the EMF metrology fundamentals is the EMF measurement in a room with the use of a log-periodic antenna, described in [5]. You can wonder what in fact has been measured?

Why is EMF metrology so important? The answer is relatively simple, because it consitutes a sine qua non condition of the activities associated with protection of electromagnetic environment, as well as of fundamental research, especially research on EMF impact on the animate matter, in particular, on human beings. Such research is an initial step leading to determination of protective regulations, pertaining both to the safety of work as well as protection of the general population. As an interesting side note, we shall remind here that in spite of the poor EMF measurement accuracy and even lesser accuracy of biomedical research based on them, the protection standards are determined with an amazing accuracy. And the EMF metrology is not counted among the easiest and the most accurate. If the achievable accuracy in the far field amounts to 1 dB, in the near field it is only 3 dB, and even 6 dB! This fact shows that the existing measurement methods need to be analysed and their accuracy increased and that new measurement techniques should be pursued, e.g. photonic sensors [1].

2. Is It Still the Near Field or Already the Far Field?

Prior to discussing the differences existing in EMF metrology in the near and in the far field, meaning of these notions should be defined. What does "the near field" mean? The authors propose two new definitions. The first one, more general and less rigorous, can be as follows: the near field is the field surrounding primary and secondary radiation sources where measurement accuracy is limited (e.g.) to 5 %, as compared with the far field. The second definition is more demanding: the near field exists everywhere where we carry out measurements. This definition results from the experience and it refers to measurements in urbanized areas where multipath propagation may occur and we have to do with interference and reflections - sometimes reflected rays can be stronger than the direct ray. This shows that it is necessary to act with due caution even during measurements in the far field, where directional antennas are used, which may not "catch" all transmitted rays. And here we encounter a paradox - a correctly calibrated meter does not ensure the expected measurement accuracy.

In the traditional approach (Figure 1), in order to dis-



Figure 1. Near and Far Field around an antenna

tinguish the specificity of measurements in the near field and in the far field, a criterion was adopted which enables delimitation of these two areas, although there is no clearly defined and discrete boundary between the near field and the far field [8].

If D is adopted as the largest antenna size and the emitted wave length is designated, the boundary (R) between the near zone and the far zone can be determined from the following relationship:

$$R \ge \frac{2D^2}{\lambda} \tag{1}$$

In order to demonstrate that the far zone can be the same for different types of antennas, operating on different frequencies, two examples will be given (Trzaska, 2002):

• Example No. 1:

For a parabolic antenna with 3 m dish diameter operating on 10 GHz frequency the far field boundary is at 600 m,

• Example No. 2:

For the antenna of the former transmission centre in Gąbin (Poland) having a height of 0.5 λ and operating on 227 kHz frequency the Far Field boundary is at 660 m.

As you can see, the Far Field zone is not something assigned permanently to a given antenna operating on the preset frequency. As the above two examples show, the same far zone boundary exists for extremely different antennas. Also the Near Field can be a function of electrical size [2].

3. Measured Quantities

In the Near Field, the mutual relationship between electric field (E) and magnetic field (H) components depends on the type of EMF source and on the distance between the source and the observation point. Therefore, determination of one of them is not sufficient for computing the other.

Situation is different in the case of the far field where knowledge of one of the field components, e.g. of electric field vector $- \mathbf{E}$, enables determination of the other (magnetic field vector $- \mathbf{H}$), using the relationship in which these two quantities are interrelated by means of the impedance of free space (Z):

$$\mathbf{n} \times \mathbf{E} = Z\mathbf{H} \tag{2}$$

In both cases, i.e. in the near and in the far field, when we know the E and H components, we are in a position to determine the power density. With this aim, the mean value of the Poynting vector (S) is determined:

$$\mathbf{S} = 0.5 \operatorname{Re}\{\mathbf{E} \times \mathbf{H}\}$$
(3)

In the Far Field metrology it is not necessary to carry out an additional measurement of quantities other than the E, H or S, contrary to the near field metrology in which the temperature increase and current density, caused by the EMF impact, are also measured.

Measurement of the temperature increment (ΔT), resulting from the EMF impact, of a material which has a given specific heat (cw), makes it possible to determine the Specific Absorption Rate SAR:

$$SAR = \frac{c_w \Delta T}{t} \tag{4}$$

The SAR is commonly used for examination of the EMF impact on human body. However, there are some limitations of its use, which are discussed in detail in [8]. In this paper we shall only note that the SAR parameter can be used for the frequencies higher than 300 MHz due to too small sensitivity. In the lower frequency ranges an essential parameter is the density of the current induced into tissues [7]. Knowing the conductivity (σ) of the examined medium and the density value of electric field (E) existing in this medium, the current density (J) can be calculated:

$$J = \sigma E \tag{5}$$

The manner of measurements of the current flowing through a human body is described in [7]. Often measurements of the currents flowing through legs or feet are presented, neglecting the currents appearing in other parts of the body, or unmeasurable eddy currents.

For electric field measurements in the near field electrically-short dipole antennas are used, while magnetic field is measured by means of small frame antennas. In the far field directional antennas are used. An essential problem faced in EMF measurements, regardless of what sensor is used, is the sensor's presence in the measured field, which causes deformation of this field and mutual interaction between the sensor and the neighbouring material objects. This interaction constitutes a serious factor affecting the measurement accuracy, both during EMF measurements and EMF sensor calibration, as well as in cases when we use exposure kits for examination of the features of any material object [6].

In the measurements polarization is also important. After all, the E, H and S vectors can have three spatial components each (quasi-ellipsoidal polarization caused, for example, by rotation of the polarization plane in space). In such a case isotropic sensors have to be used.

4. Measurement Accuracy

Measurement accuracy constitutes the biggest problem in EMF measurements. For calibration of EMF meters, as well as for examination of equipment and matter sensitivity to EMF impact, EMFs of known parameters are used - standard EMFs. For generation of a standard EMF knowledge of not only the values of generated parameters is necessary, but also of the accuracies of their generation. EMF standards are among the least accurate as compared with the standards of other physical quantities. Many of such quantities are determined with an accuracy of 10^{-10} % or higher, while the error of standard EMFs generation in renowned centres ranges from 5 % to 10 %. In other words, even before we commence field measurements, from the very beginning the measurement result is burdened with an error which amounts to 5 % in the best case, and this is not all.

The main factor which limits EMF measurement accuracy in the near field is the antenna dimensions. Point antennas would be the best to use, because otherwise an antenna causes averaging of the measured EMF values. Variations of the spatial field strength, resulting from either amplitude or phase variations, are subject to averaging. These variations depend on the curvature of the EMF field which surrounds the source [3,8]. Some examples of error graphs, both amplitude and phase errors, are shown on Figures 2 and 3 (where: Ro – the distance between the source and the measuring antenna centre, α – an exponent characterizing field curvature, h – the length of dipole arm, k – propagation constant). The presented curves refer to a dipole antenna but identical considerations are applicable to a frame antenna as well [8].

Passing over the impact of the meter used and of the person performing the measurements on the disturbances of the measured EMF, you should not forget the error



Figure 3. Phase error **\delta** f

which is contributed by the measuring person, which we shall call a "human factor". This factor also depends on the conditions in which measurements are performed and its importance is essential, as it is shown in [4]. This factor is described on the basis of two measurement series, performed by four persons in the same measuring points, by means of two meters: MEH-25 with 3AS-1 probe and PMM 8053A with EP-300 probe. This simple experiment has shown (see Table 1) how diversified the measurement results can be if the measurements are performed by different persons. Therefore, the "human factor" is a gross error but, unfortunately, it is not taken into account when measurement results are worked out.

Desition of			S	eries I					Se	eries II		
	1	2	2	4	Mean	$\delta_{(min-max)}$	1	2	2	4	Mean	$\delta_{(min-max)}$
measurements	I	2	3	4	value	[%]	I	2	3		value	[%]
1	14.9	16.5	17.6	15.5	16.1	8.3	16.6	15.8	14.2	16.9	15.9	8.7
2	17.6	16.2	18.5	19.0	17.8	8.0	16.5	16.5	19.2	18.5	17.7	7.6
3	9.2	7.3	6.6	8.8	8.0	16.5	5.8	8.8	8.2	6.1	7.2	20.5
4	9.1	8.2	10.2	8.2	8.9	10.9	9.6	10.4	9.9	8.6	9.6	9.5
5	9.9	10.2	11.6	11.0	10.7	7.9	11.0	10.4	14.0	10.4	11.5	14.8

Table 1. "Human factor" measurement results [4]

Parameter	Near Field	Far Field
measured EMF component	E, H & S	E or H, and S on mwaves
other magnitudes measurement	I, T, (SA, SAR) "HESTIA"	unnecessary
spatial components	3	1 or 2
polarization	quasi-ellipsoidal	linear or elliptical
environment	complex, multipath propagation & inter- ference	usually simple
frequency spectrum	wide, often unknown, many fringes	usually single frequency
antennas	small, omnidirectional	resonant, directional
emporal & spatial EMF alternations	significant	usually negligible
uncertainty	3, 6 or more dB	around 1 dB
temperature sensitivity	significant	unessential
susceptibility	significant	ommitable
influence of surroundings	significant	usually ommitable
procedures	complex	simple
agreement with theory	reasonable	good
measured levels	V/m, kV/m	mV/m, mV/m

Table 2. Comparison of measurements in the near field and in the far field

5. Summary

The paper presents a comparative analysis of EMF metrology in the near field and in the far field. Measurements in the near field are more difficult and burdened with a considerably larger error than measurements performed in the far field. As you can see there are many factors which have an impact on measurement accuracy and the selection of a measurement zone should involve proper selection of adequate tools and measurement techniques.

It is not feasible to present all aspects of EMF measurements in the near field and in the far field. Due to practical limitations of this paper only most important aspects of this metrology are discussed herein, supplemented by Table 2.

