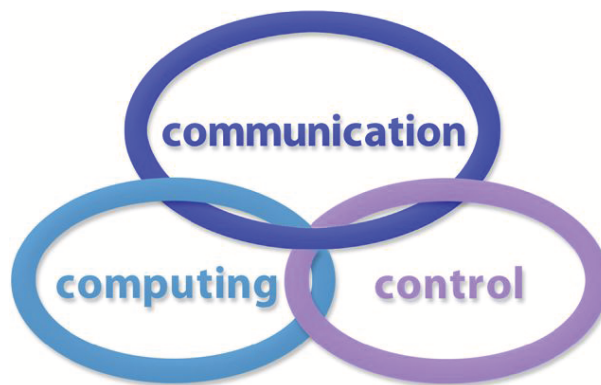


INTERNATIONAL JOURNAL
of
COMPUTERS, COMMUNICATIONS & CONTROL

ISSN 1841-9836

ISSN-L 1841-9836



A Bimonthly Journal
With Emphasis on the Integration of Three Technologies

Year: 2013 Volume: 8 Issue: 6 (December)

This journal is a member of, and subscribes to the principles of,
the Committee on Publication Ethics (COPE).



Agora University Editing House

CCC Publications

<http://univagora.ro/jour/index.php/ijccc/>

International Journal of Computers, Communications & Control



EDITOR IN CHIEF:

Florin-Gheorghe Filip

Member of the Romanian Academy
Romanian Academy, 125, Calea Victoriei
010071 Bucharest-1, Romania, ffilip@acad.ro

ASSOCIATE EDITOR IN CHIEF:

Ioan Dzitac

Aurel Vlaicu University of Arad, Romania
St. Elena Dragoi, 2, 310330 Arad, Romania
ioan.dzitac@uav.ro

&

Agora University of Oradea, Romania
Piata Tineretului, 8, 410526 Oradea, Romania
rector@univagora.ro

EXECUTIVE EDITOR:

Răzvan Andonie

Central Washington University, USA
400 East University Way, Ellensburg, WA 98926, USA
andonie@cwu.edu

MANAGING EDITOR DEPUTY MANAGING EDITOR

Mișu-Jan Manolescu

Agora University of Oradea, Romania
Piata Tineretului, 8, 410526 Oradea
mmj@univagora.ro

Horea Oros

University of Oradea, Romania
St. Universitatii 1, 410087, Oradea
horos@uoradea.ro

TECHNICAL SECRETARY

Cristian Dzitac

R & D Agora, Romania
rd.agora@univagora.ro

Emma Valeanu

R & D Agora, Romania
evaleanu@univagora.ro

EDITORIAL ADDRESS:

R&D Agora Ltd. / S.C. Cercetare Dezvoltare Agora S.R.L.
Piata Tineretului 8, Oradea, jud. Bihor, Romania, Zip Code 410526
Tel./ Fax: +40 359101032

E-mail: ijccc@univagora.ro, rd.agora@univagora.ro, ccc.journal@gmail.com
Journal website: <http://univagora.ro/jour/index.php/ijccc/>

International Journal of Computers, Communications & Control

EDITORIAL BOARD

Boldur E. Bărbat

Sibiu, Romania
bbarbat@gmail.com

Pierre Borne

Ecole Centrale de Lille
Cité Scientifique-BP 48
Villeneuve d'Ascq Cedex, F 59651, France
p.borne@ec-lille.fr

Ioan Buciu

University of Oradea
Universitatii, 1, Oradea, Romania
ibuciu@uoradea.ro

Hariton-Nicolae Costin

Faculty of Medical Bioengineering
Univ. of Medicine and Pharmacy, Iași
St. Universitatii No.16, 6600 Iași, Romania
hcostin@iit.tuiasi.ro

Petre Dini

Cisco
170 West Tasman Drive
San Jose, CA 95134, USA
pdini@cisco.com

Antonio Di Nola

Dept. of Mathematics and Information Sciences
Università degli Studi di Salerno
Salerno, Via Ponte Don Melillo 84084 Fisciano,
Italy
dinola@cds.unina.it

Ömer Egecioglu

Department of Computer Science
University of California
Santa Barbara, CA 93106-5110, U.S.A
omer@cs.ucsb.edu

Constantin Gaidric

Institute of Mathematics of
Moldavian Academy of Sciences
Kishinev, 277028, Academiei 5, Moldova
gaidric@math.md

Xiao-Shan Gao

Academy of Mathematics and System Sciences
Academia Sinica
Beijing 100080, China
xgao@mmrc.iss.ac.cn

Kaoru Hirota

Hirota Lab. Dept. C.I. & S.S.
Tokyo Institute of Technology
G3-49,4259 Nagatsuta, Midori-ku, 226-8502, Japan
hirota@hrt.dis.titech.ac.jp

Gang Kou

School of Business Administration
Southwestern University of Finance and Economics
Chengdu, 611130, China
kougang@yahoo.com

George Metakides

University of Patras
University Campus
Patras 26 504, Greece
george@metakides.net

Ștefan I. Nitchi

Department of Economic Informatics
Babes Bolyai University, Cluj-Napoca, Romania
St. T. Mihali, Nr. 58-60, 400591, Cluj-Napoca
nitchi@econ.ubbcluj.ro

Shimon Y. Nof

School of Industrial Engineering
Purdue University
Grissom Hall, West Lafayette, IN 47907, U.S.A.
nof@purdue.edu

Stephan Olariu

Department of Computer Science
Old Dominion University
Norfolk, VA 23529-0162, U.S.A.
olariu@cs.odu.edu

Gheorghe Păun

Institute of Mathematics
of the Romanian Academy
Bucharest, PO Box 1-764, 70700, Romania
gpaun@us.es

Mario de J. Pérez Jiménez

Dept. of CS and Artificial Intelligence
University of Seville, Sevilla,
Avda. Reina Mercedes s/n, 41012, Spain
marper@us.es

Dana Petcu

Computer Science Department
Western University of Timisoara
V.Parvan 4, 300223 Timisoara, Romania
petcu@info.uvt.ro

Radu Popescu-Zeletin

Fraunhofer Institute for Open
Communication Systems
Technical University Berlin, Germany
rpz@cs.tu-berlin.de

Imre J. Rudas

Institute of Intelligent Engineering Systems
Budapest Tech
Budapest, Bécsi út 96/B, H-1034, Hungary
rudas@bmf.hu

Yong Shi

Research Center on Fictitious Economy
& Data Science
Chinese Academy of Sciences
Beijing 100190, China
yshi@gucas.ac.cn
and
College of Information Science & Technology
University of Nebraska at Omaha
Omaha, NE 68182, USA
yshi@unomaha.edu

Athanasios D. Styliadis

University of Kavala Institute of Technology
65404 Kavala, Greece
styliadis@teikav.edu.gr

Gheorghe Tecuci

Learning Agents Center
George Mason University, USA
University Drive 4440, Fairfax VA 22030-4444
tecuci@gmu.edu

Horia-Nicolai Teodorescu

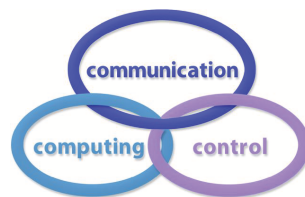
Faculty of Electronics and Telecommunications
Technical University "Gh. Asachi" Iasi
Iasi, Bd. Carol I 11, 700506, Romania
hteodor@etc.tuiasi.ro

Dan Tufiş

Research Institute for Artificial Intelligence
of the Romanian Academy
Bucharest, "13 Septembrie" 13, 050711, Romania
tufis@racai.ro

Lotfi A. Zadeh

Professor,
Graduate School,
Director,
Berkeley Initiative in Soft Computing (BISC)
Computer Science Division
Department of Electrical Engineering
& Computer Sciences
University of California Berkeley,
Berkeley, CA 94720-1776, USA
zadeh@eecs.berkeley.edu

**DATA FOR SUBSCRIBERS**

Supplier: Cercetare Dezvoltare Agora Srl (Research & Development Agora Ltd.)

Fiscal code: 24747462

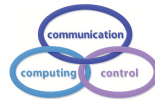
Headquarter: Oradea, Piata Tineretului Nr.8, Bihor, Romania, Zip code 410526

Bank: MILLENNIUM BANK, Bank address: Piata Unirii, str. Primariei, 2, Oradea, Romania

IBAN Account for EURO: RO73MILB000000000932235

SWIFT CODE (eq.BIC): MILBROBU

International Journal of Computers, Communications & Control



Short Description of IJCCC

Title of journal: International Journal of Computers, Communications & Control

Acronym: IJCCC

Abbreviated Journal Title: INT J COMPUT COMMUN

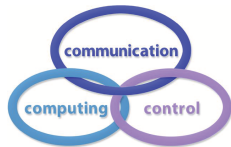
International Standard Serial Number: ISSN 1841-9836, ISSN-L 1841-9836

Publisher: CCC Publications - Agora University

Starting year of IJCCC: 2006

Founders of IJCCC: Ioan Dzitac, Florin Gheorghe Filip and Mişu-Jan Manolescu

Logo:



Publication frequency: Bimonthly: Issue 1 (February); Issue 2 (April); Issue 3 (June); Issue 4 (August); Issue 5 (October); Issue 6 (December).

Coverage:

- Beginning with Vol. 1 (2006), Supplementary issue: S, IJCCC is covered by Thomson Reuters - SCI Expanded and is indexed in ISI Web of Science.
- Journal Citation Reports(JCR)/Science Edition:
 - Impact factor (IF): JCR2009, IF=0.373; JCR2010, IF=0.650; JCR2011, IF=0.438; JCR2012, IF=0.441.
- Beginning with Vol. 2 (2007), No.1, IJCCC is covered in EBSCO.
- Beginning with Vol. 3 (2008), No.1, IJCCC, is covered in Scopus.

Scope: International Journal of Computers Communications & Control is directed to the international communities of scientific researchers in computer and control from the universities, research units and industry.

To differentiate from other similar journals, the editorial policy of IJCCC encourages the submission of scientific papers that focus on the integration of the 3 "C" (Computing, Communication, Control).

In particular the following topics are expected to be addressed by authors:

- Integrated solutions in computer-based control and communications;
- Computational intelligence methods (with particular emphasis on fuzzy logic-based methods, ANN, evolutionary computing, collective/swarm intelligence);
- Advanced decision support systems (with particular emphasis on the usage of combined solvers and/or web technologies).

Copyright © 2006-2013 by CCC Publications

Contents

Association Rule Mining using Path Systems in Directed Graphs S. Arumugam, S. Sabeen	791
Adaptive Network Coding Scheme for TCP over Wireless Sensor Networks Y.-C. Chan, Y.-Y. Hu	800
Operation Mechanism of the Driving Force System of Ecosystem of Cyber-society Based on the System Dynamics X. Guan, Z. Zhang, S. Zhang	812
Localization of Wireless Sensor Network Based on Genetic Algorithm N. Jiang, S. Jin, Y. Guo, Y. He	825
Active Queue Management of TCP Flows with Self-scheduled Linear Parameter Varying Controllers C. Kasnakoglu	838
An Improved Genetic Algorithm for the Multi Level Uncapacitated Facility Location Problem V. Korać, J. Kratica, A. Savić	845
A Improved EPC Class 1 Gen 2 Protocol with FCFS Feature in the Mobile RFID Systems X. Li, Y. Quan	854
Direct Method for Stability Analysis of Fractional Delay Systems M.A. Pakzad, M.A. Nekoui	863
Dimensionality Reduction and Generation of Human Motion S. Qu, L.D. Wu, Y.M. Wei, R.H. Yu	869
Security Ontology for Adaptive Mapping of Security Standards S. Ramanauskaitė, D. Olifer, N. Goranin, A. Čenys	878
RSEA-AODV: Route Stability and Energy Aware Routing for Mobile Ad Hoc Networks P. Srinivasan, P. Kamalakkannan	891
Design of Congestion Control Scheme for Uncertain Discrete Network Systems H. Wang, C. Yu	901
Author index	907

Association Rule Mining using Path Systems in Directed Graphs

S. Arumugam, S. Sabeen

S. Arumugam*

1. National Centre for Advanced Research in Discrete Mathematics (*n*-CARDMATH)

Kalasalingam University

Anand Nagar, Krishnankoil-626126, INDIA.

2. School of Electrical Engineering and Computer Science

The University of Newcastle

NSW 2308, Australia.

*Corresponding author: s.arumugam.klu@gmail.com

S. Sabeen

Department of Computer Applications

Jaya Engineering College

Chennai-600054, INDIA.

sabeens@rediffmail.com

Abstract: A transaction database (TDB) consists of a set I of items and a multiset \mathcal{D} of nonempty subsets of I , whose elements are called transactions. There are several algorithms for solving the popular and computationally expensive task of association rule mining from a TDB. In this paper we propose a data structure which consists of a directed graph D (loops and multiple arcs are permitted) and a system of directed paths in D to represent a TDB. We give efficient algorithms for generating the data structure, for extracting frequent patterns and for association rule mining. We also propose several graph theoretic parameters which lead to a better understanding of the system.

Keywords: Directed graphs, path system, in-degree, out-degree, association rule mining, frequent patterns, data mining.

1 Introduction

The task of association rule mining in a large database of transactions was proposed by Agarwal et al. [1]. Since then this problem has received a great deal of attention and association rule mining is one of the most popular pattern discovery methods in Knowledge Discovery from Database (KDD). A broad variety of efficient algorithms such as Apriori Algorithm [3], FP-growth [4], FP-tree [4], SETM [6], DIC [5] have been developed during the past few years. In this paper we propose a data structure consisting of a directed graph D and a multiset of directed paths in D to represent a database of transactions. We give an algorithm which generates the directed graph D and which also simultaneously computes several other measures such as in-degree, out-degree, total number of arcs, length of a largest transaction, frequency of occurrence of various nodes and the number of occurrences of each arc in D . The second algorithm generates all the patterns of the transaction database using the above data structure. This algorithm can be modified to extract frequent patterns. The third algorithm deals with association rule mining. In this process we scan the database exactly once and the digraph is constructed dynamically.

2 Directed Graphs and Path Systems

A *directed graph* $D = (V, A)$ consists of a finite nonempty set V and a multiset A of ordered pairs of elements of V . The elements of V are called *vertices* and the elements of A are called *arcs*. An ordered pair (v, v) is called a *loop* at v . We allow both loops and multiple arcs (that

is, an arc (u, v) appearing more than once) in D . If $a = (u, v)$ is an arc in D , then u is called the *tail* of a and v is called the *head* of a . For basic terminology in directed graphs we refer to the book by Chartrand and Lesniak [2]. A *directed path* in D is a sequence of distinct vertices $P = (v_1, v_2, \dots, v_k)$ such that (v_i, v_{i+1}) is an arc in D for all $i, 1 \leq i \leq k - 1$. A *directed cycle* in D is a sequence of vertices $C = (v_1, v_2, \dots, v_k, v_1)$ such that (v_i, v_{i+1}) is an arc in D for all $i, 1 \leq i \leq k - 1$, (v_k, v_1) is an arc in D and the vertices v_1, v_2, \dots, v_k are distinct. A directed graph is called *acyclic* if it contains no directed cycle. For any vertex v , the *in-degree* $id(v)$ is defined to be the number of arcs of the form (u, v) in D and the *out-degree* $od(v)$ is defined to be the number of arcs of the form (v, u) in D . We observe that a loop at v contributes 1 to both $id(v)$ and $od(v)$ and an arc (u, v) contributes 1 to $od(u)$ and 1 to $id(v)$. A vertex v such that no arc is incident with v or all the arcs incident at v are loops is called an *isolated vertex* of D . We denote by D_1 the subdigraph of D obtained by removing all the loops in D . A *path system* or a *path cover* in a directed graph D is a multiset ψ of directed paths in D such that every arc of D is in exactly one path in ψ . We adopt the convention that all loops in D are members of ψ . If D contains multiple arcs, a directed path in D may occur more than once in ψ .