REFERENCES

- K. Abramski and H. Trzaska, "FM EMF Sensors," 3rd International Symposium on Electromagnetic Compatibility, Beijing, China, pp. 222–225, May 2002.
- [2] T. Adams, Y. Lewiatan, and K. S. Nordby, "Electromagnetic near fields as a function of electrical size," IEEE Transaction on Electromagnetic Compatibility, Vol. 25, No. 4, pp. 428–432, 1983.
- [3] P. Bieńkowski and H. Trzaska, "Electromagnetic fields

metrology," National Symposium of Telecommunications, Bydgoszcz, Poland, pp. 167–178, September 1995.

- [4] P. Bieńkowski, H. Trzaska, "Interlaboratory comparisons in EMF survey measurement – methods and results," International Conference and COST 281 Workshop on Emerging EMF-Technologies, Potential Sensitive Groups and Health, Graz, Austria, 20–21 April 2006, (CD proceedings).
- [5] C. Bornkessel and M. Wuschek, "Exposure measurements in different WLAN – scenarios," International Conference and COST 281 Workshop on Emerging EMF-Technologies, Potential Sensitive Groups and Health, Graz, Austria, 20–21 April 2006, (CD proceedings).
- [6] T. Dlugosz and H. Trzaska, "Proximity effects in EMF measurements and standards," XXVIIIth General Assembly of International Union of Radio Science (URSI), 23–29 October 2005, New Delhi, India, (CD proceedings).
- [7] O. P. Gandhi, Y. Chen, and A. Riazi, "Currents induced in a human being for plane-wave exposure conditions 0-50 mhz and for RF sealers," IEEE Transactions on Biomedical Engineering, Vol. BME-33, No. 8, pp.757–767, 1986.
- [8] H. Trzaska, "EMF measurements in the near field," Noble Publ. Co., 2002.
- [9] H. Trzaska, "Limitation in the SAR use," The Environmentalist, Vol. 25, No. 2/4, pp. 81–185, 2005.


Proposed Model for SIP Security Enhancement

Munir B. Sayyad¹, Abhik Chatterjee², S. L. Nalbalwar³

¹Technology Innovation Center Reliance Communication, Maharashtra, India ²Electronics Engineering of Lokmanya Tilak College of Engineering, Mumbai University Maharashtra, India ³Department of Electronics and Telecommunication, Dr. Babasaheb Ambedkar Technological University, Maharashtra, India E-mail: powerabhik@yahoo.com.in, nalbalwar sanjayan@yahoo.com Received October 24, 2009; accepted November 12, 2009

Abstract: This paper aims to examine the various methods of protecting and securing a SIP architecture and also propose a new model to enhance SIP security in certain selected, specific and confidential environments as this proposed method cannot be generalized. Several security measures and techniques have already been experimented with, proposed and implemented by several authors as SIP security is an issue of utmost importance in today's world. This paper however, aims to summarize some of the better known techniques and propose a unique method of its own. It also aims to mathematically represent SIP fitness values graphically as well via a simulation using the popular Fuzz Data Generation Algorithm. Thus this paper not only aims to contribute to the already vast field of SIP security in an effective manner but also aims to acknowledge and represent some of the fail proof methods and encryption techniques that have helped in making SIP a more secure and less wobbly network for all of us to function in.

Keywords: SIP, SoS, VoIP

1. Introduction

Session Initiation Protocol (SIP) is the Internet Engineering Task Force (IETF) standard for IP Telephony which is making huge inroads into the Voice-Over-IP (VoIP) market, previously domineered by implementations which stuck to the rather difficult H.323 ITU-T Internet Telephony standard [4]. The apparent reality is that Voice and Data services are being quickly shifted from the legacy network to the IPbased network.

The standardization of SIP helped to realize the call control function. SIP is the present as well the future of commercial communication systems. SIP is the present as well the future of commercial communication systems.

Many carriers and providers are extensively adopting it; therefore SIP security has become a topic of high importance and priority [5]. With VoIP, voice can now be transported on a traditional IP data network, making use of the vast resources of the Internet and thus drastically lowering the cost of operation.

However in the recent past, VoIP services have been plagued and hampered by numerous security threats and issues. With Internet being the primary carrier, VoIP networks are exposed to threats and dangers that an IP data network faces e.g., IP spoofing, denial of service (DoS) etc. [5].

SIP has become the effective standard for VoIP services. It is described as "an application layer control protocol

that can establish, modify and terminate multimedia sessions (conferences) such as Internet telephony calls". It is an ASCII/text based request-response based protocol that works on a client server mode.

2. Security in Sip

SIP security is an issue of prime importance. Basically we can broadly classify the attacks on any type of system into two categories [2]:

• Passive Attacks: This threatens the confidentiality of the data/signal being transmitted.

• Active Attacks: This threatens the integrity or availability of the data/signal being transmitted.

The feasibility of a passive network primarily depends on the physical transmission media in use and its physical accessibility for any intruder. Fortunately enough, the use of switching technologies makes it harder and more difficult for an attacker to passively attack a signal segment. Now in an active attack, more often than not, the intruder manipulates the domain name system (DNS) to place himself between the sender and recipient of a message. In this situation, the intruder acts as a man-in-the-middle. A very common form of attack is to spoof signals/messages on another (or nonexistent) user's behalf.

These two types of attacks can most probably encompass all the different types of attacks and forced attempts within their broadly diversified branches. The following diagram will give a clearer picture of a SIP security



Figure 1. Protocol architecture



Figure 2. SIP security breakup

breakup.

Authentication and maintaining the integrity of data/ signaling is a matter of the highest priority. It is also important to monitor the access control and the availability of information because it will prevent malformation and spoofing of data.

We will now define a structure which will include all the important points mentioned above which are of primary importance. Authentication, integrity, confidentiality, non-repudiation, access control and availability form a framework upon which the others will be derived.

Authentication is the property by which the correct identity of an entity, such as a user or a terminal, or the originality of a message that has been transmitted, is established with a required assurance.

Authentication can basically be divided into two classes, which are peer entity authentication and data

origin authentication. Peer entity authentication assures that the communicating parties are who they claim to be. Data origin authentication assures that a message has come from a legitimate and authenticated source. Authentication is typically needed to provide safety against masquerading as well as modification.

Integrity means the avoidance of unauthorized modification of information. Integrity is an important security service that proves that transmitted data has not been tampered with. Authenticating the communicating parties is not enough if the system cannot guarantee that a message has not been altered during transmission.

Confidentiality is the avoidance of the disclosure of information without the permission of its owner. Secrecy and privacy are terms synonymous to confidentiality. Confidentiality may be ensured with encipherment of the messages.

Non_Repudiation is the property by which one of the entities or parties in a communication cannot deny having participated in the whole or part of the communication. Non-repudiation prevents an entity from denying something that actually happened.

Access Control is the denial of unauthorized use of a resource. Access control is closely related to authentication, which gives the ability to limit and control access to network systems and applications.

Availability means the accessibility of systems and information by authorized users. It is closely related to authentication and access control. An authenticated entity must have access to a system and on the other hand unauthorized entity must not prevent the usability of the system (Denial of service attacks).

3. Some Security Protocols and Applications for Sip

1) **Encryption** is a mechanism to secure information so that only receiver can use it. In encryption, a cleartext message or plaintext is hidden by using cryptographic techniques, the resulting message is known as ciphertext. The receiver recovers the original plaintext by decrypting the ciphertext.

A key is a mathematical value that modern cryptographic algorithms make use of when encrypting or decrypting a message. Cryptographic techniques are not only used to provide confidentiality, but also other services, like authentication, integrity and non-repudiation may be provided. Cryptographic techniques are typically divided into two generic types: symmetric key and asymmetric key techniques.

a) **Symmetric Encryption** means that the key can be calculated from the decryption key and vice versa. In most cases both keys are the same one and the mechanism is called secret key or single key encryption. The security in symmetric key encryption rests in the key, which must be agreed before any communication. As long as the com-

munication needs to remain secret, the key must be secret, divulging the key means that anyone could encrypt and decrypt the messages.

The Data Encryption Standard (DES) is currently the most widely used symmetric encryption scheme. DES is a symmetric block cipher that processes 64-bit blocks of plaintext producing 64-bit blocks of cipher text The key length is 64 bits, but since every eighth bit $(8, 16, \ldots, 64)$ is a parity bit for error detection, the effective key length is 56 bits.

b) **Asymmetric Encryption** also called public-key encryption, the key used for encryption is different from the key used for decryption and the decryption key cannot be calculated from the encryption key. The encryption key may be published, so that anyone could use the encryption key to encrypt the message, but only the receiver with the corresponding decryption key can decrypt the message. So the encryption key is also called the public key and the decryption key is called private key.

The RSA algorithm is perhaps the most popular public-key algorithm. It was invented by Ron Rivest, Adi Shamir and Leonard Adleman in 1977. RSA can be used for encryption / decryption, providing digital signatures and key exchange. decrypt the message.

The Diffie-Hellman algorithm was the first ever public-key algorithm, invented in 1976 by Whitfield Diffie and Martin Hellman. The algorithm can be used for key exchange but not for encryption/decryption, thus the algorithm is typically used for exchanging the secret keys.

2) **Message-Digest Algorithms** are compact "distillate" or "fingerprints" of your message or file checksum. A message-digest algorithm takes a variable length message as input and produces a fixed length digest as output. This fixed length output is called the message digest, a digest or a hash of the message. The digest, which is typically shorter than the original message, acts as a fingerprint of the inputted message. The message digest verifies your message and makes it possible to detect any changes made to the message by a forger.

4. Novel Proposed Method to Enhance Sip Security in Specific Confidential Sectors: TOUCH ME NOT

In some secure and confidential sectors such as the army (for e.g.) data and signaling leakage is highly volatile and potentially very dangerous. In such cases signal tapping is neither lawful nor desirable. Thus a new security architecture termed TOUCH ME NOT is being proposed inorder to avoid signal tapping. This proposed model is currently under test and development. Its source code has been written in Turbo C++. The testing activity has been carried out using freely available evaluation copies of several popular SIP soft phone clients. Since our testing activity is not complete, we have not informed the vendors about our produced results. Hence, in this paper we



Figure 4. N vs deviation factor

are refraining from using client names.

In this process, there will be present a main security master which will be consisting of a continuous key jumbler whose task will be to randomly jumble and reassign key values in order to prevent key cracking by an intruder.

The security master will also be consisting of a list of predefined attack cycles and algorithms so that it can detect and recognize the most common and difficult types of attacks if any. The entire signaling route from the sender to the sender to the receiver will be divided into several checkpoints. If the attacker attempts to access or tap the signal at any point, on or between the checkpoints, a pre programmed delay generator which may be an exe file will appear as a non removable pop up, displaying random gibberish values or a blank screen.

This will act a cover for the signal to self destruct, in other words the signal will be auto terminated at that point and a signal informing the sender and receiver of the interception or attempted attack on the sent signal will reach the sender as well as the receiver in due time. This will prevent the signal from being tapped with, examined or malformed. This method cannot however be generalized in all sectors as tapping is lawful in several government as well as private sectors.

5. Fuzz Data Generation

We are already aware of the Fuzz Data Generation Algorithm. Fuzz testing or fuzzing is a software testing technique used to find implementation defects using malformed or semi malformed input data [1]. We have to define a set of parameters that contribute to the overall fitness value of a given data. All these parameters need not always be used: a subset of them can be used depending on the input population and the application being fuzzed. They can be for e.g. Native size, Native type, Parent's Fitness etc. [1] The challenge will be to define more and more criterion to define a fitness value. The more the value of N, the better the fitness value. This can be verified from the graph given below as well as the set of relations provided [1].

Let N be the number of parameters chosen to contribute

to the fitness value.

Calculate the deviation factor DF+1/N (We can also calculate a weighted DF, if some of the parameters need to be given more weight compared to the others).

Calculate the deviation contribution DC=A*DF, for each parameter, where A is the deviation percentage.

Calculate total deviation contribution TDC=SUM(DC) for all N.

Final Fitness Value F= Ceiling [TDC*10]

6. Conclusions

Thus we have analyzed some of the methods which make SIP a more secure network. The proposed TOUCH ME NOT architecture is an effective way to prevent illegal tapping in selected confidential setups. Fuzzing data generation along with the simulation can be used to determine fitness values. These steps will hopefully help in making SIP a stronger and a more secure network.

REFERENCES

- [1] IEEE Paper: A SIP Security Testing Framework: Hemanth Srinivasan and Kamil Sarac.
- [2] Applied Cryptography-Second Edition-Protocols, algorithms and Source code in C: Bruce Scheneier.
- [3] SIP Tutorial: Daniel-Constantin Mierla.
- [4] IEEE Paper: SIP Security Issues: The SIP Authentication Procedure and its Processing Load: Stefano Salsano, Luca Veltri, Donald Papalilo.
- [5] IEEE Paper: Security Challenges for Peer-to-Peer SIP: Jan Seedorf.

A Model for Cu-Se Resonant Tunneling Diodes Fabricated by Negative Template Assisted Electrodeposition Technique

Meeru Chaudhri¹, A. Vohra¹, S. K. Chakarvarti²

¹Department of Electronic Science, Kurukshetra University, Kurukshetra, India ²Department of Applied Physics, National Institute of Technology (Deemed University), Kurukshetra, India E-mail: meerachaudhri@rediffmail.com Received November 16, 2009; accepted December 29, 2009

Abstract: In this paper, the authors present and discuss a model for Cu-Se nano resonant tunneling diodes (RTDs) fabricated by negative template assisted electrodeposition technique and formulate the mathematical equations for it. The model successfully explains the experimental findings.

Keywords: track-etch membrane, template synthesis, Cu-Se resonant tunneling diodes, electrodeposition

1. Introduction

For nanoelectronics to become a reality one must be able to fabricate the devices and circuits at nanometer dimensions. For this, the researchers the world over have put in efforts in three different areas: nanofabrication, quantum modeling and circuit innovations. Modeling of a device is an essential part of this effort that provides a test bench and also forms the basis for simulation tools for the device. With the help of models, one can also adjust the structural parameters and keep at bay the undesirable parameters through device design and optimization while fabrication. However, the traditional device modeling is not valid in the nanometer regime [1]. Each of these areas has their own importance. As the nano dimensioned materials lead to new phenomenon and also possibly novel devices based on quantum tunneling mechanisms [2] a device theory that can properly treat quantum transport phenomenon is, therefore called for. In our previous publications [3–6], we have discussed the fabrication and the characterization of RTDs of various diameters made by utilizing different material systems. In this paper, we have developed a model for these RTDs. Equations have been formulated for this model and the experimental results have been verified with the help of these equations.

2. Experimental

Cu-Se RTDs have been fabricated by electrodepositing Cu and Se in the pores of the polycarbonate track-etch membranes (PC TEMs) [3,4]. (PC TEMs) with pores of diameters 1 μ m, 100 nm and 40 nm were used for this purpose. The experimental set-up used to fabricate the

Cu-Se RTDs is shown in the Figure 1.

TEM foils with in situ Cu–Se binary structures were used for obtaining I–V characteristics. However for SEM characterization, membranes were dissolved in the solvent dichloromethane (CH₂Cl₂), should be leaving behind the structures. SEM view of Cu-Se RTD of diameter 1 μ m in back-scattering mode is shown in the Figure 2. This mode is used to obtain the contrast image of the object. In the figure, dark part is indicating Se and bright part is indicating Cu.

An ohmic contact was made by applying Ag based paint on the top side of the Se to obtain I-V characteristics. Figure 3 illustrates the schematic cross-section of the samples in the pores of the membranes with the silver paste.

Experimental results of I-V characteristics of Cu-Se binary structures of diameters 2 μ m, 1 μ m, 100 nm and 40 nm are shown in Figures 4 and 5.

It is clear from the Figures 4 and 5 that a prominent feature of negative differential resistance region (NDR) appear as the diameter of the Cu-Se binary structures reduces from 2μ m to 1nm. This NDR increases with further reduction in the diameters of the Cu-Se binary structures. The values of peak to valley current ratios (PVCRs) of Cu-Se RTDs of different diameters are shown in Table 1.

3. A Model for Cu-Se RTDs

The structure consists of three different layers-Cu, Se and Ag. As the quantum size effects in metals are normally seen at 1 nm [7], density of states (DOS) in Cu and Ag are expected to be continuous. It, thus behaves as a metal. Quantum size effects in semiconducting material





Figure 1. Experimental set-up for negative-template assisted electrodeposition of nano-/micro binary structures







Figure 3. Samples with silver paste

become apparent when the size of the semiconducting material is of the order of hundreds of nanometers [8]. Thus, Se material has quantized bands as shown in Figure 6 (a), with infinite potential on the both sides of it i.e. Se semiconductor at small dimensions, forms a quantum well similar to the one fabricated by exploiting the energy band discontinuities of semiconductor heterostructures. The fabricated Cu-Se-Ag structure with wire shape,



Figure 4. Experimental I-V characteristics of Cu-Se binary structures of 2 μm and 1 μm diameters



Figure 5. Experimental I-V characteristics of Cu-Se RTDs of 100 nm and 40 nm diameters

Table 1. Variation of PVCR with diameters of Cu-Se devices



Figure 6. Model utilized for explaining the I-V characteristics of Cu-Se RTDs (a) equillibrium state (b) electrons flow from Cu to Se when a suitable voltage is applied across this system

hence, forms one dimensional RTD with Cu as emitter, Ag as collector and Se as a potential well. On applying a voltage across the device, the band diagrams can be redrawn as shown in Figure 6(b). The electrons from the Cu electrode tunnel to the empty states in the conduction band of Se. The electrons in the well stay at a particular energy level until these electrons get enough energy to jump to the next higher energy level. These electron waves reflect back and forth between the two walls of the well and interfere, causing the change in the amplitude of the wave. When the energy of the electrons is equal to the energy of the quantized level in the well, the two waves interfere constructively and resonance of the electron wave takes place, which results in maximum transmission of electrons. The accumulation of electrons in the well thus results in a decrease in the current up to valley point current of the I-V curve.

Based on this model of quantization of energy levels, the I-V behavior of the Cu-Se structures has been explained. The energy levels in Cu-Se-Ag structures can be drawn as in Figure 7. The bulk behavior of the Cu-Se binary structures of 2 μ m can be explained by the energy level diagram of Figure 7 (a) where the energy levels are continuous.

However, as the dimensions of the device are reduced, quantized energy levels appear in Se semiconductor and a negative differential resistance region starts appearing. This is shown in Figure 7 (b). On reducing the diameter further, the negative differential resistance region increases and this is illustrated in Figure 7 (c). The energy band diagrams in Figures. 7 (b) and (c) show the increase in spacing in energy levels in the conduction band of the Se with decrease in the dimensions (diameter) of the fabricated binary structures.

Further, the cut-in voltages of the devices increase with decrease in diameters of the device. This indicates an increase in the Schottky barrier height due to increase in band gap of Se as the device dimensions are reduced. Various workers [9,10] have reported an increase in band gap with reduction in dimensions. A similar behavior is expected for Se as well and has been shown in Figures 7(a), (b) and (c).

4. Theoretical Analysis of Experimental Results

In this section, the authors intend to correlate some of the experimental observations. Figures 4 and 5 and Table 1 indicate that there is 1) an increase in cut-in voltage as the diameters of the device is decreased 2) The PVCR increases with decrease in diameters of the device.

4.1 Increase in Cut-In Voltage

The increase in cut in voltage as seen in Figs. 4 and 5 can be explained due to increase in band gap. Such an increase in band gap with decrease in diameter is reported in literature [10–12].