3 Digraph model for transaction database

Let $I = \{1, 2, \dots, n\}$ be a set of objects whose elements are called items. We impose on I the natural ordering of the set of positive integers. Any nonempty subset of I is called a transaction. A transaction database \mathcal{D} is a multiset of transactions in I . Thus a subset X of I may occur more than once in \mathcal{D} . We assume that for each $i \in I$, there is at least one transaction in $T \in \mathcal{D}$ with $i \in T$. In other words $\bigcup_{T \in \mathcal{D}} T = I$. We say that a transaction $T \in \mathcal{D}$ supports an

item set $X \subseteq I$ if $X \subseteq T$. The support of X is defined by $supp(X) = \frac{|T \in \mathcal{D}: X \subseteq T|}{|\mathcal{D}|}$. Thus $supp(X)$ is the fraction of transactions supporting X . If we fix a threshold value s , then any subset X of I with $supp(X) \geq s$ is called a frequent pattern. An association rule is an implication of the form $X \Rightarrow Y$ where $X, Y \subseteq I$ and $X \cap Y = \emptyset$. The support of a rule $X \Rightarrow Y$ is defined to be $supp(X \cup Y)$. The confidence of the rule is defined as $conf(X \Rightarrow Y) = \frac{supp(X \cup Y)}{supp(X)}$. The support and confidence values are usually normalized so that these values occur between 0% and 100% instead of 0 to 1.0. Normally we generate association rules for frequent patterns. For further terminology in TDB and association rule mining we refer to the book by Han and Kamber [7].

In this paper we propose a data structure which consists of a directed graph D and a system of directed paths in D for representing a transaction database. The vertex set of the directed graph is the item set I . We sort the elements of each transaction $T \in \mathcal{D}$ in the increasing order and represent T as a directed path in D . Thus the direction on each edge is given by the orientation from low-to-high. Any transaction T of the form $\{x\}$ is represented as a loop at x . We illustrate this database with a small example.

Example 1. Let $I = \{1, 2, 3, 4, 5\}$. Let $\mathcal{D} = \{T_1, T_2, T_3, T_4, T_5, T_6\}$ where $T_1 = \{1, 2, 5\}$, $T_2 = \{1, 4\}$, $T_3 = \{1, 2, 4\}$, $T_4 = \{3, 5\}$, $T_5 = \{2\}$, $T_6 = \{3\}$ and $T_7 = \{3, 4, 5\}$. The directed graph of \mathcal{D} is given in Figure 1.

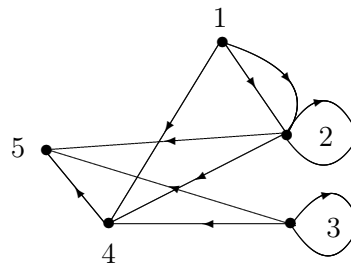


Figure 1

Since every transaction T in \mathcal{D} gives a directed path in D , it follows that the multiset \mathcal{D} of all transaction forms a path system in D . We observe that the directed graph D_1 which is obtained by removing all the loops in the directed graph D representing a TDB is an *acyclic directed graph*.

4 Directed graph of a transaction data base (DGTDB)

In this section we propose an algorithm to construct the directed graph representing a TDB. The algorithm scans the data exactly once, dynamically constructs the digraph D and simultaneously computes several parameters such as frequency of occurrence of each node, number of loops at each node, number of occurrence of each arc uv , total number of arcs in D , the maximum length of a transaction, in-degree and out-degree of each node.

The algorithm first creates all nodes of D , one node for each item, with support count 0. Then each transaction is scanned and the directed path in D representing the transaction is constructed. If (i_1, i_2, \dots, i_k) is a transaction, the arc (i_j, i_{j+1}) is represented as a linked list. The header of this list has two fields. One field is used to store the list of vertices $(i_1, i_2, \dots, i_{j+1})$, which is called the label of the arc (i, j) and the other field is used to store the frequency of occurrence of the arc (i, j) . For example in the digraph given in Figure 3, the arc $(3, 5)$ occurs with values 2, $(1,3,5)$ and also with values 1 $(2,3,5)$, indicating that it occurs twice as part of $(1,3,5)$ and occurs once as a part of $(2,3,5)$. In general if an arc (i, j) has labels $k, (i_1, i_2, \dots, i_r, i, j)$ it means that the arc (i, j) appears k items as part of $(i_1, i_2, \dots, i_r, i, j)$. Dynamic memory allocation method is used for storing these values. The pseudo code for the construction of DGTDB is given in Figure 3. This algorithm also generates further informations which are given in the output.

Algorithm: Construction of DGTDB

Input: Transaction Database TDB, n : Number of distinct items in TDB, m : Number of Transactions.

Output: DGTDB: Directed Graph of TDB.

Method:

Create a node for each item and initialize the values $f(i)$, Out-Degree (i) , In-Degree (i) , In-Edge (i) , Out-Edge (i) $L_c(i)$ to zero.

Initialize $predecessor = \emptyset$; and $X = \emptyset$;

- 1: **for** each transaction T_j and for each each item i in T_j , **do**
- 2: $X = \bigcup\{i\}$; $f(i) + +$;
- 3: **if** $(|X| > K)$ **then**

```

4:    $K = |X|$ ;
5:   end if
6:   if ( $|T_j| = 1$ ) then
7:      $L_c(i) ++$ ;
8:     if ( $L_c(i) = 1$ ) then
9:       // If an edge from  $i$  to  $i$  does not exist so far
10:      CreateEdge( $i, i$ );
11:      Set Label ( $e_c, i$ ) =  $\langle X \rangle$ ;
12:       $f(e_c, i) = 1$ ;
13:       $n(E) ++$ ;
14:      Out-Degree( $i$ ) ++;
15:      In-Degree( $i$ ) ++;
16:      In-Edge( $i$ ) ++;
17:      Out-Edge( $i$ ) ++;
18:     end if
19:   end if
20:   if ( $L_c(i) > 1$ ) then
21:      $f(e_c, i) ++$ ;
22:     Out-Degree( $i$ ) ++;
23:     In-Degree( $i$ ) ++;
24:   end if
25:   if ( $predecessor \neq \emptyset$ ) and ( $X \notin$  any Label ( $e_c, i$ )) then
26:     CreateEdge ( $predecessor, i$ );
27:      $n(E) ++$ ;
28:     Set Label( $e_c, i$ ) =  $\langle X \rangle$ 
29:      $f(e_c, i) = 1$ ;
30:     Out-Degree( $predecessor$ ) ++;
31:     In-Degree( $i$ ) ++;
32:     In-Edge( $i$ ) ++;
33:     Out-Edge( $predecessor$ ) ++;
34:   end if
35:   if ( $X \in$  Label ( $e_j, i$ )) then
36:      $f(e_j) ++$ ;
37:     Out-Degree( $predecessor$ ) ++;
38:     In-Degree( $i$ ) ++;
39:   end if
40:    $predecessor = i$ ;
41: end for
42: return DGTDB;

```

Figure 2. Pseudo code for DGTDB construction

We illustrate the algorithm DGTDB with a transaction database consisting of 12 items and 30 transactions, which is given in Table 1.

TID	PRODUCTS IN EACH TRANSACTION	TID	PRODUCTS IN EACH TRANSACTION
T001	1, 3, 7, 8	T016	2, 4, 6, 8, 10, 12
T002	1, 2, 3, 8	T017	1, 3,5, 7
T003	2, 4, 5, 6	T018	9, 11
T004	2, 3, 5, 6	T019	10, 11,12
T005	1, 4, 5, 11	T020	6, 7, 8,12
T006	3, 4, 5, 12	T021	3
T007	1, 2, 3, 4	T022	10, 11, 12
T008	4, 5, 11	T023	2, 4,6, 8, 10,12
T009	1, 4, 5, 8	T024	1, 3, 5, 7, 9, 11
T010	3, 4, 10	T025	2, 4, 6, 8, 10, 12
T011	2, 3, 4	T026	3
T012	3	T027	2, 4, 6, 8, 10, 12
T013	3, 6, 9, 12	T028	3
T014	1, 10, 12	T029	12
T015	3, 6, 9, 12	T030	3

Table 1. Transaction Database, TDB

The DGTDB for the above transaction database is given in Figure 3.

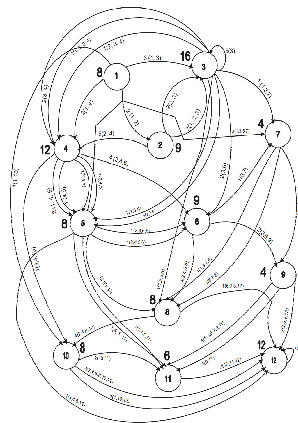


Figure 3. Directed Graph of TDB in Table 1.

The number given at each node represents the frequency of the corresponding item.

5 Algorithm for extracting frequent patterns (XoFP)

In this section we present an algorithm for extracting the set \mathcal{L} of all frequent patterns from DGTDB constructed in Section 4. For each node i , the algorithm generates all frequent patterns L_1 with i as the last item of L_1 . If in-edge $(i) = 0$, then $\{i\}$ is the only pattern with i as the last item. Otherwise for each edge e_c with i as head we consider all subsets X of label (e_c, i) such that $|X| \geq 2$ and $i \in X$. Then the frequency of X is the sum of the frequencies of all the edges e_c with i as head for which $X \subseteq label(e_c, i)$. If this frequency is greater than or equal to the minimum support threshold, then X is added to the set of frequent patterns. The algorithm XoFP is given in Figure 4.

Algorithm: XoFP, Extraction of Frequent Patterns from DGTDB.

Input: s : Minimum support threshold; DGTDB: Directed Graph of TDB.

Output: \mathcal{L} : The set of all frequent patterns mined from the DGTDB.

Method:

```

1:  $\mathcal{L} = \emptyset; m = 0;$ 
2: // Initialize set of frequent patterns
3: for each node  $i$  do
4:   if ( $f(i) \geq s$ ) then
5:      $\mathcal{L} = \mathcal{L} \cup \{i\}$ 
6:     return  $\{i\}, f(i)$ 
7:      $m++;$ 
8:   end if
9:   if ( $\text{In-Edge}(i) > 0$ ) then
10:     $f = 0;$ 
11:    for each Label  $(e_c, i), 1 \leq c \leq \text{In-Edge}(i)$  do
12:       $f = f(e_c, i);$ 
13:       $W = \text{Label}(e_c, i);$ 
14:      if ( $W \notin \mathcal{L}$ ) then
15:        for each  $x \subseteq W, |x| \geq 2$  and  $i \in x$  do
16:          for each Label  $(e_y, i)$  do
17:            if ( $x \subseteq \text{Label } e_y(i)$ ) then
18:               $f = f + f(e_y, i);$ 
19:            end if
20:          end for
21:        end for
22:        if ( $f \geq s$ ) then
23:           $\mathcal{L} = \mathcal{L} \cup \{x\};$ 
24:          return  $\{x\}, f$ 
25:         $m++;$ 
26:        end if
27:      end if
28:    end for
29:  end if
30: end for
31: return ( $\mathcal{L}$ )

```

Figure 4. XoFP, Extraction of patterns from DGTDB

Example 2. By applying XoFP to the DGTDB given in Figure 3, the set of frequent patterns obtained using the items 1,2,3 and 4 along with the respective frequencies are given in Table 2.

Node id	Frequent Patterns	Frequency
1	{1}	8
2	{2}	9
	{1, 2}	2
3	{3}	16
	{1, 3}	5
	{2, 3}	4
	{1, 2, 3}	2
4	{4}	12
	{1, 4}	3
	{2, 4}	7
	{3, 4}	4
	{2, 3, 4}	2

Table 2. Extracted frequent patterns from the nodes 1, 2, 3 and 4 of DGTDB where $s = 2$.

6 Algorithm for generating association rules (GEAR)

In this section we present an algorithm for generating association rules and strong association rules from the set of frequent patterns mined from the given TDB. An association rule which satisfies both minimum support threshold and minimum confidence threshold is called a strong association rule. For each frequent pattern X and for each nonempty proper subset Y of X the algorithm computes the support and confidence of the association rule $Y \Rightarrow X - Y$.

Algorithm: GEAR, Generating association rules from the frequent patterns

Input: \mathcal{L} : set of all frequent patterns; c - Minimum confidence threshold of rule; s - Minimum support threshold.

Output: R : set of all strong association rules; R' : set of all association rules not in R .

Method:

```

1:  $R = \emptyset; R' = \emptyset; N(R) = 0; N(R') = 0;$ 
2: for each  $X, Y$  where  $X \in \mathcal{L}$  with  $|X| > 1$  and  $Y \subseteq X, Y \neq \emptyset, Y \neq X$  do
3:   Generate Rule ( $Y \Rightarrow (X - Y)$ )
4:   Compute  $conf = \frac{f(X)}{f(Y)}$ 
5:   //  $f(X)$  is the frequency of pattern  $X$ .  $f(Y)$  is the frequency of pattern  $Y$ .
6:   if ( $conf \geq c$ ) then
7:      $R = R \cup \{Y \Rightarrow (X - Y) : support=s, confidence = conf\};$ 
8:      $n(R) ++;$ 
9:   else
10:     $R' = R' \cup \{Y \Rightarrow (X - Y) : support=s, confidence = conf\};$ 
11:     $n(R') ++;$ 
12:   end if
13: end for
14: return ( $R$ ) and return ( $R'$ );

```

Figure 5. GEAR, Algorithm for Generating Association Rules

Example 3. From Table 2 we have $X = \{1, 2, 3\}$ is a frequent pattern with frequency 2. The set of all association rules generated from this pattern by using GEAR along with the confidence and support for each rule is given in Table 3. We have taken the minimum confidence threshold c and the minimum support threshold s as 50 and 6 respectively.

S.No	Association Rules	Confidence of the Rule	Support of the Rule	R or R'
1	$\{1\} \Rightarrow \{2, 3\}$	25	6	R'
2	$\{2\} \Rightarrow \{1, 3\}$	22.2	6	R'
3	$\{3\} \Rightarrow \{1, 2\}$	12.5	6	R'
4	$\{1, 2\} \Rightarrow \{3\}$	100	6	R
5	$\{1, 3\} \Rightarrow \{2\}$	66.67	6	R
6	$\{2, 3\} \Rightarrow \{1\}$	50	6	R

Table 3. Association Rules mined from the frequent pattern $\{1, 2, 3\}$.

For the TDB given in Table 1 we have generated 188 frequent patterns. The number of frequent patterns generated with various support counts is given in Table 4. The total number

of associate rule generated is 1588. The number of association rule with various support counts is given in Table 5. The breakup of the number of association rules generated with various levels of confidence is given in Table 6.

S.No	Support %	No. of Patterns Generated	S.No	Support %	No. of Patterns Generated
1	3	188	10	30	3
2	6	113	11	33	3
3	9	84	12	36	3
4	12	72	13	39	1
5	15	25	14	42	1
6	18	14	15	45	1
7	21	12	16	48	1
8	24	9	17	51	0
9	27	5	18	54	0

Table 4. Number of patterns with various support counts

S.No	Support % \geq	No. of Association Rules Generated	S.No	Support % \geq	No. of Association Rules Generated
1	3	1588	10	30	0
2	6	746	11	33	0
3	9	484	12	36	0
4	12	292	13	39	0
5	15	38	14	42	0
6	18	16	15	45	0
7	21	10	16	48	0
8	24	2	17	51	0
9	27	0	18	54	0

Table 5. Number of association rules with various support counts

S.No	Confidence %	No. of Association Rules Generated
1	< 50	465
2	50-59	162
3	60-69	81
4	70-79	36
5	80-89	95
6	≥ 90	749

Table 6. Number of strong association rules with various confidence levels

7 Conclusion

In this paper we have proposed a new data structure consisting of a directed graph D and a path system in D for representing a TDB. We have presented algorithms for constructing D , for generating frequent patterns using D and for generating association rules. During the entire process the data is scanned exactly once. Further it is possible to get several information about the TDB by using graph theoretic parameters. For example if $\text{in-degree}(i) = 0$ in the directed

graph D_1 obtained from D by removing all the loops, then the item i always appears as the first item in every transaction T with $i \in T$. Similarly if $\text{out-degree}(i) = 0$, then the item i always appears as the last item in every transaction T with $i \in T$. Our algorithm can be used to identify all such items. Use of other graph theoretic parameters to extract new knowledge about the TDB and the comparison of the performance of this algorithm with other existing algorithms in the literature using real data set will be reported in a subsequent paper.