As a metal is brought in contact with a semiconductor, a barrier will be formed at the metal-semiconductor interface. The height of the barrier is governed by metal work function and the electron affinity of the semiconductor. The voltage required to increase the energy of electrons on the metal side to overcome the barrier is cut-in voltage. The cut-in voltage and the band gap of the semiconductor are related as [13].

$$q\phi_b = E_g - q (\phi_m - \chi)$$
(1)

where

 E_g is band gap of the semiconductor

 $q\phi_m$ is work-function of the metal $q\phi_b$ is Schottky barrier height at the metal- semiconductor

contact

 $q\chi$ is electron affinity of the semiconductor



Figure 7. Energy band diagrams and corresponding I-V characteristics of Cu-Se resonant tunneling diodes of different diameters illustrating the emergence of quantum size effects (a) bulk effect (b & c) quantum size effects

(2)

From the Equation 1, it is clear that the cut-in voltage is directly dependent upon the band gap of the semiconductor material i.e. higher the band gap, higher will be the cut-in voltage. Klimov while studying the absorption spectra of CdSe material in bulk and in quantum dot form [14], found the appearance of quantized bands and an increase in band gap of CdSe in quantum dots. Further, the researcher also obtained an expression for the size dependent energy gap using the spherical "quantum box" model which is given below.

 $E_{a}^{band}(d) = E_{g}^{band}(bulk) + h^{2}/8m_{eh}d^{2}$

where

 $m_{eh} = m_e m_h / (m_e + m_h)$

and d is the diameter of circular/cylindrical material me is the effective mass of electron

 m_h is the effective mass of hole

In Equation 2, the parameter 'd' introduces the size based effects. Equation 2 can be written as [15].

$$E_{g}^{band}(d) = E_{g}^{band}(bulk) + K/d^{2}$$
(3)

It is clear from Equation 3 that second term in Equation 3 tends to increase as the diameter of the device decreases. It implies that the value of band gap will increase as the diameter of the device is reduced. As the value of band gap increases, following Equation 1, the cut-in voltage will also increase.

Hence, the reduction in diameter of the device leads to an increase in the band gap of Se, which is indicated by an increase in cut-in voltage of the device.

4.2 Tunneling Current

Tunneling current I can be expressed by Equation 4 [16,17]

$$I = \int T(E)\eta_m(E)\eta_s(E)[F_m(E)-F_s(E)] dE, \quad (4)$$

where

T(E) is tunneling probability between the occupied level in the

Cu metal and the unoccupied level in the Se semiconductor

 $\eta_m(E)$ or $\eta_s(E)$ is Density of states (DOS) of the metal and semiconductor, respectively

 F_m and F_s is Fermi-distribution function in metal and Semiconductor respectively

From the WKB (Wentzel-Kramers-Brillouin) approximation, the tunneling probability can be approximated as [17]

$$T(E) \approx \exp(-2kt)$$
 (5)

where

k is wave vector

t is width of the barrier

$$F_{\rm m}$$
 and $F_{\rm s} = 1/(1 + \exp \left(\frac{(E-E_{\rm f}/kT)}{f}\right) [18]$ (6)

where, E_f is Fermi energy of metal or semiconductor

Density of states in metals can be estimated by parabolic approximation, resulting in an $E^{1/2}$ dependency of the density of states [18].

 $\eta_m = 3.14/2 \times \text{volume} \times (8m_e \div h^2)^{3/2} \times E^{1/2}$ (7) where m_e is effective mass of an electron in the metal h is the Planck's constant

However, the small size of Se semiconductor implies the presence of (Columbic) charging energy states in addition to the density of states of the particles. Taking into account the size distribution of the materials, we can express the density of states of the Se material as [17].

$$\eta_{s} = \eta_{s}^{0} \sum_{n=0}^{n=\infty} \exp[-(E - nE_{c})2/(2\sigma^{2})/(\sqrt{2\pi\sigma^{2}})]$$
(8)

 η_s^0 is density of states without the charging states and is given by [19]

$$\eta_s^0 = \sqrt{2\pi m_e} / \sqrt{h^2} \times 1 / \sqrt{E}$$

where η_s is density of states of the Se

 E_c is charging energy of the Se

 $\boldsymbol{\sigma}$ is size-dependent standard deviation in energy space

As the spacing between the energy levels increases, σ will increase as the size of the semiconductor decreases. The various parameters for Cu-Se binary structures are given in Table 2.

The values of the various other parameters are given below.

Free mass of the electron (m) = 9.1×10^{-31} kg Planck's constant (h) = 6.602×10^{-34} J-s E_f(Cu) = 7.0 eV [23] E_f(Se) = 5.6 eV [21]

4.3. Calculation of Charging Energies for Se

Charging energy is the energy required to put a charge q on a conductor of capacitance C_0 and is given by [24,25].

$$E_c = e^2/2C_0 \approx e^2/4\pi \in e_0 \in d$$
(7)

Table 2. Parameters of	Cu, Ag and	Se materials
------------------------	------------	--------------

	Work function (eV)	Electron affinity (eV)	Effective mass (Kg)	References
Cu	4.7		1.46 m	[20]
Se	5.11	3.0	0.22 m	[20-22]
Ag	4.73			[20]

where

- \in_0 is permittivity of free space = 8.854×10^{-12} F/m
- \in_r is relative permittivity or dielectric constant of Se material = 6.1 [26]
- d is diameter of the semiconductor
- e is charge on an electron = 1.6×10^{-19} C

Substituting the values of diameters of different Cu-Se devices in Equation 7, correspondingly, charging energy comes out to be 0.0002 eV (1µm diameter), 0.002 eV (100 nm diameter) and 0.006 eV (40 nm diameter). Since capacitance C_0 is dependent directly on the diameter of the device, clearly the charging energy will increase as the diameter is reduced. Hence, an electron will be able to enter into the nanomaterial if it has enough charging energy and if it is in resonance with an empty state of the well. Substituting the values of parameters of Cu, Se and Ag in Equations 5, 6, 7 and 8, η_m , η_s , T(E) and F_(m or s) are calculated for different values of energies. Substituting the values of these terms in Equation 4, the variation of tunneling current in Cu-Se RTDs of diameters 1 µm, 100 nm and 40 nm for various values of voltages are calculated and plotted. The calculated I-V curves for different diameters are shown in Figures. 8, 9 and 10.

From Figures 8, 9 and 10, we observe that the theoretical I-V characteristics of the Cu-Se binary structures of diameters 1 μ m, 100 nm and 40 nm show a behavior similar to the one as seen in the experimental observations. However, the current values as calculated are very small as compared to experimental values of the currents. This type of behavior is expected, since the calculated I-V characteristics is for a single Cu-Se RTD, whereas, the experimental results are due to the collective behavior of a large number of Cu-Se RTDs in parallel.

Further, the PVCR of the Cu-Se RTDs are calculated for different diameters and are shown in tabular form in Table 3.

Calculated values of PVCR of the Cu-Se RTDs of diameters 1μ m, 100 nm and 40 nm are 2.01, 2.67 and 3.14, which show an increase in PVCR with decrease in diameters. This behavior is similar to that seen in the I-V curves obtained experimentally.

Hence, it can be inferred that the results obtained from theoretical model of RTD show a behavior similar to that obtained experimentally. However, the theoretical device currents are small in values because, in experimental set-up, several devices are working in parallel while, in theoretical equations, the current for a single device is calculated.

Table 3. Variation of calculated PVCR with diameter.

Diameter (nm)	PVCR
40	3.14
100	2.67
1000	2.01



Figure 8.Theoretical I-V characteristics of Cu-Se RTD of 1 μm diameter



Figure 9. Theoretical I-V characteristics of Cu-Se RTD of 100 nm



Figure 10. Theoretical I-V characteristics of Cu-Se RTD of 40 nm diameter

5. Conclusions

A suitable model for the template synthesized Cu-Se RTDs is proposed. Experimental results have been verified with the help of equations formulated for this model. The results obtained from theoretical model of RTD show a behavior similar to that obtained experimentally. However, the theoretical device currents are small in values and PVCRs show deviation in their values from the experimental values. This is because, in experimental set-up, several devices are working in parallel while, in theoretical equations, the current for a single device is calculated.

REFERENCES

[1] J. P. Sun, G. I. Haddad, P. Mazumde, and J. N. Schulman,

Proceedings of the IEEE, Vol. 86, No. 4, pp. 641, 1998.

[2] O. I. Mićić and A. J. Nozik, in Hari Singh Nalwa (Ed.), Colloidal Quantum Dot of III-V Semiconductors, Handbook of Nanostructured Mateials and Nanotechnology,

Academic Press, 2000.

- [3] M. Chaudhri, A. Vohra, S. K. Chakarvarti, and R. Kumar, J. mater. Sci., Mater Electron, Vol. 17, pp. 189, 2006.
- [4] M. Chaudhri, A. Vohra, and S. K. Chakarvarti, J. Mater. Sci., Mater Electron, Vol. 17, pp. 993, 2006.
- [5] M. Chaudhri, A. Vohra, and S. K. Chakarvarti, Physica E, Vol. 40, pp. 849, 2008.
- [6] M. Chaudhri, A. Vohra, and S. K. Chakarvarti, Mater. Sci. Engg. B, Vol. 149, No. 7, pp. 641, 2008.
- [7] H. D. Vladimir Gavryushin, Functional Combinations in Solid States, 2002. http://www.mtmi.vu.lt/pfk/funkc_dariniai/index.html.
- [8] V. V. Moshchalkov, V. Bruyndoncx, L. L. Van, M. J. Van Bael, Y. Bruynseraede, and A. Tonomura, in Hari Singh Nalwa (Ed.), Quantization and Confinement Phenomena in Nanostructured Superconductors, Handbook of Nanostructured Mateials and Nanotechnology, Academic Press, 2000.
- [9] Y. J. Choi, I. S. Hwang, J. H. Park, S. Nahm, and J. G. Park, Nanotechnology, Vol. 17, pp. 3775, 2006.
- [10] J. Heremans, C. M. Thrush, Y. M. Lin, S. Cronin, Z. Zhang, M. S. Dresselhaus, and J. F. Mansfield, Phys. Rev. B, Vol. 61, pp. 2921, 2000.
- [11] M. Li and J. C. Li, , Mater. Lett. Vol. 60, pp. 2526, 2006.
- [12] S. Cronin, Z. Zhang, and M. S. Dresselhaus, Phys. Rev. B, Vol. 61, No. 4, pp. 2921, 2000.
- [13] S. M. Sze, Physics of Semiconductor Devices. New York,

Wiley, 1981.