Acknowledgment

The first author is thankful to the Department of Science and Technology, New Delhi for its support through the n-CARDMATH Project No. SR/S4/MS:427/07.

Bibliography

- [1] R. Agrawal, T. Imielinski and A. Swami, Mining association rules between sets of items in large database, In *Proc. of the ACM SIGMOD International Conference on Management of Data (ACM SIGMOD 93)*, Washington, USA, 22(2)207-216, May 1993.
- [2] G. Chartrand and L. Lesniak, *Graphs and Digraphs*, Chapman and Hall, CRC, 4th edition, 2005.
- [3] R. Agrawal and R. Srikant, Fast Algorithms for mining association rules, In *Proc. of the 20th International Conference on Very Large Database (VLDB' 94)*, Santiago, Chile, 487-499, June 1994.
- [4] J. Han, J. Pei and Y. Yin, Mining Frequent Patterns without Candidate Generation, In *Proc. of the 2000 ACM SIGMOD International Conference on Management of Data*, Dallas, Texas, USA, 29(2):1-12, May 2000.
- [5] S. Brin, R. Motwani, J.D. Ullman and S. Tsur, Dynamic itemset counting and implications rules for market basket data, In *Proc. of the ACM SIGMOD International Conference on Management Data*, 26(2):255-264, 1997.
- [6] M. Hontsma and A. Swami, *Set oriented mining for association rules in relatrend database*, The technical report RJ9567, IBM Almaden Research Centre, San Jose, California, October 1993.
- [7] J. Han and M. Kamber, *Data mining, Concepts and Applications*, Elsevier Inc., (2006).

Adaptive Network Coding Scheme for TCP over Wireless Sensor Networks

Y.-C. Chan, Y.-Y. Hu

Yi-Cheng Chan*, Ya-Yi Hu

Department of Computer Science and Information Engineering,
National Changhua University of Education
No.2, Shi-Da Road, Changhua City 500, Taiwan
ycchan@cc.ncue.edu.tw, m9954016@mail.ncue.edu.tw

*Corresponding author: ycchan@cc.ncue.edu.tw

Abstract: The purpose of this paper is to develop a network coding scheme to enhance TCP performance in wireless sensor networks. It is well known that TCP performs poorly over wireless links which suffer from packet losses mainly due to the bad channel. To address this problem, it is useful to incorporate network coding into TCP, as network coding can offer significant benefits in terms of throughput, reliability, and robustness. However, the encoding and decoding operations of network coding techniques will bring an additional delay that has a negative effect on applications of wireless sensor networks. In this paper, we propose an adaptive network coding (ANC) scheme which contains two major aspects: the adjustment of the redundancy factor R and the adjustment of the coding window size CW . We dynamically adjust these two parameters depending on the measured packet loss rate, so that the proposed ANC can effectively mask packet losses and reduce the decoding delay of network coding. The performance of our scheme is evaluated by simulations using NS-2 simulator. Compared to other schemes, the ANC not only achieves a good throughput but also has the lowest average delay and the lowest maximum delay in all experimental environments.

Keywords: network coding, TCP, delay, wireless sensor networks.

1 Introduction

The Transmission Control Protocol (TCP) is the main transport protocol that provides reliable transmission in the current Internet. It has been developed for many years. Most applications on the Internet depend on TCP to ensure safe delivery of data, such as FTP (File Transfer Protocol), HTTP (HyperText Transfer Protocol), SMTP (Simple Mail Transfer Protocol), and so on. TCP performs well in wired networks where packets losses mainly occur due to congestion. However, the performance of TCP degrades very fast in wireless networks. In a wireless environment, there is not only congestion but also numerous other reasons for packet losses exist. Traditional TCP treat a packet loss event as the indication for network congestion, and then decrease its congestion window size. This may severely impair the throughput when the TCP runs on a wireless channel. Therefore, it is very important to improve the performance of TCP in wireless networks. And this goal can be achieved by combining network coding with TCP.

Network coding is a new transmission paradigm originally proposed by Ahlswede et al [1]. In recent years, it has received much attention and has been generated huge research in communication networks [2]. TCP/NC [3] is the first one that incorporates network coding into TCP with minor changes to the protocol stack. They present a solution which embeds the network coding operation in a separate layer below transport layer and above network layer on the source and receiver side. The idea of network coding is that, instead of transmitting individual packets, the sender takes several packets and combines them together for transmission. Consequently, successful reception of information does not depend on receiving specific packet content but rather

on receiving a sufficient number of combinations [4]. Then it can compensate for the packets in the presence of random losses, as long as the sent redundant combinations are enough. This characteristic is very attractive for TCP to improve the robustness and effectiveness of data transmission over lossy wireless networks.

Wireless sensor networks (WSNs) are usually composed of a large number of radio-equipped sensor devices to cooperatively monitor physical or environmental conditions, such as temperature, sound, vibration, pressure, motion, or pollutants. These devices typically include some strong constraints in terms of energy, memory, computational speed and communication bandwidth. In fact, most applications on wireless sensor networks prefer faster and reliable packet delivery to higher throughput [5]. Generally, these applications are not only loss-sensitive that require successful transmission of all packets or at a certain success ratio but also delay-sensitive that require timely delivery of data [6]. For these applications, packet loss will lead to retransmission and the inevitable consumption of additional battery power. In addition, if the delay time is too long, the monitored data may be outmoded. To overcome these problems it is beneficial to use network coding by transmitting redundant packets to mask packet losses. Thus, the number of retransmissions and timeouts can be reduced.

In this paper, we propose a new network coding scheme to improve TCP performance over wireless sensor networks which is called adaptive network coding (ANC). The concept of ANC is divided into two parts. In the first part, we adjust the redundancy parameter dynamically according to the network situation. As a result, the event of packet loss can be masked from the congestion control algorithm of TCP by sending enough redundant combinations. In the second part, we change the coding window size contingent on the measured packet loss rate to obtain the optimal throughput-delay trade-off. The results of simulation show that our scheme achieves higher throughput and lower delay in difference network scenarios.

The remainder of this paper is organized as follows. In Section 2 we provide an overview of related work. Section 3 we propose our schemes and algorithms focus on the redundancy parameter and the coding window size of network coding. In Section 4 we present an experimental evaluation of the performance, and finally, our conclusions as well as the future work are discussed in Section 5.

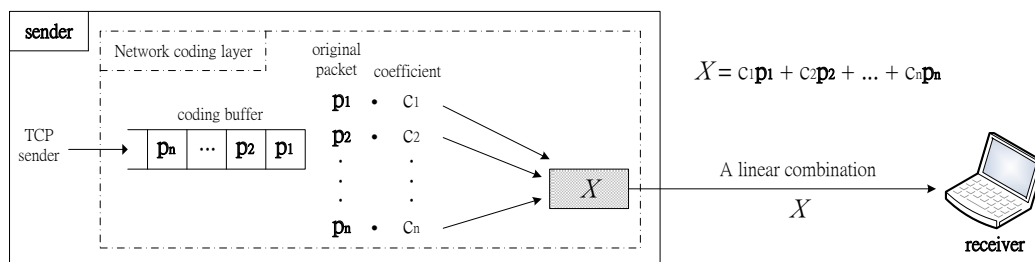


Figure 1: An example of random linear network coding

2 Related Work

Random linear network coding is one of the major techniques in the network coding. The encoding coefficients are randomly chosen from a set of coefficients of a finite field. A linear equation of packets is then performed to generate a coded packet, and the receiver only needs to receive a sufficient number of linear equations in the forms of the coded packets to successfully decode the original packets. For instance (see Figure 1), suppose that the sender buffers the first

n packets p_1, \dots, p_n in its coding buffer, then the sender chooses n coefficients c_1, \dots, c_n from a finite field and mates n coefficients with n packets. Next, the sender encodes these pairs into a linear combination X , $X = c_1p_1 + c_2p_2 + \dots + c_np_n$, then sends it to the receiver. If the receiver receives enough number of these combinations then it can decode the original packets. In recent years, there have been many studies make use of random linear network coding to improve TCP performance.

The main aim of TCP/NC [3] is to mask losses from TCP using random linear coding. The source transmits random linear combinations of packets currently in the congestion window, and the receiver acknowledges every innovative linear combination it receives, even if it cannot decode an original packet immediately. This scheme gives a new interpretation of ACKs, and brings a new concept that “seen packet” which is defined in [7] as an abstraction for the case in which a packet cannot yet be decoded but can be safely removed from the coding buffer at the sender. TCP-Vegas is chosen for the transport layer protocol in TCP/NC, as it is more compatible with their modifications. When losses are effectively masked, TCP-Vegas can infer congestion from increased *RTTs*. Additionally, TCP/NC uses a constant redundancy parameter R to compensate for the loss rate of the channel, for every packet arrives from TCP, R linear combinations are sent to the IP layer on average. However, in a wireless environment, the packet loss rate is very likely to not a constant value. A fixed redundancy parameter R may damage network performance, since it is not always suitable for all network conditions. If the value of parameter R is too small, then the losses are not effectively masked from the TCP layer and will result in timeout. On the other hand, if the value of parameter R is too large, the source may send too many linear combinations that consume network resources. Thus, the value of R should be dynamically adjusted depending on the estimated packet loss rate.

The feedback based network coding (FNC) [8] uses the implicit information behind of the seen scheme to find out the exact number of packets needed by receiver to decode all data. The receiver computes the *DIFF* value by two variables: the number of seen packets and the largest packet index in the coefficient matrix, and then embeds this value into the ACK header. When the sender receives an ACK, it uses this value to decide how many random linear combinations should be retransmitted and how many original packets should be combined in a linear combination. In this way, the FNC retransmission scheme can reduce the decoding delay and the number of redundant retransmissions. But the FNC retransmission scheme highly relies on feedback, this will greatly impair the throughput when the network with a large round-trip time. Furthermore, if the ACK is lost, the sender cannot repair the loss of linear combinations through instant retransmissions.

The SANC-TCP protocol [9] is designed primarily to optimize the TCP/NC protocol. The redundancy factor R of TCP/NC is constant, while SANC-TCP adjusts the redundancy factor R adaptively based on existing network conditions. In order to implement this approach, SANC-TCP adds some information in the ACK header to indicate the current network state, thus enable the sender to dynamically change the value of R according to the real system. Our scheme is similar to the SANC-TCP that with dynamic adjustment of redundancy factor R , but in different ways. Moreover, we limit the coding window size and dynamically change it depending on the packet loss rate to reduce the delay of network coding.

3 Proposed Method

In this section, we describe our adaptive network coding scheme (ANC) which consists of two parts. First, we adjust the redundancy factor R dynamically by estimating the packet loss rate in current network, so that the value of R can represent the actual network state. Second, in order to reduce the decoding delay of network coding, we limit the coding window size and

dynamically adjust it contingent on the measured packet loss rate.

3.1 Adjustment of the redundancy factor R

Before adjusting the redundancy factor R , we have to find the packet loss rate of current network. In our scheme, we insert two variables into the header of ACK, which are $send_count$ and $seen_count$. The variable $send_count$ refers to the number of packets that have been sent from the sender and the variable $seen_count$ refers to the number of seen packets at the receiver. Then, the sender can use these and other related variables to calculate packet loss rate. For example, when the ACK arrives from receiver, the sender retrieves the variables $send_count$ and $seen_count$ from ACK header and compares $seen_count$ with a threshold T which is the variable used to determine if it is time to adjust the value of R . If $seen_count$ is more than or equal to T , the sender starts to calculate the packet loss rate by the following equations:

$$diff_send = send_count - send_old, \quad (1)$$

$$diff_seen = seen_count - seen_old, \quad (2)$$

$$loss_rate = \frac{diff_send - diff_seen}{diff_send}, \quad (3)$$

where $send_old$ and $seen_old$ are the previous value of $send_count$ and $seen_count$ respectively. The initial $send_old$ and $seen_old$ is set to 0. The $diff_send$ refers to the number of packets that have been sent during this calculation cycle while $diff_seen$ refers to the number of seen packets during this calculation cycle. Therefore, the $loss_rate$ means the packet loss rate of current network. Through the above equations, we can obtain the accurate packet loss rate to conduct the correct adjustment of redundancy factor R .

After calculating the packet loss rate, the sender adjusts the redundancy factor R accordingly. If current T is equal to $init$ means this is the first time for the sender to adjust the value of R . In the simulation, we choose $init = 20$ which is the same as initial coding window size. The equation below shows the calculation of the redundancy factor R for the first time:

$$R_{current} = \frac{1}{1 - loss_rate}, \quad (4)$$

where $R_{current}$ is current redundancy factor R . On the other hand, if it is time to adjust the value of R after the first calculation, the value of R can be obtained as follow:

$$R_{current} = a * \frac{1}{1 - loss_rate} + b * R_{old}, 0 \leq a, b \leq 1, a + b = 1. \quad (5)$$

The above equation is calculated using moving average to take into account the previous trends on R_{old} , where R_{old} is the last value of redundancy factor R . The values of variables a and b can be set based on the network situations. If the random loss rate of the network environment tends to change frequently, the value of a should be set greater than b . In contrast, if the random loss rate of the network environment tends to change occasionally, the value of a should be set less than b .

3.2 Adjustment of the coding window size

The coding window size represents that the largest number of original packets can be encoded in a linear combination. To reduce the decoding delay of network coding, we limit the coding window size and adjust it according to the packet loss rate.

After the sender computed the redundancy factor R , the next process for the sender is to adjust the coding window size. The following equation describes the adjustment of coding window size for the first time:

$$CW = 20 + \lfloor loss_rate * 10 \rfloor, \quad (6)$$

where CW refers to the coding window size. In our scheme, the coding window size is dynamically adapted depending on the packet loss rate. This means that the coding window size increases as the packet loss rate increases, and the coding window size decreases as the packet loss rate decreases. After the first calculation, the value of coding window size can be obtained as follow:

$$CW = 20 - \left\lfloor \frac{1 - R_{current}}{R_{current}} * 10 \right\rfloor. \quad (7)$$

The above equation shows that the packet loss rate estimation for adjusting the coding window is derived from the value of $R_{current}$ rather than the measured packet loss rate at this period. That is because the value of $R_{current}$ is used to predict the packet loss rate in the next period.

Finally, we must update several related variables, such as R_{old} to $R_{current}$, $send_old$ to $send_count$ and $seen_old$ to $seen_count$. The value of T also should be reset by the following equation:

$$T = CW + seen_count. \quad (8)$$

Table 1: Coding window size adjustment in difference loss rates

CW	Loss rate 0%		Loss rate 10%		Loss rate 20%		Loss rate 30%		Loss rate 40%	
10	585.1	0.06828	617.2	0.07108	515.7	0.07460	346.3	0.07648	91.1	0.07476
15	877.6	0.06828	849.3	0.08667	733.7	0.09681	602.5	0.10356	341.6	0.10394
20	998.6	0.07997	898.6	0.11153	793.3	0.13131	681.6	0.14362	509.7	0.15326
25	998.6	0.09997	899.9	0.14440	798.9	0.17561	695.2	0.19556	589.3	0.20850
30	998.6	0.11993	899.9	0.17740	799.0	0.21640	696.1	0.25118	592.8	0.27962
35	998.6	0.13986	899.9	0.21272	799.0	0.21640	696.1	0.31536	592.8	0.27962
40	998.6	0.15976	899.9	0.24271	799.0	0.21640	696.1	0.37298	592.8	0.27962
ANC	998.6	0.07997	899.9	0.13830	799.0	0.18977	696.1	0.23014	592.4	0.26464
CW	20		24		27		28		29	

3.3 Analysis of the coding window size

The initial value of coding window size in our scheme is acquired by the experimental results presented in Table 1. For this experiment, the topology and experimental parameters are the same as that in the subsection 4.1.1. Table 1 shows the data of throughput (TP) and delay which are obtained by adjust the coding window size in difference loss rates. The data that have an outline border is the optimal trade-off value when compared with others data in the same loss rate. The last two rows in Table 1 are the experimental results and the coding window adjustments of our scheme ANC. For the initial setting, we assume the packet loss rate at beginning of network is zero, so we set the initial coding window size to 20 that is the optimal value in Table 1 when the loss rate is 0%. After this, the sender adjusts the coding window size in accordance with the measured packet loss rate. This experimental result (see Table 1) indicates that the throughput

of our scheme at last row can reach the optimal value in most cases while has a shorter delay, that is, the adjustment of the coding window size in our scheme is suitable.

As mentioned before, in order to reduce the decoding delay of network coding, we limit the coding window size in our scheme. When the coding window size is limited, the coding coefficient matrix can become smaller and the complexity of decoding can become lower. Hence, the time of packets remain in the decoding buffer is decreased, that is, packets can be discarded earlier from the decoding buffer. For these reasons, our scheme can effectively reduce the decoding delay of network coding. Nevertheless, a smaller coding window size does not mean to be better. If the size of coding window is too small, the coding window will restrict the use of available network resources which could cause a poor throughput. As a consequence, we must take into account the trade-off between throughput and delay to get higher throughput and lower delay.

The algorithm of our scheme is specified using pseudo-code that shown in Tables 2 and 3.

Table 2: The operations of network coding layer at the sender

Event	Pseudo-code
Initialization:	1)Set $NUM, send_count, send_old$ and $seen_old$ to 0. 2)Set $init$ to 20. Let $T = init$.
When the packets arrive from TCP sender:	1)If the packet is a control packet used for connection management, deliver it to the IP layer then doing nothing; else, move to state 2). 2)If the packet not already in the coding buffer, add it to the coding buffer. 3)Set $NUM = NUM + R$. ($R =$ redundancy factor) 4)Repeat the following $\lfloor NUM \rfloor$ times: 1.Generate a random linear combination of the packets in the coding window. 2.Count the $send_count = send_count + 1$. 3.Add the network coding header specifying the set of packets in the coding window and the coefficients used for the random linear combination. 4.Add the variable $send_count$ to the network coding header. 5.Deliver the packet to the IP layer. 5)Set $NUM =$ fractional part of NUM .
When the ACK arrives from receiver:	1)Pick up the variables $seen_count$ and $send_count$ from ACK header. 2)If $seen_count \geq T$ start to adjust the values of R and CW ; else, move to state 3). 1.Compute the $diff_send = send_count - send_old$ and the $diff_seen = seen_count - seen_old$. 2.Compute the $loss_rate = (diff_send - diff_seen)/diff_send$. 3.If $T = init$: a) $R_{current} = 1/(1 - loss_rate)$. b) $CW = 20 + \lfloor loss_rate * 10 \rfloor$. c)Move to state 5. 4.If $T > init$: a) $R_{current} = a * 1/(1 - loss_rate) + b * R_{old}$. b) $CW = 20 - \lfloor (R_{current} - 1)/R_{current} * 10 \rfloor$. 5.Update R_{old} to $R_{current}$, $send_old$ to $send_count$ and $seen_old$ to $seen_count$. 6.Reset $T = CW + seen_count$. 3)Remove the ACKed packet from the coding buffer and hand over the ACK to the TCP sender.

4 Performance Evaluation

Our adaptive network coding scheme is evaluated by means of the Network Simulator (ns-2) under various conditions such as different link bandwidths, packet sizes, and random loss rates. Evaluation metrics in performance test are the throughput, the average delay, and the maximum

Table 3: The operations of network coding layer at the receiver

Event	Pseudo-code
Initialization:	1)Set <i>seen_count</i> to 0.
When a packet arrives from sender:	1)Count the $seen_count = seen_count + 1$. 2)Remove the network coding header, then retrieve the coding vector and the variable <i>send_count</i> . 3)Add the coding vector as a new row to the existing coding coefficient matrix, and perform Gauss-Jordan elimination to update the set of seen packets. 4)Add the payload to the decoding buffer. Perform the operations corresponding to the Gauss-Jordan elimination, on the buffer contents. If any packet gets decoded in the process, deliver it to the TCP sink and remove it from the buffer. 5)Generate a new ACK with sequence number equals to that of the oldest unseen packets and add two variables <i>send_count</i> and <i>seen_count</i> to the ACK header.
When an ACK arrives from the TCP sink:	1)If the ACK is a control packet for connection management, deliver it to the IP layer; else, ignore the ACK.

delay. We compare our scheme ANC with SANC-TCP and TCP-Vegas under both fixed and unfixed loss rate. We use TCP-Vegas as the transport layer protocol, and set the parameters of TCP-Vegas to $\alpha = 28$, $\beta = 30$, $\gamma = 2$. This setting is identical with TCP/NC [3]. The values of a and b in our scheme are set to $a = 0.2$ and $b = 0.8$ individually.

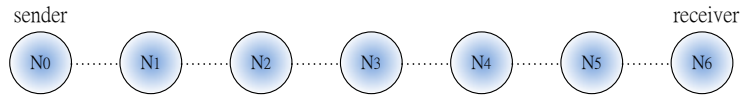


Figure 2: A tandem network consisting of 6 hops

4.1 Results for fixed random loss rates

The topology in this simulation is a tandem network consisting of 6 hops, as shown in Figure 2. The sender and the receiver are at opposite sides of the chain.

In contrast to traditional wireless networks, wireless sensor networks generally have a lower bandwidth for data transmission. We experiment with two different bandwidth settings: (a) bandwidth = 1 Mbps and (b) bandwidth = 250 kbps which is the channel bandwidth of IEEE 802.15.4. In the following experiments we assume that the source always has data to send.

Set the link bandwidth to 1 Mbps

In this experiment, each link has a bandwidth of 1 Mbps, and a propagation delay of 10 ms. The buffer size on the links is set to 200 packets. The packet size is 1000 bytes, and the random loss rate is varied from 0% to 20% on each link. The simulation time is 1000 seconds.

The throughput obtained corresponding to different random loss rates is plotted in Figure 3. Our scheme ANC and SANC achieve a similar throughput, but the throughput of Vegas decrease rapidly when the random loss rate increases. The adjustment of the redundancy factor R that we proposed in subsection 3.1 makes the sender to send out the proper number of redundant linear combinations so that the random loss rate can be masked effectively. As a result, we can get a fairly well throughput.

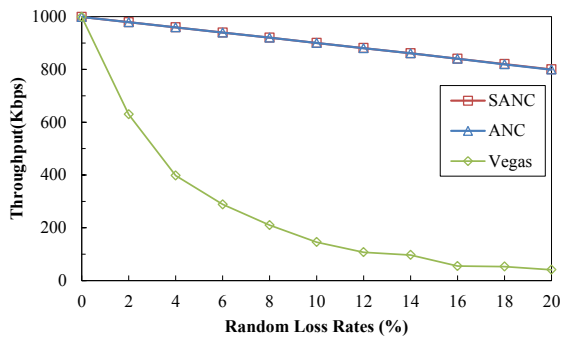


Figure 3: Throughput

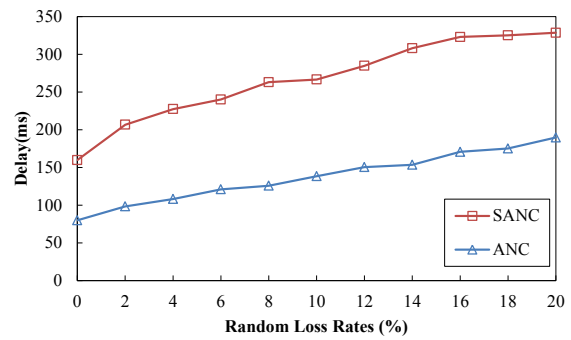


Figure 4: Average delay

Figure 4 shows the average delay of ANC and SANC. Compared with the SANC, our scheme have lower average delay in every case, moreover, the average delay of our scheme is almost half of the average delay of SANC. This is because we limit the coding window size which would decrease the complexity of decoding and reduce the time that packets remain in the decoding buffer. Thus, we can have a lower delay time than that of SANC.

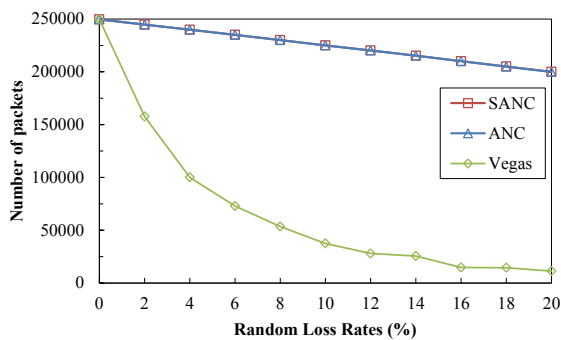


Figure 5: The total number of packets

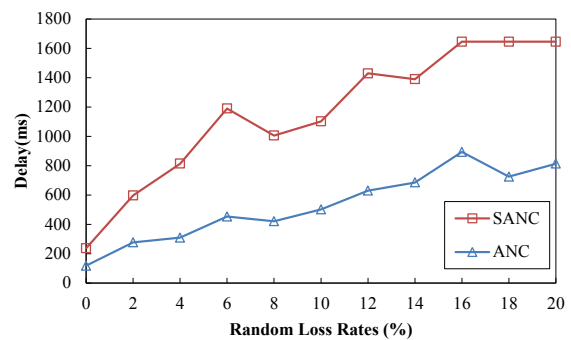


Figure 6: Maximum delay

Figure 5 shows the total number of packets that have been send for SANC, ANC, and Vegas. This figure is essentially identical to Figure 3 except the Y-axis represents the total number of sent packets. As shown in Figure 5, the metric of total number of sent packets for ANC and SANC outperform Vegas as the loss rate increases.

With regard to the simulation of delay, we also retrieve the maximum delay at various loss rates. As illustrated in Figure 6, the maximum delay of SANC rises significantly with loss rate, while our scheme rise slightly. Furthermore, the maximum delay of SANC is at least two times more than ANC in most cases.

Set the link bandwidth to 250 kbps

The channel bandwidth of IEEE 802.15.4 is 250 kbps [10]. Many WSNs adopt IEEE 802.15.4 for communicating among nodes. The characteristics of IEEE 802.15.4 technology includes low rate, low transit distance, low power, low cost, simple architecture, small size, etc. All these characteristics are applied to the applications of WSN.

The experimental parameters in this simulation are identical to the setting used in subsection 4.1.1, except that the bandwidth of each link is 250 kbps, and the buffer size on the links is set to 50 packets. Besides, we simulate the effect of two different packet sizes on the performance. The packet sizes are set to 1000 bytes and 200 bytes.

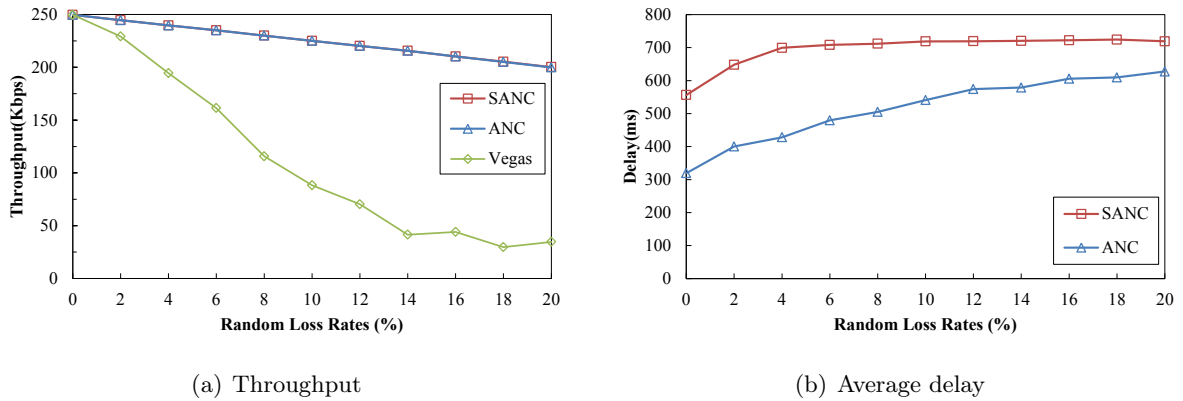


Figure 7: The throughput and delay when packet size is 1000 bytes

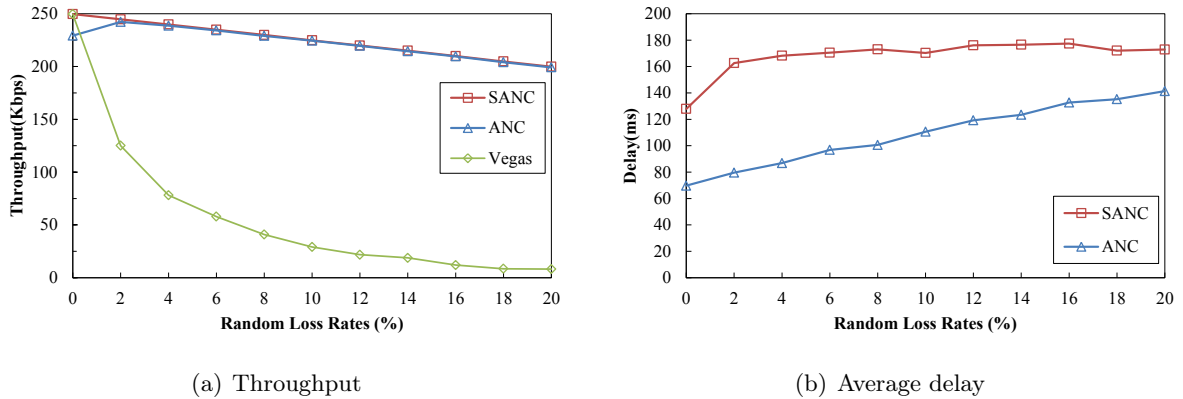


Figure 8: The throughput and delay when packet size is 200 bytes

As the link capacity and the buffer size become smaller, our scheme ANC still has high throughput and low delay, as shown in Figure 7. However, it must take a longer delay to transmit a large packet (1000 bytes) if the link capacity becomes smaller. This is demonstrated through the simulation results, when the loss rate is 20%, ANC's average delay is 189.77 ms in Figure 4 where the bandwidth is 1 Mbps, while ANC's average delay is 627.58 ms in Figure 7(b) where the bandwidth is 250 kbps.

Then, we change the packet size to 200 bytes. Figure 8(a) demonstrates that when the loss rate is 0%, the throughput of our scheme ANC is lower than SANC and Vegas. This is because that the coding window size we adjust according to subsection 3.2 is not big enough in such a network environment, thus the available network resources is restricted by coding window and result in lower throughput. That is to say, when the packet size is very small, we must have a sufficient coding window size to take full advantage of network resources. In Figure 8(b), the average delay of our scheme is lower than SANC in every case. In addition, the average delay in Figure 8(b) is lower when compared with the Figure 7(b) because the smaller packet size has the lower transmission time.

4.2 Results for unfixed random loss rates

All previous simulations focus on the behavior of ANC under the fixed loss rate. Now, we evaluate its performance in an unknown environment with unfixed loss rate. We use the same topology and parameters as 4.1.1. The background loss rate in this scenario is 2%. The loss rate

is changed to 6% from 200 to 400 second, and changed to 10% from 600 to 800 second. Total simulation time is 1000 seconds.