- [14] V. I. Klimov, Vol. 28, pp. 215, 2003.
- [15] S. Ogut, J. R. Chelikowsky, and S. G. Louie, Phys. Rev. Lett., Vol. 79, pp. 1770, 1997.
- [16] A. Sigurdardottir, V. Krozer, and H. L. Hartnage, Appl. Phys. Lett., Vol. 67, No. 22, pp. 3313, 1995.
- [17] S. H. Kim, G. Markovich, S. Rezvani, S. H. Choi, S. H., K. L. K. L. Wang, and J. R. Heath, Appl. Phys. Lett., pp. 317, 1999.
- [18] B. G. Streetman, "Solid state electronic devices," Prentice-Hall of India Private Limited, New Delhi, 1994.
- B. V. Zeghbroeck, "Principles of semiconductor devices," 2004. http://ece-www.colorado.edu/~bart/book.
- [20] P. A. Tipler and R. A. Liewellyn, Modern Physics, 3rd Editon, W. H. Freeman, 1999.
- [21] K. Barbalace, 2006. http://Klbprouctions.com/Periodic Table of Elements-Selenium-Se, Environmental Chemistry.com, 1995–2006. Accessed online: 7/13/2006. http:// Environmental Chemistry.com//yogi/Periodic/Se.html.
- [22] C. M. Fang, R. A. De Groot, and G. A. wiegers, Journal of Physics and Chemistry of Solids, Vol. 63, pp. 457, 2002.
- [23] N. W. Ashcroft and N. D. Mermin, Solid State Physics, Saunders, 1976.
- [24] A. J. Quinn, P. Beecher, D. Iacopino, L. Floyd, G. De-Marzi, E. V. Shechenko, H. Weller, and R G. edmond, Small 1, 613. Vol. 1, pp. 613, 2005.
- [25] S. Möller, H. Buhmann, S. F. Godijn, and L. W. Molenkamp, Phys. Rev. Lett., Vol. 81, No. 23, pp. 5197, 1998.
- [26] Dielectric Constant References Guide: http://www.asiinstr.com/technal/Dielectric%20Constants.htm.



Live Video Services Using Fast Broadcasting Scheme

Satish Chand

Computer Engineering Division, Netaji Subhas Institute of Technology, Dwarka, New Delhi, India E-mail: schand86@hotmail.com Received November 17, 2009; accepted December 29, 2009

Abstract: The Fast Broadcasting scheme is one of the simplest schemes that provide video services. In this scheme, the video is divided into equal-sized segments depending upon the bandwidth allocated by the video server. If the video length is not known, then this scheme cannot be applied as the number of video segments cannot be determined. In a live video wherein the video size is unknown, especially the ending time of the live broadcast, e.g., cricket match, this scheme cannot be applied. In this paper, we propose a model that helps the Fast Broadcasting scheme to support live video broadcasting. The basic architecture of the system consists of a live system with one video channel that broadcasts the live video and a video server that broadcasts the already broadcast live video to users.

Keywords: fast broadcasting scheme, live video channel, channel transition

1. Introduction

Video-on-Demand (VOD) services are one of the important classes that has several potential applications, such as entertainment, advertisement, distance education, etc. In spite of vast usability, these applications could not gain much attention because the earlier network technologies were not sufficient enough to support high data rate and same was with the storage technologies. These technologies have been developed significantly and are further being enhanced. New applications are also being explored that again put high demand on these resources. Therefore, efficient utilization of the resources especially the bandwidth is must. Several schemes exist in literature. but all are meant for the stored videos. These schemes require the video length to construct its segments. If the video length is not known in advance, which is generally the case for live videos, these schemes can be applied. To develop a broadcasting scheme for live videos is the motivation of this work. In this paper, we propose a mechanism that helps the Fast Broadcasting scheme to support the live video services. The Fast Broadcasting scheme is one of the simplest schemes that provides the video services. For broadcasting a video, there is a need to have a video server to transmit the video data. Designing a video server has many issues, e.g., efficient system design [1–3], storage management [4–7], broadcasting techniques. The broadcasting techniques are very important as they help utilizing the bandwidth efficiently. Some of the important broadcasting techniques are discussed in [8-9].

The basic model in this paper consists of a system (we call it as a live system) with a video channel that is called as the live channel. This system is active for the duration of the live video and broadcasts the live video data only once using its channel. There is another system that stores the video data from the live channel in its buffer and then broadcasts. This system simultaneously stores the new data from the live channel into its buffer and broadcasts the already stored data by using its channels. This process continues for the duration of the live video transmission. When the live video is over, the data is broadcast by the video server only. If the live video broadcast is still there and the video server has no free channel to broadcast the newly downloaded data, then the last channel of the video server is made free by transferring its data to other channels and the last channel can broadcast the newly downloaded data. This mechanism is called the channel transition mechanism. The important issue in channel transition is that the current as well as future users should get continuous data delivery. In literature, there does not seem to appear any work that discusses live video broadcasting. It is perhaps that the broadcasting schemes existing in literature require the video in terms of prespecified number of segments depending on the allocated bandwidth. To determine the number of video segments of a live video is not possible because the video length is not known. In this paper, we develop a mechanism so that the Fast Broadcasting scheme can support the live videos. The remaining of the paper is organized as follows. Section 2 reviews the previous work. Section 3 proposes a system model to support the live video transmission. The main concept in this paper is channel transition, which is of two types intermediate and final channel transition. In intermediate channel transition the live video is not over, whereas in the final channel transition the live video is over. Section 4 presents the results and finally Section 5 concludes the paper.

2. Previous Work

The broadcasting schemes are an important class of protocols supporting the video-on-demand services. Some schemes use a segment as the basic data unit for transmission and a video channel a basic transmission unit. Some schemes further divide the segments into subsegments and the video channels into subchannels. A subchannel transmits all the subsegments of a particular single segment. Taking subsegment as a basic data unit and subchannel as a basic transmission unit reduce the resources requirement. This however increases the complexity. Some of the important broadcasting schemes include harmonic scheme [11], geometrico-harmonic scheme with continuous delivery [12], Pyramid Broadcasting scheme [13], Fast Broadcasting scheme [14]. The harmonic and geometrico-harmonic schemes have their basic data and transmission units as subsegments and subchannels, respectively. The pyramid and Fast Broadcasting schemes employ the segments and video channels as their basic data and transmission units, respectively. A video channel is a logical channel that has bandwidth equal to the consumption rate of the video. The Fast Broadcasting scheme is one of the simplest schemes for providing the video services. In this scheme, the bandwidth is equally-divided into channels, each having bandwidth equal to the video viewing rate. The video is also uniformly divided into segments. The first video channel transmits the first segment, repeatedly, and the second video channel transmits the second & third segments, alternately and periodically. The *i*th video channel transmits 2^{i-1} segments from $S_{2^{i-1}}$ to $S_{2^{i-1}}$, sequentially and periodically. The segment size determines the user's waiting time. The video transmission time is also divided into fixed time units, called time slots. In a time slot, a segment can be exactly viewed at the viewing rate. The number of segments of the video of length D for allocating K video channels is 2^{K} -1. Figure 1 shows transmission of the video segments over four video channels, denoted by C₁, C₂, C₃, and C₄, in the Fast Broadcasting scheme.

3. Live Fast Broadcasting Scheme

In live video broadcasting, the video size is not known in the beginning and hence its number of segments cannot be determined. The Fast Broadcasting scheme needs the number of segments in advance, so this scheme cannot

Time	slots														
	T_1	T_2	${ m T}_3$	T_4	T_5	T_6	T_7	T_8	Т9	$T_{10} \\$	$T_{11} \\$	$T_{12} \\$	T ₁₃	$T_{14} \\$	
C ₁	S1	S ₁	S ₁	S ₁	S1	S ₁									
C ₂	S ₂	S ₃	S ₂	S ₃	S ₂	S3	S ₂	S ₃							
C3	S ₄	S 5	S ₆	\mathbf{S}_7	S4	S ₅	S ₆	S ₇	S ₄						
C ₄	S ₈	S9	S ₁₀	S ₁₁	S ₁₂	S ₁₃	S ₁₄	S ₁₅	S ₈	S9	S ₁₀	S ₁₁	S ₁₂		

Figure 1. Data transmission in Fast Broadcasting scheme

be applied. In order to use this scheme for live video transmission, we make a simple modification in its ar chitecture. We transmit N number of segments using K video channels, which are related by

$$N = 2^{K} - 2 \tag{1}$$

The basic system model for supporting the live videos consists of a system, called the live system. This system broadcasts the live video by using its video channel (live channel). There is another system; we call it as the video server. The video server downloads the data from the live channel into its buffer in terms of fixed time durations, which are time slots. The data stored in a time slot is referred to as a video segment. The video server performs two activities. It stores new segments from the live channel into its buffer and simultaneously broadcasts the already stored segments. This server supports the video data to any user request, which misses the initial portion of the live video. A user request received after the live video has been started gets future data from the live channel and the initial missing data from the video server. The video server is always tuned to the live channel for new segments to store into its buffer. The stored segments are broadcast by the video server according to the following broadcasting procedure.

3.1 Broadcasting Procedure

The number of video segments for allocating K video channels is 2^{K} - 2. Denote the video segments by S_i (*i*=1,2,...,*N*). The first video channel transmits the first segment S₁, repeatedly, and the second video channel transmits the segments S₂ and S₃, alternately and repeatedly. The *i*th video channel transmits the segments $S_{2^{i-1}}$,

 $S_{2^{i-1}+1}$, $S_{2^{i-1}+2}$, ..., $S_{2^{i}-1}$, sequentially and periodically, provided this *i*th channel is not the last video

channel. The last video channel, i.e., *K*th video channel transmits the segments $S_{2^{K-1}}$, $S_{2^{K-1}+1}$, $S_{2^{K-1}+2}$, ..., $S_{2^{K}-2}$, sequentially and periodically.

The video server is always connected to the live channel for downloading the new segments. When a segment has been downloaded, the video server broadcasts that segment by using its channels along with other segments. This process continues till there is a free channel with the video server and the live video transmission is going on. Since the bandwidth allocated to the video by the video server is of fixed amount, this will get exhausted at some point of time. In that case, the video server would not be able to broadcast the newly downloaded segments. One solution to this problem is to increase the bandwidth, which may not always be possible. A better solution is to make the last channel free by transferring its segments to other video channels so that this video channel can broadcast the new segments. This process is called channel transition. The important issue while doing a channel transition is that all (present or future) user requests must get timely video data delivery. There are two types of channel transition: intermediate channel transition in which the live video is not over and the final channel transition in which the live video is over. We first discuss the intermediate channel transition.

3.2 Intermediate Channel Transition

The intermediate channel transition is performed when the video server is downloading new segments from the live system, but it does not have a free channel to broadcast them. We transmit $(2^{K}-2)$ segments using K channels, not $(2^{K}-1)$ as required in the Fast Broadcasting scheme.

This is done because the capacity of the last channel is equal to total capacity of all other channels. Here the 'capacity' of a channel refers to the maximum number of segments transmitted by that channel. While transferring segments of the last video channel to other channels to make it free, we simply need to double the segments' size. If $S_1, S_2, \dots, S_k \dots$ are the old segments, then the new segments, denoted by S_{1}^{1} , S_{2}^{1} ,..., are given by $S_{1}^{1} = S_{1}$ $+ S_2, S_2^1 = S_3 + S_4, \dots, S_k^1 = S_{2k-1} + S_{2k}, \dots, and so forth.$ After a channel transition the segment size (i.e., new waiting time) becomes twice of the old one (old waiting time). To avoid the increase in the waiting time, we should delay the channel transition as much as possible, which can be delayed till all channels have their capacity full. It is not difficult to show that by delaying a channel transition maximally, all user requests (before and after the channel transition) get the video data in time.

Figure 2 shows the first channel transition at thick black line for allocating four video channels to the video by the video server. The gray-colored channel signifies the live channel and remaining are the video server's channels. The problem of data availability may be for a request that begins viewing the video just before the channel transition. This request, denoted by R_{13} in Figure 2, begins downloading the data from the time slot TS_{14} . This request R_{13} downloads the segments S_{15} onward from the live channel and S_1 to S_{14} from the video server. The segments available for downloading from the video server using the 1^{st} , 2^{nd} , 3^{rd} , and 4^{th} video channels in the time slot TS_{14} are S_1 , S_2 , S_6 , and S_{14} , respectively. The segment S_1 can be viewed while downloading from the video server and does not requires storage space. Just

						Reques	St IC13						
Time	slots						Ļ		s1	[¹ 1 →	ST ¹ 2	•	
ST_0	ST_1	ST_2		ST_{10}	ST_{11}	ST_{12}	ST_{13}	ST_{14}	ST_{15}	ST_{16}	ST ₁₇ ST ₁₈		
S ₁	S ₂	S ₃		S ₁₁									
				-									
S_0	S_1	S_1		S_1	S1	S_1	S_1	S1	S^1	L	S ¹ 1	S ¹ 1	
		S_2		S ₂	S ₃	S ₂	S ₃	S ₂	S_2^1	2	S ¹ ₃	S ¹ ₂	
			1						-		1		
				S_6	S_7	S_4	S_5	S ₆	S^1_4	l .	S ¹ ₅	S^{1}_{6}	
									-				
				S ₁₀	S ₁₁	S ₁₂	S ₁₃	S ₁₄	S^{1}_{s}	3	S ¹ 9	S ¹ 10	

Doguost D

Figure 2. First channel transition

Γime sl ST₀ ^{t-1}	ots ST1 ^{t-1}	ST2 ⁶⁻¹	$\mathrm{ST_3}^{\mathfrak{l}\text{-}1}$	ST4 ⁶⁻¹	ST5 ⁶⁻¹		$\begin{matrix} R_8 \\ \downarrow \\ ST_8^{t-1} \end{matrix}$	$\begin{matrix} R_9 \\ \downarrow \\ ST_9^{\ell-1} \end{matrix}$	ST10 ⁶⁻¹	ST_1^{ℓ}	ST ₂	ST ₃	ST	t ST:	^c ST _c	5 [€] ST7	¢	
$S_1^{\ell-1}$	S2 ^{ℓ-1}	S3 ⁶⁻¹	S4 ^{ℓ-1}	S5 ⁶⁻¹	S ₆ ^{ℓ-1}	 	S9 ⁶⁻¹											
	S1 ⁶⁻¹	S1 ⁶⁻¹	S1 ^{ℓ-1}	S1 ⁶⁻¹	$S_1^{\ell-1}$	 	S ₁ ^{ℓ-1}	S1 ^{ℓ-1}	S1 ⁶⁻¹	$S_1{}^\ell$	$\mathbf{S}_1^{\mathfrak{l}}$	$\mathbf{S}_1^{\mathfrak{l}}$	$\mathbf{S_1}^{\mathfrak{l}}$	$\mathbf{S}_1^{\mathfrak{l}}$	S_1^{ℓ}	$\mathbf{S}_1^{\mathfrak{l}}$		
		S2 ⁶⁻¹	S3 ⁶⁻¹	S2 ⁶⁻¹	S3 ⁶⁻¹	 	S2 ⁽⁻¹	S3 ⁶⁻¹	S2 ⁶⁻¹	${S_2}^\ell$	$\mathbf{S_3}^{\mathfrak{l}}$	S_2^{ℓ}	$\mathbf{S_3}^{\ell}$	$\mathbf{S_2}^{\mathfrak{l}}$	$\mathbf{S_3}^{\ell}$	S_2^{ℓ}		
				S4 ⁶⁻¹	S5 ⁶⁻¹	 	S4 ^{ℓ-1}	S5 ⁶⁻¹	S6 ^{t-1}	${S_4}^\ell$	$\mathbf{S}_5^{\mathfrak{l}}$	$\mathbf{S_6}^{\mathfrak{l}}$	${\mathbf S_7}^{\mathfrak l}$	$\mathbf{S_4}^{\mathfrak{l}}$	$\mathbf{S}_5^{\mathfrak{l}}$	$\mathbf{S_6}^{\mathfrak{l}}$		
						 	$S_8^{\ell-1}$	S ₉ ^{ℓ-1}	S3 ⁶⁻¹	${\mathbf S_8}^\ell$	$\mathbf{S_9}^{\mathfrak{l}}$	$\mathbf{S}_{10}^{\ \ell}$	$\mathbf{S}_{11}^{\ \ell}$	\mathbf{S}_{12}^{ℓ}	S ₁₃ ^ℓ	$\mathbf{S}_{14}^{\mathfrak{l}}$		

r igure 5. r inai channel transiti	Figure 3	3. Final	channel	transitio
------------------------------------	----------	----------	---------	-----------

after the channel transition, R_{13} needs the segment S_2 for viewing and it is already in its buffer because it had been stored just before the channel transition. After viewing S_2 , R_{13} needs S_3 segment which is the first half part of the second segment S_2^1 (= $S_3 + S_4$) transmitted by the second channel after channel transition. Thus, the data of the segment S_3 can be made available to R_{13} . Using similar discussions, we can show that a request received in any time slot can receive timely video data. We now discuss the final channel transition.

3.3 Final Channel Transition

In order to carry out the final channel transition, the video size must be known and it is possible only when the live video transmission is over. Since the video size is known, we can construct the video segments. The live video can be over at any point of time. If the data broadcast by the live channel does not form a complete segment, we can add dummy data to make it a complete segment. The live video can be over in any time slot, which means that the last video channel can have any number of segments ranging from one to its maximum capacity. If its capacity is full, then nothing is required to do. To carry out the final channel transition, the last channel must have at least one segment and at least one empty time slot. To occupy empty time slots with the video data, we need to increase the number of segment to $(2^{K} - 2)$ so that all time slots of all the video channels K are occupied. Increasing the number of segments decreases the segment size as the video size is of fixed length. The first channel C_1 transmits $S_1^{\,\ell-1}$ and the second channel C_2 transmits the segments $S_2^{\,\ell-1}$ and $S_3^{\,\ell-1}$ just before the final transition. Here the superscript ℓ of a segment (time slot) signifies the segment (time slot) after the

final (last) channel transition and (ℓ -1) for a segment (time slot) just before the final channel transition. After the final channel transition, the segment size decreases, i.e., some last portion of $S_1^{\ell-1}$ segment is added to the beginning of $S_2^{\ell-1}$ and some last portion $S_2^{\ell-1}$ is added to the beginning of $S_3^{\ell-1}$, and so forth. The number of new segments is ($2^K - 2$) and they can occupy all the time slots on all channels. While carrying out this process, timely video data delivery to the current as well as future requests must be ensured.