In Figure 9, the X-axis represents the simulation time, and the Y-axis represents the average throughput. It clearly shows that ANC and SANC can quickly and effectively handle the sudden bursty losses by adjusting the redundancy factor R dynamically, but Vegas is seriously affected by bursty losses. In Figure 10, we show the comparison of delay between ANC and SANC. The delay of SANC increases rapidly when the loss rate is suddenly changed at 200 second and 600 second. In contrast to SANC, ANC can still maintain quite low latency under the changed loss rate. This is because our scheme ANC can constantly measure the packet loss rate of network and further adjust the coding window size to get low delay. Based on the above experimental results, we can demonstrate that even under the unknown environment with bursty losses, our scheme ANC still can send enough linear combinations to compensate the lost packets. Thus, ANC can reach high throughput and keep low delay through the adjustment of R and CW .

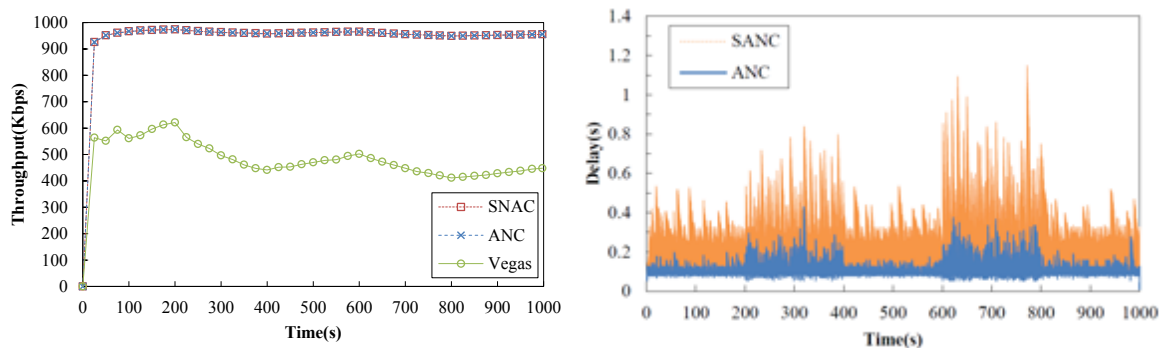


Figure 9: The average throughput under bursty losses Figure 10: The delay evolution under bursty losses

4.3 Results in a more realistic environment

In this subsection, we compare our scheme with SANC-TCP and TCP-Vegas in a wireless topology shown in Figure 11. There are five nodes (four hops) in this topology. The distance between each node is 30 meters. Table 4 is a parameter table which shows the parameter settings that used in the performance test. The simulation time is 100 seconds. The FTP flow starts at 1 second. In this experimental environment, some packets can be lost for reasons other than congestion.

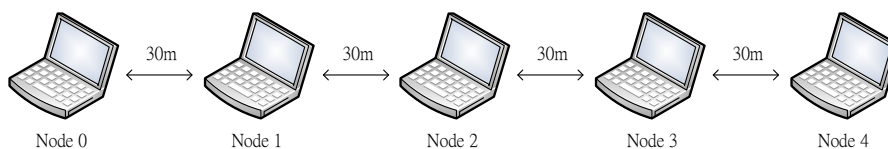


Figure 11: Wireless topology

The experimental results show that the throughput of ANC and SANC are 272.0 kbps and 271.9 kbps respectively. Both of them outperform Vegas whose throughput is 260.5 kbps. The comparison of RTT between ANC and SANC is shown in Figure 12. At the beginning of this experiment, we can see the RTT for both schemes are rising rapidly, then the RTT of ANC is maintained between 9 to 14 milliseconds that is lower than the RTT of SANC which is maintained between 15 to 20 milliseconds. Figure 13 presents the average RTT and maximum

Table 4: Parameter table

Parameter	Value
Packet Size	200 bytes
Buffer Size	50 packets
Buffer Management Scheme	DropTail
Link Bandwidth	2 Mbps
Transmission Range	40 meters
Carrier Sensing Range	90 meters
MAC Protocol	CSMA/CA
Routing Algorithm	DSR

RTT of ANC and SANC. It is clear ANC has lower average RTT and maximum RTT than SANC.

As mentioned before, this is because we limit the coding window size, and then the complexity of decoding can be decreased. Thus, packets will be discarded earlier from the decoding buffer. According to the experimental results, our scheme ANC can also have a good performance in such an environment.

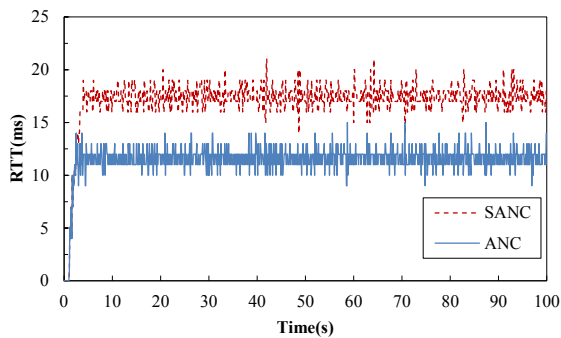


Figure 12: RTT of ANC and SANC

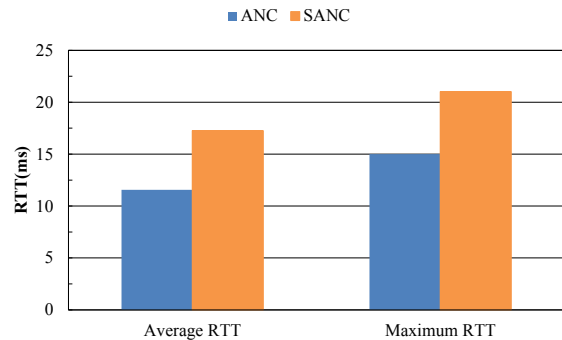


Figure 13: Average RTT and Maximum RTT of ANC and SANC

5 Conclusions

In this paper, we introduced a network coding scheme called adaptive network coding (ANC) that can be applied on wireless sensor networks. The objectives of our study are to mask packet losses by transmitting redundant linear combinations and to reduce the decoding delay of network coding by limiting the coding window size. We compare our scheme with SANC-TCP and TCP-Vegas in difference random loss rates, different link bandwidths, and difference packet sizes. From the experimental results obtained, our scheme has a lower delay than SANC-TCP in all experimental environments without sacrificing throughput. Besides, both the throughput of our scheme and SANC-TCP are significantly better than that of TCP-Vegas in most cases. It is important to note that the adjustment of coding window size can be investigated in more detail so as to improve the throughput when TCP runs on a harsh network environment.

A possible direction for further study is to allow intermediate nodes to perform the encoding and decoding whenever they receive packets. That means encoding and decoding operations are done in a hop-by-hop manner. In this way, the network coding can be used for a wide variety of topologies of wireless sensor networks.

Bibliography

- [1] R. Ahlswede, N. Cai, S.-Y. R. Li, and R. W. Yeung, Network information flow. *IEEE Trans. on Information Theory*, 46(4):1204-1216, Jul. 2000.
- [2] D. Silva and F. R. Kschischan, Universal Secure Network Coding via Rank-Metric Codes, *IEEE Trans. on Information Theory*, 52(2): 1124-1135, Feb. 2011.
- [3] J. K. Sundararajan, D. Shah, M. Medard, M. Mitzenmacher, and J. Barros, Network Coding Meets TCP, *2009 Proceedings of IEEE INFOCOM*, 280-288, Apr. 2009.
- [4] C. Fragouli, J.-Y. Le Boudec, and J. Widmer, Network coding: An instant primer, *ACM SIGCOMM Computer Communication Review*, 36(1): 63-68, Jan. 2006.
- [5] Yao-Nan Lien, Hop-by-Hop TCP for Sensor Networks, *International Journal of Computer Networks & Communications*, 1(1):1-16, Apr. 2009.
- [6] C. Wang, K. Sohraby, B. Li, M. Daneshmand, and Y. Hu, A survey of transport protocols for wireless sensor networks. *IEEE Network Magazine*, 20(3): 34-40, Jun. 2006.
- [7] J. K. Sundararajan, D. Shah, and M. Medard, ARQ for network coding, in *Proc. of IEEE International Symposium on Info. Theory (ISIT)*, 1651-1655, Jul. 2008.
- [8] J. Chan, L. Liu, X.Hu, and W. Tan, Effective retransmission in network coding for TCP. *Int J Comput Commun*, ISSN 1841-9836, 6(1):53-62, 2011.
- [9] S. Song, H. Li, K. Pan, J. Liu, and Shuo-Yen Robert Li, Self-adaptive TCP protocol combined with network coding scheme, In *International Conference on Systems and Networks Communications (ICSNC)*, 20-25, Oct. 2011.
- [10] F. Xia, A. Vinel, R. Gao, L. Wang, and T. Qiu, Evaluating ieee 802.15.4 for cyber-physical systems, *EURASIP J. Wireless Commun. and Networking*, 1-15, Feb. 2011.

Operation Mechanism of the Driving Force System of Ecosystem of Cyber-society Based on the System Dynamics

X. Guan, Z. Zhang, S. Zhang

Xiaolan Guan

School of Economics and Management
Beijing Institute of Graphic Communication
1 Xinghua Avenue (Band Two), Daxing, Beijing 102600, China
*Corresponding author: 08113101@bjtu.edu.cn

Zhenji Zhang, Shugang Zhang

School of Economics and Management
Beijing Jiaotong University
No. 3 Shang Yuan Cun, Hai Dian District, Beijing 100044, China
zhjzhang@bjtu.edu.cn, 10113132@bjtu.edu.cn

Abstract: Operation of the driving force system of Ecosystem of Cyber-society needs a scientific mechanism of intervention and regulation to solve the integration problem of a variety of organizations and forces within the Ecosystem of Cyber-society, shorten the process from uncoordination to coordination, promote the orderly operation of the driving force system of Ecosystem of Cyber-society, make the system play a strong force, in order to promote the formation and rapid development of Ecosystem of Cyber-society. We analyze the driving force system of Ecosystem of Cyber-society using the theory of System Dynamics and propose a theoretical framework, and then present its operation mechanism systematically.

Keywords: Ecosystem of Cyber-society, Driving Force System, Operation Mechanism, Life Cycle, System Dynamics.

1 Introduction

The operation of the driving force system needs certain operation mechanism, which refers to create interaction relationships between the driving forces, establish a fulcrum, platform and regulation for the driving forces, so that the driving force system can rely on the environment of the driving force system, and make the driving forces function effectively. Similarly, the operation of the driving force system of Ecosystem of Cyber-society also needs a scientific mechanism of intervention and regulation to solve the integration problem of a variety of organizations and forces within the Ecosystem of Cyber-society, shorten the process from uncoordination to coordination, promote the orderly operation of the driving force system of Ecosystem of Cyber-society, make the system play a strong force, in order to promote the formation and rapid development of Ecosystem of Cyber-society. However, since the driving force system consists of a lot of driving forces, its hierarchy is very complex with multi-levels, and different structures may emerge unexpected functions, so we need to find out an effective way to analyze it.

Currently, among the literature can be found around the world, although some scholars have put forward the concept of ecological networks from different perspectives and gradually realized the importance of this research area, there is still not much comprehensive and systematic research about the cyber-society from the perspective of ecosystem. And there is no literature that has proposed the concept of driving force system of Ecosystem of Cyber-society and studied its mechanism systemically from the perspective of the ecosystem for the development of cyber-society. It is a new research area.

The related research mainly focuses on the sociology of network and social ecology. And among them, different scholars have different views about the concept of sociology of network.

About the research on the sociology of network, some scholars advocated classifying it into the area of information science, while some other scholars emphasized the unique features of the network society in the media and dissemination, and claimed to study the evolution of the network society from the perspective of media and dissemination, and one of the representatives is an American social scientist Manuel Custer. He thought that the network society is a new social form of society that different from the reality, and pointed out the network constructs the new form of our society, and the spread of network logic has changed the operating result of production, experience, rights and culture. And it is characterized by social forms than the superiority of social action [1]; While Gong Qi considered the network society as a virtual space of human and humans life, study and work in this alternative space [2].

And the research on the social ecology started earlier. For example, Brown who comes from American Cornell University had analyzed the human activities in social life from the perspective of the natural ecology [3]. Developed so far, the research of this discipline gradually divided into two directions, one is to analyze the phenomenon of human social life using related theories and methods of natural ecology, such as the human activities, social development and so on [4][5], and the other one is to research the coordinated development between the society, economics and resources quantitatively and qualitatively based on the theory of ecology [6][7][8]. The concept of driving force system of Ecosystem of Cyber-society that proposed in this paper is a cross research between the above two disciplines. It not only emphasizes on analyzing the network space from a sociological perspective, but also reflects the theory of system dynamics in the ecology. Therefore, although the concept of driving force system of Ecosystem of Cyber-society is new, since that the research is a frontier that based on the current disciplines and research achievements, the development status and trends of related disciplines are still significant to our research [9].

The organization of this paper is as follows. In Section 2 we propose a theoretical framework of the driving force system of Ecosystem of Cyber-society, and analyze its structure, functions and characteristics. Then, in Section 3, we analyzed the operation environment and formation process of the driving forces. And finally in Section 4 we present the operation mechanism of the driving force system systematically.

2 The Driving Force System of Ecosystem of Cyber-society

2.1 Structure of the Driving Force System

The driving forces that affect the formation and development of Ecosystem of Cyber-society are not chaotic. However, they are an organic whole with a multi-level structure. The formation and development of Ecosystem of Cyber-society are driven by the driving forces that interact and constraint with each other, and ultimately form the multi-layered driving force system [10]. At the same time, the driving force system that affects the formation and development of Ecosystem of Cyber-society is also a unified whole and there is ongoing exchange of material, energy and information between the driving forces. If they can develop harmoniously and coordinately, then they will create a stronger driving force together to facilitate the rapid development of Ecosystem of Cyber-society. Otherwise, they will delay the historical process of formation and development of Ecosystem of Cyber-society. There are three major subsystems in the driving force system of Ecosystem of Cyber-society, namely subsystem of driving forces, stimulate subsystem of driving forces and carrier subsystem of driving forces, as shown in Figure 1. These subsystems interact with each other, and there are continuous exchange of material, energy and information between them. Meanwhile, as a constituent of the driving force system, when the driving forces function on the formation and development of Ecosystem of Cyber-society, their function mode and strength

would be different for different driving forces, and thus form the unique structure of the driving force system [11].

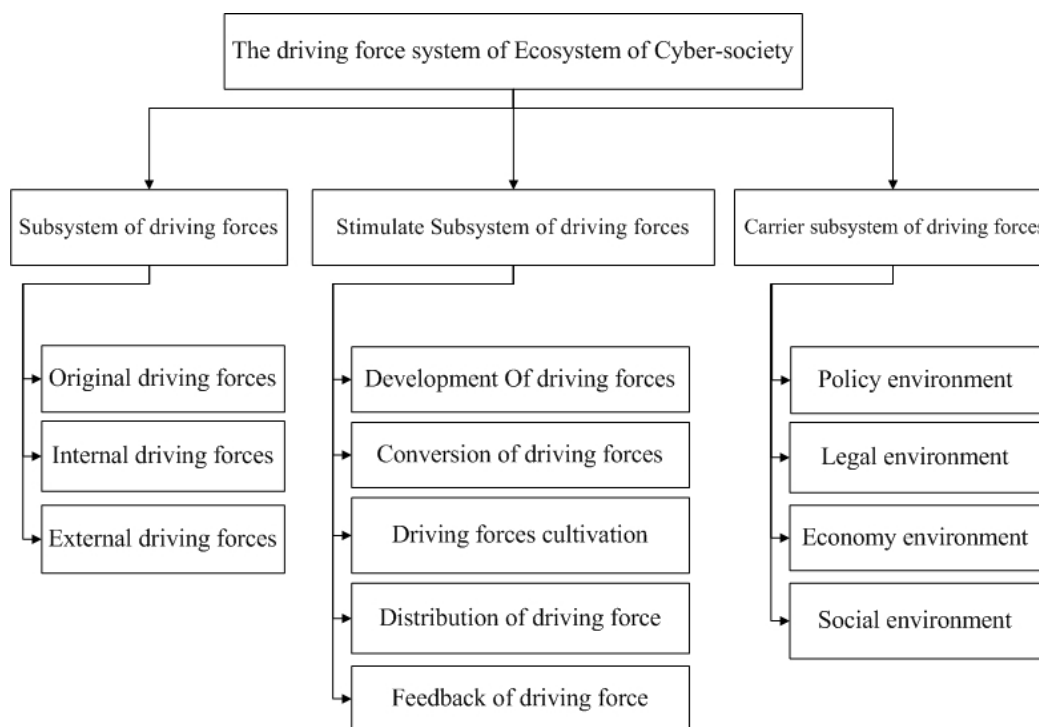


Figure 1: A theoretical framework of the driving force system of Ecosystem of Cyber-society

2.2 Functions of the Driving Force System

The driving force system that affects the formation and development of Ecosystem of Cyber-society has the following functions [12].