We now discuss how the video data delivery can be timely maintained to all user requests by taking an example. Consider that the video server has allocated four video channels to the video. The last (i.e., fourth) video channel can accommodate the segments $S_8,\ S_9,\ldots,\ S_{14}.$ The live video can be over in any time slot ${\rm ST_i}^{\ell-1}$ (7 < i < 14). Let i=8, i.e., the last video channel has to transmit the last segment as $S_9^{\ell-1}$ in the time slot $ST_8^{\ell-1}$. The segment $S_9^{\ell-1}$ will be available at the video server for transmission in the time slot $ST_9^{\ell-1}$ (refer Figure 3). We can carry out the channel transition immediately after the time slot $ST_9^{\ell-1}$. The request R_8 that begins to download the data from the time slot $ST_9^{\ell-1}$ onward can download, if required, the segments $S_1^{\ell-1}$, $S_3^{\ell-1}$, $S_5^{\ell-1}$, and $S_9^{\ell-1}$ in that time slot (refer Figure 3). The segment $S_1^{\ell-1}$ is viewed while downloading in the time slot $ST_9^{\ell-1}$ and in the next time slot (i.e., after channel transition) this request needs $S_2^{\ell-1}$. The segment $S_2^{\ell-1}$ is distributed among the segments S_2^{ℓ} and S_3^{ℓ} after the channel transition. The segment S_2^{ℓ} is available for downloading just after the channel transition, but the initial portion of segment S_2^{ℓ} comprises some last portion of $S_1^{\hat{t}-1}$ and the remaining is some initial portion of the segment $S_2^{\ell-1}$. It means that just after channel transition the initial portion of S_2^{ℓ} that

is from the segment $S_1^{\ell-1}$ will be available, not the segment $S_2^{\ell-1}$. Thus, the request R_8 will not get the data of the segment $S_2^{\ell-1}$ in time. To circumvent this problem, we delay the final channel transition one time slot after the live video is over. In the instance case, we carry out final channel transition at the end of the time slot $ST_{10}^{\ell-1}$, not $ST_9^{\ell-1}$. By doing so, we have one empty time slot (shown as gray-colored on the last video channel in Figure 3) on the last video channel that can be used to broadcast $S_2^{\ell-1}$ or $S_3^{\ell-1}$ depending upon the last segment broadcast by the live channel is even or odd. The segments available for downloading to the request R_8 in the time slot $ST_{10}^{\ell-1}$ are $S_2^{\ell-1}$ and $S_6^{\ell-1}$. The request R_8 has segments $S_2^{\ell-1}$ and $S_3^{\ell-1}$ in its buffer and will need $S_4^{\ell-1}$ for viewing in the time slot $ST_{12}^{\ell-1}$, which is not in the buffer storage. After the channel transition, the segment $S_4^{\ell-1}$ is distributed among S_5^{ℓ} , S_6^{ℓ} , and S_7^{ℓ} and the time slot $ST_{12}^{\ell-1}$ is distributed in the time slots ST_2^{ℓ} , ST_3^{ℓ} , and ST_4^{ℓ} . The segment S_5^{ℓ} can be downloaded in the time slot ST_2^{ℓ} and the segments S_6^{ℓ} and S_7^{ℓ} in the time slots ST_3^{ℓ} and ST_4^{ℓ} , respectively. It may be noticed that the segment S_5^{ℓ} is available at the beginning of the time slot ST_2^{ℓ} , but some initial portion of ST_2^{ℓ} comprises the last portion of the time slot $ST_{11}^{\ell-1}$. Thus, the request R₈ can get the video data in time. Consider another request R9 that begins downloading the sider another request R_9 that begins downloading the video data from the time slot $ST_{10}^{\ell-1}$ onward and can download, if required, the segments $S_2^{\ell-1}$, $S_6^{\ell-1}$, and $S_3^{\ell-1}$. The request R_9 requires the segment $S_4^{\ell-1}$ in the time slot $ST_{13}^{\ell-1}$. The segment $S_4^{\ell-1}$ is distributed among S_5^{ℓ} , S_6^{ℓ} , and S_7^{ℓ} and the time slot $ST_{13}^{\ell-1}$ becomes a part of the time slots ST_4^{ℓ} and ST_5^{ℓ} . The segments S_5^{ℓ} , S_6^{ℓ} , and S_7^{ℓ} can be downloaded in the time slots ST_2^{ℓ} , ST_3^{ℓ} , and ST_4^{ℓ} , respectively. Thus, the segment $S_4^{\ell-1}$ can be made available in time to request R₉. Using similar discussions, we can show that all requests can be provided the required data in time.

We now find out those segments S_i^{ℓ} that have the data of the segment $S_4^{\ell-1}$ after the final channel transition. The indices of the segments S_i^{ℓ} that have the data of segment $S_4^{\ell-1}$ are I_L , I_{L+1} , ..., I_H , where I_L and I_H are given by $I_L =$ n_1 such that $\min_{n_1} \left\lfloor \frac{n_1 * p}{N} \right\rfloor \ge 3$ and $I_H = n_2$ such that

 $\min_{n_2} \left\lfloor \frac{n_2 * p}{N} \right\rfloor \ge 4 \text{ where } p \text{ is the index of the last segment}$ broadcast by the live channel and N is the number of

segments that is given by (1). $\sum_{i=1}^{n} d_{i} = 0$

For the request R₉ that receives the data from ST₁₀^{*l*-1} time slot onward, assuming that the last segment broadcast by the live channel is S₉^{*l*-1}, the segment S₄^{*l*-1} would be distributed among the segments S₅^{*l*}, S₆^{*l*}, and S₇^{*l*} as I_L = 5 and I_L = 7 because for n₁ = 5, $\min_{n_1} \left\lfloor \frac{n_1 * 9}{14} \right\rfloor \ge 3$ and for n₂ = 7, $\min_{n_2} \left\lfloor \frac{n_2 * 9}{14} \right\rfloor \ge 4$ hold. This can easily be verified as follows.

$$\begin{split} S_{1}^{\ell} &= 9^{*}S_{1}^{\ell-1}/14; \ S_{2}^{\ell} &= 5^{*}S_{1}^{\ell-1}/14 + 4^{*}S_{2}^{\ell-1}/14; \\ S_{3}^{\ell} &= 9^{*}S_{2}^{\ell-1}/14; \ S_{4}^{\ell} &= S_{2}^{\ell-1}/14 + 8^{*}S_{3}^{\ell-1}/14; \\ S_{5}^{\ell} &= 6^{*}S_{3}^{\ell-1}/14 + 3^{*}S_{4}^{\ell-1}/14; \ S_{6}^{\ell} &= 9^{*}S_{4}^{\ell-1}/14; \ S_{7}^{\ell} &= 2^{*}S_{4}^{\ell-1}/14 + 7^{*}S_{5}^{\ell-1}/14. \end{split}$$

The video channels allocated by the video server to the video transmit new segments S_i^{ℓ} (i=1,2,...,N) in the new time slots ST_i^{ℓ} (i=1,2,...,N,...) according to the broadcasting procedure. Thus, we have discussed channel transition. In next Section, we discuss the results.

4. Results

The Fast broadcasting scheme is one of the simplest broadcasting schemes. That's why we have considered it for live video broadcasting. In its architecture, we have made a simple modification so that the number of segments transmitted by its last video channel is equal to the total segments transmitted by its remaining channels. This modification makes the final channel transition at the optimal time point. In other words, the channel transition is performed after all time slots of all the video channels are full with the video segments. Consider again Figure 2. The request R_0 arrived in the 0th time slot ST_0 begins downloading the segments S_2 onward from the live channel and S_1 from the video server. The buffer requirement for this request is one segment. The request R_1 arrived in the 1th time slot ST₁ begins downloading the segments S₃ onward from the live channel into its buffer and S1 and S2 from the video server. Its buffer requirement is of two segments. We describe the buffer computation for an arbitrary request. Consider an arbitrary request, say, R₉, that arrives in the 9th time slot ST₉. This request begins downloading the data from the time slot ST_{10} onward. The segments S_{11} onward are downloaded from the live channel and the segments S_1 to S_{10} from the video server. The segments available for downloading, actually downloaded, and that required for viewing, in different time slots are shown in Table 1.

In last column '+' and '-' signs signify that the segment is stored in and read from the buffer, e.g., +2-1+1=2 means 2 segments are stored from the video server into the user's buffer, one is read from the user's buffer, and one is stored from the live channel into the user's buffer. Thus, the net buffer requirement is of two segments. The segment S₁ is viewed while downloading.

Using similar discussions, we can compute the buffer requirement for any request. Table 2 shows the buffer requirement for different requests (referring Figure 2), assuming that four video channels have been allocated to the video by the video server.

In Table 2, for the requests R_7 , R_8 , R_9 , R_{10} , and R_{11} , there are two different storage requirements. The first requirement is one when the live video is going on. The second requirement refers to one when the live video is over after the 14th segment. The maximum user's waiting

Time slot	Segments available for storing from Video Server	Segments stored from Video Server	Segments stored from live channel	Segment required for viewing	Total segments required				
ST_{10}	$S_1 + S_2 + S_6 + S_{10}$	S_2	S ₁₁	\mathbf{S}_1	+1+1=2				
ST_{11}	$S_{3}+S_{7}$	S_3	S_{12}	S_2	+1-1+1=1				
ST_{12}	S_4	\mathbf{S}_4	S_{13}	S_3	+1-1+1=1				
ST_{13}	S_5	S_5	S_{14}	S_4	+1-1+1=1				
ST_{14}	S_6	S_6	S_{15}	S_5	+1-1+1=1				
ST_{15}	$S_4^{-1}/2 = S_7$	S_7	S_{16}	S_6	+1-1+1=1				
ST_{16}	$S_4^{1/2} = S_8$	S_8	S_{17}	S_7	+1-1+1=1				
ST_{17}	$S_5^1/2 = S_9$	S_9	S_{18}	S_8	+1-1+1=1				
ST_{18}	$S_5^1/2 = S_{10}$	S_{10}	S19	S_9	+1-1+1=1				
	Buffer Storage required for request $R_9 = 10S$								

Table 1. Segments stored from the live channel and the video server by request R₉

Table 2. Buffer Requirement for	different Requests	allocating four	video channels
---------------------------------	--------------------	-----------------	----------------

Request	Buffer Requirement	Request	Buffer Requirement
R ₀	S	R_6	78
R_1	28	R ₇	8S or 7S
R_2	38	R_8	9S or 5S
R ₃	4S	R ₉	10S or 5S
R_4	58	R ₁₀	11S or 5S
R_5	6S	R ₁₁	12S or 5S

time in this scheme is pre-decided for the initial users and remains the same till the channel transition time. After every channel transition except the final one the user's waiting time becomes double of the previous one. The final channel transition is done only when the live video transmission is over. After the final channel transition the user's waiting time is stabilized and the scenario becomes like a stored video. The size of a segment (time slot) after a channel transition except the final one becomes double. If S_i, ST_i and S¹_i, ST¹_i are the *i*th segments and time slots, respectively, before and after the first channel transition (assuming four video channels), then we have the following relations:

 $S_{i}^{1} = S_{14+2i-1} + S_{14+2i}$ $ST_{i}^{1} = ST_{14+2i-1} + ST_{14+2i}$ In general, we have

$$S_{i}^{k} = S_{14+2i-1}^{k-1} + S_{14+2i}^{k-1}$$
, for $k = 1, 2, ..., \ell-1$, (2)

$$ST_{i}^{k} = ST_{14+2i-1}^{k-1} + ST_{14+2i}^{k-1}$$
, for $k = 1, 2, ..., \ell-1$, (3)

where S_{i}^{0} , ST_{i}^{0} denote the very first *i*th segment and time slot, respectively, i.e., $S_i^0 = S_i$, $ST_i^0 = ST_i$ and ℓ denotes the final channel transition.