(1) To promote the formation of a joint force that affects the formation and development of Ecosystem of Cyber-society.

The formation and development of Ecosystem of Cyber-society is affected by a variety of driving forces together, and the driving forces vary in size, strength, and thus have positive or negative effect. Due to the existence of the driving force system, it makes the driving forces coordinate with each other, and form a joint force that advantageous to the formation and development of Ecosystem of Cyber-society.

(2) Ensure that the driving forces that affect the formation and development of Ecosystem of Cyber-society have a continuous nature.

The short-term characteristics and changes of the driving forces make the formation and development of Ecosystem of Cyber-society in the status of instability. Well operation of the driving force system can ensure the formation and development of Ecosystem of Cyber-society have a continuous driving force.

(3) To promote the optimal combination of the elements of Ecosystem of Cyber-society.

There are many driving forces that interrelate and interact with each other affect the Ecosystem of Cyber-society, and different combinations of the driving forces will play different roles in promoting. The driving force system can promote a better combination of the driving forces to some extent.

(4) To improve the competitiveness of the subjects of Ecosystem of Cyber-society.

The driving forces are the external forces that promote the formation and development of Ecosystem of Cyber-society. They combine into a unified whole through the driving force system and input negative entropy flow into the system continuously, thereby enhancing the competitiveness of the subjects of Ecosystem of Cyber-society.

2.3 Characteristics of the Driving Force System

The driving force system that affects the formation and development of Ecosystem of Cyber-society is an organic whole with multi-levels and multi-factors, and it has the following characteristics.

(1) Integrity

The integrity of the driving force system that affects the formation of Ecosystem of Cyber-society refers to the unified whole that constitutes by various driving forces that interrelate and interact with each other according to some law or in order to achieve a particular purpose, rather than the chaotic random combination of the driving forces. At the same time, the driving force system is not the simple summation of the driving forces, however, it may often produce new functions as a whole that the individual driving force doesn't have. It can be expressed in mathematical language as the following:

$$D = a \sum_{i=1}^n d_i \quad (1)$$

Here,

d represents the function value of the driving forces in the driving force system;

D represents the function value of the overall driving force system;

a represents the function coefficient of the driving forces.

(2) Structure hierarchy

The structure hierarchy of the driving force system that affects the formation of Ecosystem of Cyber-society refers to the hierarchical order of the status, role, structure and function of the system due to the various differences between the driving forces. The structure hierarchy is relative and changes constantly. When the driving forces form a whole with a reasonable structure, then the system will have new functions, and the functions will be larger than the overall functions of individual driving forces. Otherwise, when the driving forces form a whole with a poor structure, then it will damage the functions of the whole system. The expression that uses mathematical language is:

$$D = f(d_1, d_2, d_n) \quad (2)$$

Here,

$f(d_k), k = 1, 2 \dots n$ represents the structure function of the driving force system, and different levels of the structure correspond to different functional forms.

a represents the function coefficient.

(3) Orderliness

The orderliness of the driving force system that affects the formation of Ecosystem of Cyber-society refers to the characteristic of reasonable structure that combined through the coordination and promotion between the various driving forces. The orderliness is relative to the disorder, it shows the structure rationality of the driving force system, and also determines whether the system can function normally. The larger the orderliness of the driving force system, the smaller the obstruction to the formation and development of Ecosystem of Cyber-society and the faster the pace of development. At the same time, the orderliness of the driving force system is not

static, it shows the characteristic of dynamic changes in the long term along with the positive and negative effects on the system between the driving forces, but its overall characteristic and performance maintain stable.

(4) Openness

The openness of the driving force system that affects the formation of Ecosystem of Cyber-society refers to that the driving forces not only have to coordinate and cooperate with each other, but also have to exchange material, information and energy with the external environment continuously at the same time, in order to maintain the specific structure and functions of the driving force system. This open process allows the environment to influence and transform the driving forces, and at the same time the driving force system will also transform the external environment and make it evolve toward its own beneficial development direction. The openness is the guarantee to the orderliness of the driving force system, and just because of the openness of the driving force system, it makes the system can absorb negative entropy flow from the external environment in order to maintain or decrease the entropy within the system continuously, and thus keep the orderliness of the system and further enhanced.

3 Formation of the Driving Forces of Ecosystem of Cyber-society

3.1 Operation Environment of the Driving Force System of Ecosystem of Cyber-society

The so-called operation environment of the driving force system of Ecosystem of Cyber-society refers to the carrier subsystem of driving forces. While the environment is always relative to a center thing, the environment will be different due to different centers, and varies with the changes of the central things. The things associated with the system constitute a collection called the environment of the system.

Similarly, the operation of the driving force system of Ecosystem of Cyber-society also needs certain environmental conditions. The environment of the driving force system of Ecosystem of Cyber-society refers to the sum of the various driving forces that affect the operation and development of Ecosystem of Cyber-society. The various driving forces form an internal system with compact structure through the interaction with each other, and there is the real society outside with the sub-environment of economics, politics and culture. The environment of the driving force system of Ecosystem of Cyber-society represents the potential long-term impact of the real society on Ecosystem of Cyber-society and it is the manifestation of indirectly, more vague, and broad-brush driving forces, and provides environmental support for generating driving forces. The formation and development of Ecosystem of Cyber-society has a strong feature of dynamic, and due to the dual intertwine between the real society and the environment of Ecosystem of Cyber-society, it makes the formation and development of Ecosystem of Cyber-society in an uncertain environment. The driving force system is constrained by these sub-environment systems, and its function on the various driving forces is an important guarantee for the well operation of the driving force system.

3.2 Formation Process of the Driving Forces of Ecosystem of Cyber-society

The so-called formation process of the driving forces of Ecosystem of Cyber-society is just the process that the stimulate subsystem of driving forces functions. The function process of the stimulate subsystem of driving forces consists primarily of the following five sessions, namely develop, transform, cultivate, allocate and feedback, as shown in Figure 2.

(1) Source develop of driving forces

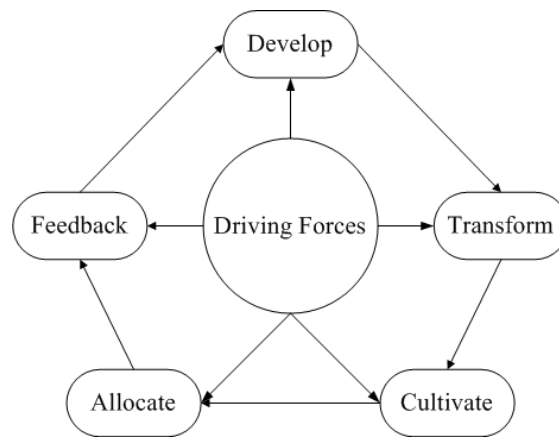


Figure 2: The function process of the stimulate subsystem of driving forces of Ecosystem of Cyber-society

The subjects' participation and use of Ecosystem of Cyber-society are the main reasons that lead to the formation and rapid development of Ecosystem of Cyber-society. In other words, there would be no Ecosystem of Cyber society without the subjects' participation and use no matter how quickly the development of modern science and technology is. The motives of the subjects' behavior are reflected in the interest and benefit, and both of these two aspects are also the major incentives that attract the subjects of Ecosystem of Cyber-society to involve in the network social life. The network economy has also been considered as the eyeball economy, and it means that one will be able to succeed in the competition of Ecosystem of Cyber-society if he can attract the attention of Internet users, seize the interest of Internet users, and meet the benefit of Internet users. The saying of eyeball economy may be not correct, but there is one thing for sure that the Internet users' interest plays a huge role in the formation and change of the current network pattern of Ecosystem of Cyber-society.

(2) Transform of driving forces

The interest and benefit are the sources of driving forces of Ecosystem of Cyber-society, and they are just the potential driving forces. Only if they are transformed into real driving forces, then they can promote the network social behavior of the subjects to meet their needs, and really play an important role as driving forces.

(3) Cultivate of driving forces

The driving forces not only need to be developed and transformed, they are also need to be cultivated and maintained, in order to facilitate the formation and sustainable development of Ecosystem of Cyber-society. Meanwhile, for the driving forces in different levels and different types, the focus of cultivation is also different. It will greatly promote the rapid and healthy development of Ecosystem of Cyber-society to seize the key influencing forces.

(4) Allocate of driving forces

In order to promote the formation and health development of Ecosystem of Cyber-society, we not only have to develop, transform and cultivate appropriate driving forces, but we also have to allocate these relatively independent driving forces to the parts of Ecosystem of Cyber society that need them most, and make them related with each other and combined into a rational structure, then the driving force system can promote the formation and development of Ecosystem of Cyber-society as an organic whole role.

(5) Feedback of driving forces

As a closed loop feedback system, the session of feedback can guarantee the good functioning of the driving force system. After the first few sessions, the Ecosystem of Cyber-society has

obtained some driving forces. However, they are still need the session of feedback to coordinate and make sure that they are adequate. The session of feedback will collect the information about whether these driving forces promote the benign operation and healthy development of Ecosystem of Cyber-society, and then make judgments and give feedback to guide the driving force system to adjust the direction and strength of the driving forces or adjust the combinations between the driving forces, and thus maximize the overall objective of the formation and development of Ecosystem of Cyber-society.

4 Operation Mechanism of the Driving Force System of Ecosystem of Cyber-society

4.1 Operation Life Cycle of the Driving Forces of Ecosystem of Cyber-society

The resources within Ecosystem of Cyber-society can be allocated optimally and achieve the Pareto optimal through the good operation of the driving force system of Ecosystem of Cyber-society, and thus promote the continuous escalating and innovation of the formation and development stages of Ecosystem of Cyber-society. The operation life cycle of the driving force system can be divided into the following three stages, namely initial stage, integration stage and mature stage, and each stage has a different operation mechanism which plays a leading role, as shown in Figure 3 [13].

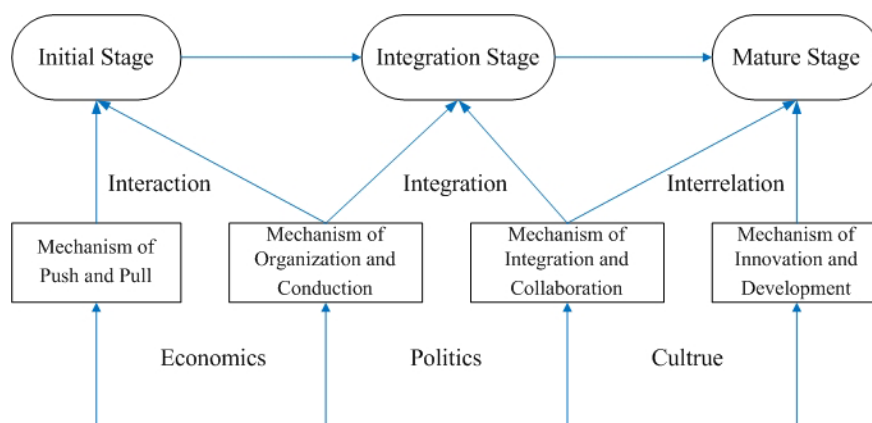


Figure 3: The operation lifecycle of driving force system of Ecosystem of Cyber-society

It has been discussed earlier in this paper about the promotion role that the sub-environment of economics, politics and culture plays on the formation and development of Ecosystem of Cyber-society. It needs to be noted that the force of economics, politics and culture has to infiltrate into the system through the driving force system during the formation and development process of Ecosystem of Cyber-society and then function continuously. The promotion force that the driving force system giving to the Ecosystem of Cyber-society under the interaction between economics, politics and culture should be appropriate and it can be understood specifically from the direction, size and cooperation between the driving forces [14].

(1) Direction of the driving forces. To judge whether the direction of the driving forces is appropriate, we can see whether it is consistent with the overall goal and direction of the development of Ecosystem of Cyber-society. Positive driving force is conducive to the formation and rapid development of Ecosystem of Cyber-society, while negative driving force will inhibit the formation and development of Ecosystem of Cyber-society. Therefore, the force of economics, politics and culture that based on the deep roots of development needs of the subjects of Ecosys-

tem of Cyber-society should be positive, should be conducive to meet the basic survival needs of the subjects of Ecosystem of Cyber-society, should be conducive to increase the high-level needs and satisfaction degree of the subjects of Ecosystem of Cyber-society, and also should be conducive to the coordination between the individuals and groups of Ecosystem of Cyber-society.

(2) Size of the driving forces. As to the positive driving force, its size and strength determine the formation and development speed of Ecosystem of Cyber-society to some extent, that is to say, the larger the value is, the faster the system develops, and vice versa. The "interest" and "benefit" of the subjects promote the formation and rapid development of Ecosystem of Cyber-society, but the increase of interest level and the expectation of benefit pursuit must adapt to the reality of the development of Ecosystem of Cyber-society, and this needs the driving forces to adjust. The Ecosystem of Cyber-society with inadequate driving forces is dull and lack of energy, while the unlimited and vigorous expanded pursuit to the interest and benefit would cause disorder and concussion of Ecosystem of Cyber-society. An appropriate driving force can combine the satisfaction of the subjects interest and the good running order of Ecosystem of Cyber-society together. It can not only activate the driving force of reasonable interest pursuit of the subjects of Ecosystem of Cyber-society, but it can also control the various driving forces and activities that affect the formation and development of Ecosystem of Cyber-society within reasonable limits.

(3) Cooperation between the driving forces. Only if the force of economics, politics and culture can complement with each other, coordinate their actions and optimize the allocation of various needs and interest pursuit of the subjects of Ecosystem of Cyber-society, then it can give appropriate positive driving force to the formation and development of Ecosystem of Cyber-society, and thus effectively promote the formation and development of Ecosystem of Cyber-society.

4.2 Mechanism of Push and Pull

The formation and development of Ecosystem of Cyber-society is a dynamic process, and influenced by various driving forces during the process of dynamic development. Combined with the mechanism analysis of each driving force, the joint effect in the process of formation and development of Ecosystem of Cyber-society can be shown as Figure 4. Here, the ball represents the Ecosystem of Cyber-society, the x -axis represents the time stage of the formation and development of Ecosystem of Cyber-society, the y -axis represents that the formation and development of Ecosystem of Cyber-society is a tough climb process with load, and F_i represents the various driving forces, $i = 0, 1, 2, \dots, n$ [15][16].

The roles that each driving force plays in the formation and development of Ecosystem of Cyber-society have different directions and sizes. According to the theory of Mechanics, the Ecosystem of Cyber-society will be in the state of fast forward development when the positive driving force is larger than the negative driving force, while the Ecosystem of Cyber-society will be in the state of uniform development when the positive driving force is equal to the negative driving force. And the Ecosystem of Cyber-society will be in the state of stagnation or even retrogression development when the positive driving force is smaller than the negative driving force. Thus, it will greatly promote the formation and development of Ecosystem of Cyber-society by increasing the positive driving force and making full use of various driving forces.

4.3 Mechanism of Organization and Conduction

The mechanism of organization and conduction refers to the formation process of the joint force that affect the formation and development of Ecosystem of Cyber-society through the cross

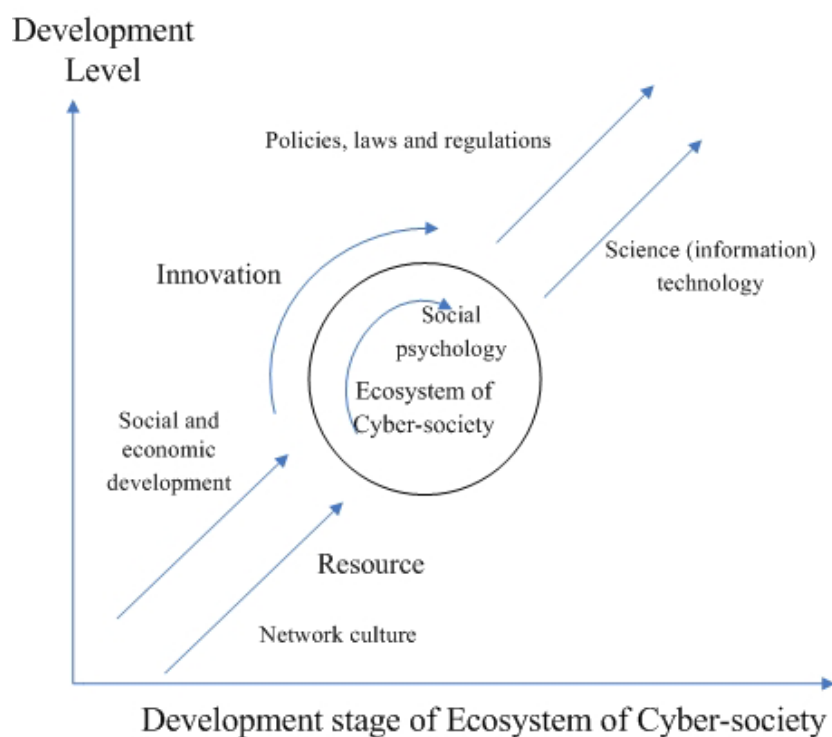


Figure 4: Composition forces model of the formation and development of Ecosystem of Cyber-society

combination and strength coordination between the various driving forces within the driving force system.

The various driving forces interact and influence with each other, thus form a joint force for the formation and development of Ecosystem of Cyber-society. During the process of the formation and development of Ecosystem of Cyber-society, the joint force it gets is not the simple accumulation of each driving force, however, there will be a new force with new quality due to the interaction and its size is often larger than the simple sum of the driving forces. As shown in Figure 5, the various driving forces first act on the Ecosystem of Cyber-society in accordance with their respective different mechanisms, and thereby combine into a new force through the interaction between them, then this new force will combine with other forces, and the cycle continues, eventually form the integrated force that promotes the formation and development of Ecosystem of Cyber-society.

During the process of interaction between the various driving forces, there will be part of the forces being reduced or offset because of various practical conditions limit. Therefore, in order to speed up the formation and development of Ecosystem of Cyber-society, the relationships between the various driving forces have to be coordinated as much as possible in order to increase the joint effect, which need to optimize the organizational structure of the driving force system or use regulations, as shown in Figure 6. By optimizing the hierarchy of the driving force system, making the driving forces all in their places and out of their functions, the joint force can conduct layer by layer orderly, effectively reduce the deadweight loss during the process of force conduction, and thus make the formation and development of Ecosystem of Cyber-society got stable promotion.

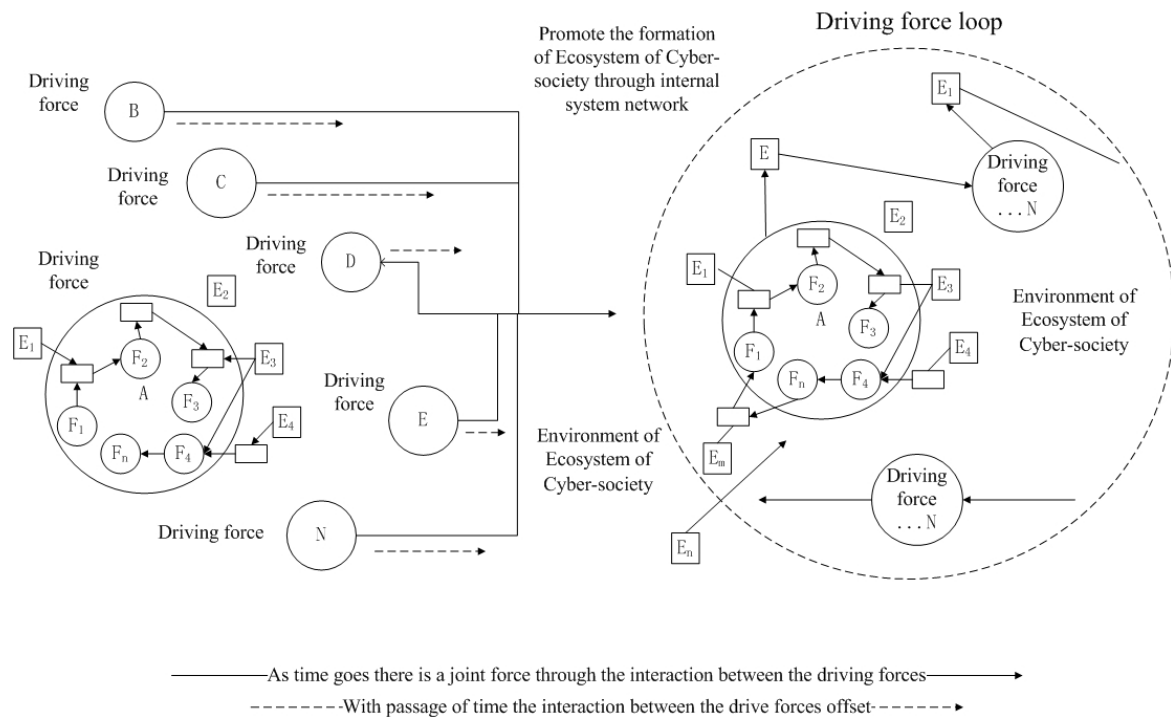


Figure 5: Cross combination of driving forces to the formation of Ecosystem of Cyber-society

4.4 Mechanism of Integration and Collaboration

The collaboration generally refers to the relationship of mutual adaptation between two or more interacting species in the process of evolution. To understand from a broad concept, the collaboration can refer to the coordination in the process of long-term mutual adaptation between species and species, or between species and environment. Similarly, this collaboration relationship also widely exists in the Ecosystem of Cyber-society. During the process of the formation and development of Ecosystem of Cyber-society, there are interrelationships between the various subjects, or between the subjects and the environment. And only to form a coordinated and mutually reinforcing relationship between them, then they can promote the formation and development of Ecosystem of Cyber-society effectively.

The formation and development of Ecosystem of Cyber-society is not a simple question, and we should consider the formation and development objectives of Ecosystem of Cyber-society from the height of the overall socio-economic development, and take the development of Ecosystem of Cyber-society as an important element of real economic and social development. The so-called mechanism of integration and collaboration of driving force system of Ecosystem of Cyber-society refers to that make clear the positive role of Ecosystem of Cyber-society in the real economic and social development, take the Ecosystem of Cyber-society as an important part of the real society, make the development of Ecosystem of Cyber-society integrated into the overall development of the real society, strengthen the linkage between Ecosystem of Cyber-society and the real society, maximize the efficiency of resource allocation and utilization, and thus form a integrated force to facilitate the formation and development of Ecosystem of Cyber-society.

From the perspective of the functions of the driving force system of Ecosystem of Cyber-society, the mechanism of integration and collaboration includes the overall coordination between the driving force system and Ecosystem of Cyber-society, while from the perspective of the collaboration between various driving forces, the mechanism of integration and collaboration

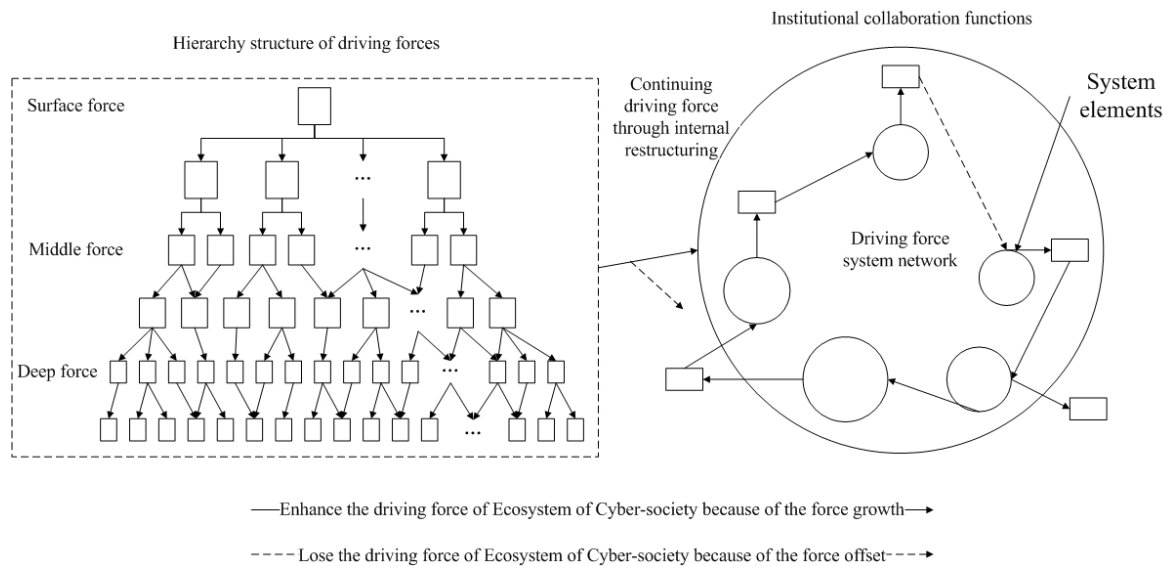


Figure 6: Organizational conduction of driving forces to the formation of Ecosystem of Cyber-society

also includes both of the collaboration within the system and outside the system, as shown in Table 1.

Table 1: Contents of the mechanism of integration and collaboration of driving force system of Ecosystem of Cyber-society (EC)

Scope of integration and collaboration	Contents of integration and collaboration
Internal collaboration	Social psychology and EC
	Innovation and EC
	Network culture and EC
External collaboration	Resources and EC
	Network culture and EC
Overall collaboration	Social and economic development and EC
	Policies, laws and regulations and EC

5 Conclusions

We propose a basic theoretical framework of the driving force system of Ecosystem of Cyber-society and present its operation mechanism using the theory of System Dynamics in this paper. The driving force system of Ecosystem of Cyber-society is constituted by the subsystem of driving forces, stimulate subsystem of driving forces and carrier subsystem of driving forces, which is a virtual manual system with the characteristics of integrity, structure hierarchy, orderliness and openness. Its operation mechanism mainly includes the mechanism of push and pull, mechanism of organization and conduction, mechanism of Integration and collaboration, and mechanism of innovation and development.

The formation and development of Ecosystem of Cyber-society is the function result of a variety of driving forces together. There is no one single driving force can independently control the formation and development of Ecosystem of Cyber-society. Only through the coordination

between the driving forces within and outside the system, then it can generate new driving force structures and joint force, improve the operation efficiency of the driving force system, and thus promote the fast forward development of Ecosystem of Cyber-society with its powerful functions.

Acknowledgements

This work is partially supported by the General Research Project (18190113002) and the Quality Improvement Project of Personnel Training (03150113016) of Beijing Municipal Education Commission, and also the Institute Level Project (E-b-2012-20) and the Course Construction Project (22150112088) of Beijing Institute of Graphic Communication.

Bibliography

- [1] Custer, M. (2001); Rise of the Network Society, Social Science Academic Press (China): 59-76.
- [2] Qi, G. (2003); Essence of the Networking Society: A Digitalization Social Relationship Structure, *Journal of Chongqing University (Social Science Edition)*, 9(1).
- [3] Brown (1979); Ecology of human development.
- [4] Iansiti, M., Levien, R. (2004); Strategy as Ecology, *Harvard Business Review*, 82(3).
- [5] Iansiti, M., Levien, R. (2004); The Keystone Advantage: What the New Dynamics of Business Ecosystems Mean for Strategy, Innovation, and Sustainability, *Harvard Business School Press*: 255.
- [6] Feng J. (2010); Research on the Coordinated Development between Internet Uses and Network Environment of the Ecosystem of Cybersociety, *Beijing Jiaotong University*.
- [7] Cubillos, C., Donoso, M., Rodriguez, N., Guidi-Polanco, F., Cabrera-Paniagua, D. (2010); Towards Open Agent Systems Through Dynamic Incorporation, *International Journal of Computers, Communications & Control*, 5(5): 675-683.
- [8] Neghina, D., Scarlat, E. (2013) Managing Information Technology Security in the Context of Cyber Crime Trends, *Int J Comput Commun*, ISSN 1841-9836, 8(1): 97-104.
- [9] Zhang Z., Zhang R. (2008); Ecosystem of Cyber-society, *Beijing: Publishing House of Electronic Industry*.
- [10] Guan X. (2011); Research on the Formation Mechanism of Ecosystem of Cyber-society, Doctor Dissertation, *Beijing Jiaotong University*.
- [11] Barrett, L., Henzi, S., Lusseau, D. (2012); Taking sociality seriously: the structure of multi-dimensional social networks as a source of information for individuals, *Philosophical Transactions of the Royal Society b-biological sciences*, 367(1599): 2108-2118.
- [12] Yang X. (2002); Dynamic Mechanism of Structure, Function and Operation Process of the Social Development, *Journal of the Party School*, 6(4): 28-33.
- [13] Wang X. (2008); The Theoretical and Empirical Research On the Driving Mechanism of Urban Tourism Development, Doctor Dissertation, *Tianjin University*.

- [14] Yang X. (2002); Dynamic mechanism of social development, *Social Sciences in Guangdong*, 6: 87-94.
- [15] Li B. (2003); Research on the Enterprise Development Dynamics, Doctor Dissertation, *Harbin Engineering University*.
- [16] Xu J., Ma R. (2007); Dynamic mechanism model of the enterprise ecological development, *Productivity Research*, 17: 116-118.

Localization of Wireless Sensor Network Based on Genetic Algorithm

N. Jiang, S. Jin, Y. Guo, Y. He

Nan Jiang*, Sixin Jin, Yan Guo

College of Information Engineering, East China Jiaotong University
Nanchang, Jiang Xi 330013, P.R.China
jiangnan@ecjtu.jx.cn, aquajsx@gmail.com, guoyan997@gmail.com

*Corresponding author: jiangnan@ecjtu.jx.cn

Yueshun He

College of Information Engineering, East China Institute of Technology
Nanchang, Jiang Xi 330013, P.R.China
hys8418@163.com

Abstract: This paper proposes a novel localization approach based on genetic algorithm for Wireless Sensor Networks. In this method, we use a new way to approximate the distance between anchor node and unknown node which is out of the anchor nodes' communication radius. In addition, we use self-adapting genetic algorithm into localization, which will ensure it can produce the result as similar as its real position in any environment. The experiments on various network topologies show the better results. In comparison, we find that previous anchor node free localization approach cannot work well in the unidealization environment. The demonstration explains that the approach can help unknown nodes obtain high accuracy position whether in open space, the environment with obstruction, or even the unconnected well environment.