We can determine the size of a segment (time slot) after any channel transition in terms of the original segments (time slots). For example, consider the third channel transition (assuming it is not the final channel transition). Then, (2a) & (2b) give

$$S_{i}^{3} = S_{14+2i-1}^{2} + S_{14+2i}^{2}$$
 (4)

$$ST_{i}^{3} = ST_{14+2i-1}^{2} + ST_{14+2i}^{2}$$
 (5)

We can take i=1 because after any channel transition all the segments (time slots) are of same size. We will derive the relation for the segments only, and for the time slots exactly the same relation will hold. For i=3, we have from (3a) as

 $S_1^3 = S_{15}^2 + S_{16}^2$

We need to find out S_{15}^2 and S_{16}^2 , which are given by $S_{16}^{2} = S_{14+29}^{1} + S_{14+30}^{1} = S_{43}^{1} + S_{44}^{1}$ $S_{16}^{2} = S_{14+31}^{1} + S_{14+32}^{1} = S_{45}^{1} + S_{46}^{1}$ Again we need to find out S_{43}^{1} , S_{44}^{1} , S_{45}^{1} , and S_{46}^{1} ,

which are given by

hich are given by $S_{43}^{1} = S_{14+85}^{0} + S_{14+86}^{0} = S_{99}^{0} + S_{100}^{0}$ $S_{44}^{1} = S_{14+87}^{0} + S_{14+88}^{0} = S_{101}^{0} + S_{102}^{0}$ $S_{45}^{1} = S_{14+89}^{0} + S_{14+90}^{0} = S_{103}^{0} + S_{104}^{0}$ $S_{46}^{1} = S_{14+91}^{0} + S_{14+92}^{0} = S_{105}^{0} + S_{106}^{0}$ The segments S_{0}^{0} is are the original segments. Thus, we

have

 $S_{1}^{3} = S_{99} + S_{100} + \ldots + S_{106}$ For the time slots, we have

 $ST_{1}^{3} = ST_{99} + ST_{100} + \ldots + ST_{106}$

The segment (time slot) size after the final channel transition (i.e., ℓ th) is Θ times of the segment (time slot) size of that of the $(\ell - I)$ th channel transition, where 0.5 < $\Theta < 1$. Here $\Theta = 1$ signifies that the live video transmission is over when all time slots of all the video channels are full. In that case, nothing is done and the user's waiting time is unchanged. For $\Theta \neq 1$, the last channel after the final but one channel transition must have at least one

CN

The user's waiting time changes after a channel transition depending upon the segment size and this is fixed for the stored videos for a given bandwidth. In the proposed scheme, the initial user's waiting time is decided by the service provider and it is also equal to the segment size. This decision is necessary for arranging the video data that would be downloaded in terms of the segments and made on the basis of bandwidth allocated to the video. The video size increases after a new segment from the live channel has been downloaded, but this change affects broadcasting only after the channel transition. After the final channel transition, the user's waiting time generally decreases. The basic requirement to carry out the final channel transition is that there must be at least one segment and at least one empty time slot on the last video channel.

5. Conclusions

In this paper, we have proposed a technique so that the Fast Broadcasting scheme can be used for broadcasting the live videos. The Fast Broadcasting scheme does not support the live video services in its original form because it requires the number of segments of the video beforehand and for that the video length should be known, which is generally not known in advance in case of the live videos. In this proposed scheme, the live video data is stored in buffer at the video server in terms of fixed time durations, called time slots. The data downloaded in a time slot is considered as a segment. The downloaded segments are broadcast by the video server till there is free channel available with the video server. When no free channel is available and live video is there, channel transition is required. In the proposed scheme, the channel transition can be delayed maximally, i.e., up to the point when all time slots of all the video channels have been completely occupied. This study may be useful for designing a VOD system that can support the live video, such as cricket match, interactive education session, etc.

REFERENCES

 K. M. Ho and K. T. Lo, "Design of a decentralized videoon-demand system with cooperative clients in multicast environment," Advances in Multimedia Information Processing, Lecture Notes Computer Science, 4810, pp. 401–404, 2007.

- [2] Y. B. L. Jack, "Channel folding an algorithm to improve efficiency of multicast video-on-demand systems," IEEE Transactions on Multimedia, Vol. 7, No. 2, pp. 366–378, 2005.
- [3] W. F. Poon, K. T. Lo, and J. Feng, "A hierarchical video-on-demand system with double-rate batching," JVCIR, Vol. 16, No. 1, pp. 1–18, 2005.
- [4] D. J. Gemmell, H. M. Vin, D. Kandlur, P. V. Rangan, and L. A. Rowe, "Multimedia storage servers: a tutorial," Computer, Vol. 28, No. 5, pp. 40–49, 1995.
- [5] N. J. Sarhan and C. R. Das, "Caching and scheduling in nad based multimedia servers," IEEE Transactions on Parallel and Distributed Systems, Vol. 5, No. 10, pp. 921–933, 2004.
- [6] A. L. Chervnak, D. A. Patterson and R. H. Katz, "Choosing the best storage system for video service," In Proceeding of third ACM international conference on Multimedia, San Francisco, USA, pp. 109–119, 1995.
- [7] A. Dan and D. Sitaram, "Buffer management policy for an on-demand video server," IBM Watson Research Center, RC 19347, 1994.
- [8] C. C. Aggarwal, J. L. Wolf, and P. S. Yu, "A permutation based pyramid broadcasting scheme for video on demand systems," In Proceeding International Conference on Multimedia Computing and Systems, pp. 118–126, 1996.
- [9] L. Gao, J. Kurose, and D. Towsley, "Efficient schemes for broadcasting popular videos," In Proceeding of NOSS-DAV'98, pp. 183–194, 1998.
- [10] A. L. N. Reddy, J. Wyllie, and K. B. R. Wijayratne, "Disk scheduling in a multimedia I/O system," ACM Transactions on Computer Communication and Application (TOMCCAP), Vol. 1, No. 1, pp. 37–59, 2005.
- [11] L. S. Juhn and L. M. Tseng, "Harmonic broadcasting scheme for video-on-demand service," IEEE Transactions on Consumer Electronics, Vol. 43, No. 3, pp. 268–271, 1997.
- [12] H. Om and S. Chand, "Geometrico-harmonic broadcasting scheme with continuous redundancy," IEEE Transactions on Multimedia, Vol. 9, No. 1, pp. 410–419, 2007.
- [13] S. Viswanathan and T. Imielinski, "Metropolitan area video-on-demand service using pyramid broadcasting," ACM Multimedia System, Vol. 4, No. 4, pp. 197–208, 1996.
- [14] L. S. Juhn and L. M. Tseng, "Fast data broadcasting and receiving scheme for popular video service," IEEE Transactions on Broadcasting, Vol. 44, No. 1, pp. 100–105, 1998.

TABLE OF CONTENTS

Volume 2, Number 1, February 2010

Using the Power Control and Cooperative Communication for Energy Saving	
C. J. Chen, C. Jin, D. M. Li, J. C. Wang, J. N. Fang	1
A Comparison Study of Input ESD Protection Schemes Utilizing NMOS, Thyristor, and Diode Devices J. Y. Chio	11
On Possible A-Priori "Imprinting" of General Relativity itself on the Performed Lense- Thirring Tests with LAGEOS Satellites L. Iorio	26
Neural Network Performance for Complex Minimization Problem T. Wibg	. 31
Joint Power Control and Spectrum Allocation for Cognitive Radio with QoS Constraint Z. J. Zhao, Z. Peng, Z. D. Zhao, S. L. Zheng	38
A Historical Narrative of Study of Fiber Grating Solitons X. L. Li, Y. S. Jiang, L. J. Xu	44
ADPF Algorithm for Target Tracking in WSN C. H. Song, H. Zhao, W. Jing, D. Liu	50
Designing Intrusion Detection System for Web Documents Using Neural Network H. Om, T. K. Sarkar	54
On Solvable Potentials, Supersymmetry, and the One-Dimensional Hydrogen Atom R. P. Martínez-y-Romero, H. N. Núñez-Yépez, A. L. Salas-Brito	62
How to Measure in the Near Field and in the Far Field T. Dlugosz, H. Trzaska	65
Proposed Model for SIP Security Enhancement M. B. Sayyad, A. Chatterjee, S. L. Nalbalwar	69
A Model for Cu-Se Resonant Tunneling Diodes Fabricated by Negative Template Assisted Electrodeposition Technique M. Chaudhri, A. Vohra1, S. K. Chakarvarti	73
Live Video Services Using Fast Broadcasting Scheme S. Chand	79
Copyright©2010 SciRes	CN