Keywords: Wireless Sensor Networks(WSNs), Genetic Algorithm(GA), Localization

1 Introduction

Localization is one of key supporting technologies to Wireless Sensor Networks (WSNs). It could provide accurate position information for kinds of expanding application. In this paper, we study the localization problem for a 2D large-scale WSNs with obstruction. If we equip GPS for every sensor node, it will raise the cost much more, and it also brings a new problem of nodes' power support. Can we only use less GPS for large-scale WSNs to obtain the position of nodes, in obstruction or unconnected well environment?

The major challenge is reducing the influence of obstruction between pairs of nodes. How to achieve high accuracy and reliability localization scheme for WSNs that is always the research subject to researcher around the world, in the environment with obstruction or unconnected well. At the beginning of the localization study on WSNs, Received Signal Strength Indicator(RSSI) [1], Time of Arrival(TOA) [2], Time Difference of Arrival(TDOA) [2] [3], Angle of Arrival(AOA) [4] and other range-based localization have been proposed. In addition, centroid algorithm [5], Ad-Hoc positioning system(APS) [6], Amorphous [7], APIT [8] and other range-free localization also have been proposed. Through literature [1] [2] [3] [4] [5] [6] [7] [8], we found that these algorithms only could provide low accuracy position of nodes for us, and couldn't meet the needs of expanding application on WSNs.

There are many approaches with only network connectivity to calculate the position of nodes such as multidimensional scaling (MDS) [9]. It gives us an O.K. localization results on a well-connected dense network (show in Figure 7). But it does not have any well results in the environment with obstruction (show in Figure 9).

More recently, there are many researchers put graph rigidity theory into localization of WSNs [10] [11] [12], they have made a lot of contribution on it. The methods proposed by S. Lederer

et al. just use the network connectivity information to get the position of unknown nodes. In order to obtain the nodes' position, they partition the network into Voronoi cells with using Delaunay triangles closest to the shape of region of concern. These kinds of methods are good choice of localization plan for underwater, underground or indoor environment, but in these way we can't localization in the environment with obstruction or unconnected well. Apart from this, the locations of landmarks should know before localization. How to obtain the nodes' position in obstruction or unconnected well environment also is the research hotspot for us.

Our contribution. We propose a novel localization algorithm, it uses Genetic Algorithm(GA) [13] [14] into localization of WSNs, the approach causes nodes in the environment with obstruction or unconnected well could get high accuracy position by themselves.

We assume WSNs has been deployed in a region with obstruction, in the network only fewer nodes equipped GPS which could localization by themselves, but others are not. All nodes only could communicate with others in its communication radius, unknown nodes could get the position of anchor node and distance between with anchor node. When there are obstruction stop communicating pairs of nodes, unknown nodes could make use of other nodes in its communication radius to get the position of anchor nodes, and they could obtain distance with anchor nodes through the method provide in Section 4. At last, every unknown node could estimate the position by the global search feature of GA.

Through the simulation result in this paper, we prove that the method provided in this paper has a well property in the environment with obstruction than any other methods, especially it's could localization well in the environment where is separated into several parts like Figure 1.



Figure 1: Localization environment with obstruction, these pictures come from Google Maps.

The outline of the paper is as follows: Section 2 presents the theoretical foundations of genetic algorithm. The localization approach based on self-adapting genetic algorithm is proposed in Section 3. Section 4 presents the evaluation of system and compare simulate results use a different approach. Section 5 concludes this paper.

2 Theoretical Foundations

Genetic Algorithm (GA) is a bionics optimize algorithm. It has been put forward by John Holland on 1975 [13] [14]. It is based on Darwin's biological evolutionism which is "survival of the fittest" and Mendel's genetic variation principle which is "biological genetic evolution mainly take place in the chromosome". GA through selection, crossover and mutation to die out the low-fitness individuals and maintain high-fitness individuals like in the real nature, it can make the population towards to optimize and converge to a global optimum individual. TABLE 1 shows the pseudocode of the genetic algorithm.

In this paper, we use the self-adapting GA which has been proposed in literature [15], the algorithm could change the operate points by the different multiformity of group. We should set up some parameters like code length L , group size N , cross probability P_c , mutation probability

Table 1: Pseudocode of genetic algorithm

Step 1: Initialization
The number of individuals, NIND;
The maximum number of generation, MAXGEN;
The precision of variables, PRECI;
The generation gap, GGAP;
Initialized the times of generation gen= 0;
<hr/> Step 2: Coding
Initialize the group by Gray, P(t);
<hr/> Step 3: Evaluate
Evaluate the fitness of each individual in the group P(t);
<hr/> Step 4: Genetic Iteration
While t < MAXGEN
Selecting the individuals for Crossover;
Crossing the selected individuals by certain probability;
Mutating the individuals of group by certain probability;
Evaluating the fitness of the new group;
Producing a new group after the evolution, P(t);
Partial Best = min(Evaluate P(t));
Update the times of generation, t = t + 1;
End
<hr/> Step 5: Obtain the result
Global Best = Partial Best;

P_m and so on. Using GA into the node localization approach of WSNs, we make use of the partial search ability and the global search ability of GA. The partial search ability will help nodes in class I (shown in Section 3.1) to obtain its position; the global search ability will help nodes in class II (shown in Section 3.1) to obtain its position, these two kinds of abilities are supplied by cross probability and mutation probability.

3 Algorithm Description

The localization approach based on genetic algorithm (GAL) proposed in this section is divided into two phases: firstly unknown nodes obtain the distance to adjacent anchor nodes; secondly unknown nodes use genetic algorithm to estimate unknown node's position.

3.1 Obtain the Distance

After all sensor nodes been deployed, all anchor nodes begin broadcasting their position information. Then unknown nodes could obtain messages with anchor nodes' position directly which is in its radius, and they also could obtain other anchor nodes' position through other nodes by routing forwarding. At last, unknown nodes metric distance with adjacent anchor nodes through RSSI. After these steps, unknown nodes are classified into two categories:

Definition 1. Consider an unknown node has obtained the distance to three or more anchor nodes called Class I nodes. We say that this kind of nodes can use genetic algorithm to calculate its own position directly.

The calculate method will proved in section 3.2.

Definition 2. Consider an unknown node hasn't obtained the distance of which to three or more anchor nodes called Class II nodes. We say that this kind of nodes cannot use genetic algorithm to calculate its own position directly.

Now we are ready to explain how to estimate the distance for Class II nodes. Shown in Figure 2.

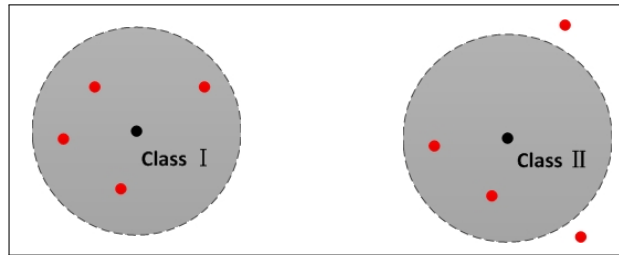


Figure 2: Two types of unknown nodes in the WSNs. The red points mean anchor nodes and the black means unknown nodes.

Definition 3. There are two kinds of nodes, one is anchor node such as node A, the other is unknown node such as node C. Node K is either kind of nodes in the network. Node K can communicate with node A and C. Let node K to node A farther than to node C. Node B also is an assume node on the circle K, the radius of this circle is the distance of between node K and node C. Shown in Figure 3.

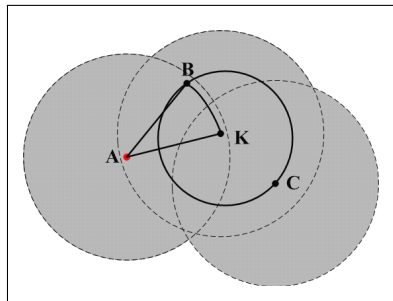


Figure 3: How to estimate the distance to anchor nodes for Class II nodes. Every node has the same communication radius, and it is marked by gray area.

Definition 4. Node D an assumed node which is the intersecting point of the extended line of line AK and circle K. θ is the included angle of line AK and KB, then $\varphi = \pi - \theta$. Shown in Figure 4.

Lemma 5. For isocetes triangle ABK, it has

$$\sin \frac{\theta}{2} = \sqrt{\frac{AK^2 + (CK - AK)^2}{4AK^2}} \tag{1}$$

Lemma 6. For the small circle of concentric circles K, it has

$$\widehat{BD} = \varphi \cdot BK \tag{2}$$

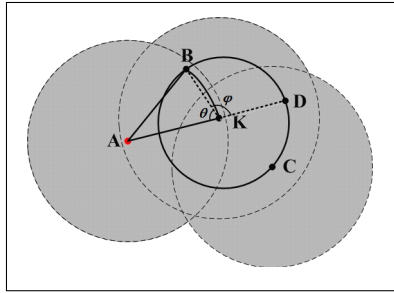


Figure 4: How to estimate the distance to anchor nodes for Class II nodes.

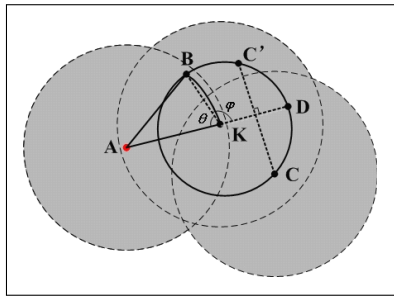


Figure 5: Class II unknown nodes get the distance to anchor nodes.

Based on Lemma 5. and Lemma 6. we can conclude that:

Theorem 7. For Class II nodes like the node C in the Figure 5, the distance between node A and C is equal to the distance between node A and C'. In this paper, we use the sum of AB and the length of \widehat{BC} to similar the distance between node A and C'. As point C is an uncertain node, point C' is also, it could appear on any position of \widehat{BD} . Because of this, we use the half length of \widehat{BD} to measure \widehat{BC}' . According this, we estimate the distance between anchor node and Class II unknown node as following

$$AC = AC' \approx AB + \widehat{BC}' \approx AB + \frac{1}{2}(\widehat{BD}) \tag{3}$$

Proof: Node C' is the mirror image of node C, it subjects to uniform distribution on \widehat{BD} . So if $C \sim U [B, D]$, then the probability distribution of node C subinterval $[b, d]$,

$$P(b \leq C \leq d) = \int_b^d \frac{1}{B - D} dx = \frac{d - b}{B - D} \tag{4}$$

In order to estimate the distance between node A and C', we use the expectation of the length \widehat{BC}' similar as the longer part of AC' than AB. The expectation of the length \widehat{BC}' is

$$E(\widehat{BC}') = \frac{0 + \widehat{BD}'}{2} \tag{5}$$

□

3.2 Compute Unknown Node's Position by Genetic Algorithm

This paper measure the multiformity of group by $F(t) \times \rho(t)$, to guarantee GA group's multiformity and it could expand search range during the previous period of iterate, moreover, optimize well during the later period of iterate. In this way, the self-adapting GA makes this localization approach can deal with the unconnected environment well.

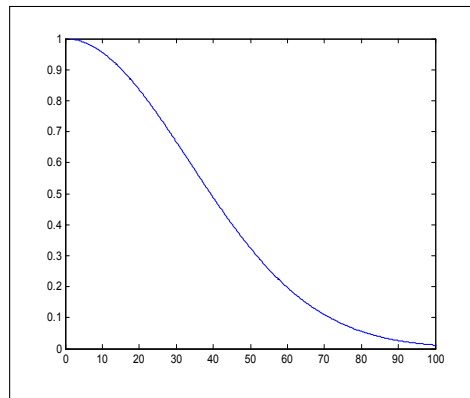


Figure 6: The function of $\rho(t)$, $T=100$.

Definition 8. In this paper, the element of measuring the group's multiformity is multiformity function $F(t)$, it is considered that

$$F(t) = 1 - \frac{H(t)}{L} \quad (6)$$

Where the average Hamming distance is

$$H(t) = \frac{\sum_{i=1}^{N-1} \sum_{j=i+1}^N H(X_i(t), X_j(t))}{\sum_{i=1}^{N-1} (N-i)} \quad (7)$$

Where the Hamming distance between individual $x_{ik}(t)$ and $x_{jk}(t)$ is

$$H(X_i(t), X_j(t)) = \sum_{k=1}^L |x_{ik}(t) - x_{jk}(t)| \quad (8)$$

In the equation (8), $x_{ik}(t)$ means the t generation's k ($k = 1, 2, \dots, L$) bit of individual i ($i = 1, 2, \dots, N$).

Definition 9. In the multiformity measure, the other element is $\rho(t)$, it is considered that

$$\rho(t) = \exp(-t^2/2\sigma^2), \sigma = T/3 \quad (9)$$

In the equation (9) of Definition 9, T means the last generation of GA iteration, when $T=100$, the figure of $\rho(t)$ like Figure 6.

Moreover, we also can't use the encoding method for simple genetic algorithm. The encoding method of simple genetic algorithm is simply binary code, as this way, it will make great hamming distance between neighbor encode elements, the distance is called Hamming Cliff [16]. The Hamming Cliff is too difficult to crossover or mutation, therefore we initialize encoding by Gray code to avoid this problem.

The most important part of genetic algorithm is fitness function. It is used to evaluate the fitness of every node's estimate position.

Definition 10. In this paper, consider the fitness function

$$F(j) = |dist(i, j) - dist(n, i)| + F(j)_{last-iteration} \quad (10)$$

Table 2: Pseudocode of localization based on genetic algorithm

Step 1: Initialization
The number of individuals, NIND;
The maximum number of generation, MAXGEN;
The precision of variables, PRECI;
The generation gap, GGAP;
Initialized the times of generation gen= 0;
The number of anchor nodes, NodeNum;
The number of aimnodes, AimNum;
The distance between anchor and aimnode, D;

Step 2: Coding
Generate the position of aimpoint by Gray code randomly, AimP(t);

Step 3: Evaluate
for n = 1; n < AimNum; n++
 t = 0;
 for j = 1; j < NIND; j++
 F(j) = 0;%F(j) is the fitness function to the group
 for i = 1; i < NodeNum; i++
 F(j) = dist(NodeP(i) - AimP(j)) - D(n,i) + F(j);
 end
 end
end

Step 4: Select, Crossover, Mutation, Reproduce
while t < MAXGEN
 Evaluating the fitness of the new group;
 Selecting the individuals for Crossover;
 Crossing the selected individuals by certain probability;
 Mutating the individuals of group by certain probability;
 Evaluating the fitness of the new group;
 Producing a new group after the evolution, AimP(t);
 Partial Best = min(F(j));
 Update the times of generation, t = t + 1;
end

Step 5: Obtain the result
AimP = Decode(Partial Best);

$dist(i, j)$ be the distance between anchor node i and unknown node j , calculating it needs the coordinate of node i and j . How to get the coordinate of anchor node i has proposed in section 3.1, and the coordinate of unknown node j could estimate by genetic algorithm. $dist(i, j)$ be defined as:

$$dist(i, j) = \sqrt{(x_{NodeP(i)} - x_{NodeP(j)})^2 + (y_{NodeP(i)} - y_{NodeP(j)})^2} \quad (11)$$

The Pseudocode of the localization method proposed in this paper is showed in TABLE 2.

4 Simulation

In this paper, we develop MATLAB toolbox to verify the effectiveness of our approach, and then contrast the experiment results by using different methods in different environment. During the follow experiment, there are 300 nodes deployed in the area $1000\text{m} \times 1000\text{m}$, the communication radius of anchor nodes and unknown nodes both are 200m. In this paper, the values of main parameters as follow, the GGAP is 0.7, cross probability is 0.85 and the mutation probability is 0.1, it is similar in other application as usual.

4.1 Localization Effect in Different Environment

From the following results, we could easily find that using MDS to get the position of unknown nodes is a good choice for WSNs, which is deployed in the open space environment (show in Figure 7). While we using GAL to localization for this environment, we can find that unknown nodes also could obtain its position accurately, if there are enough anchor nodes around it (show in Figure 7).

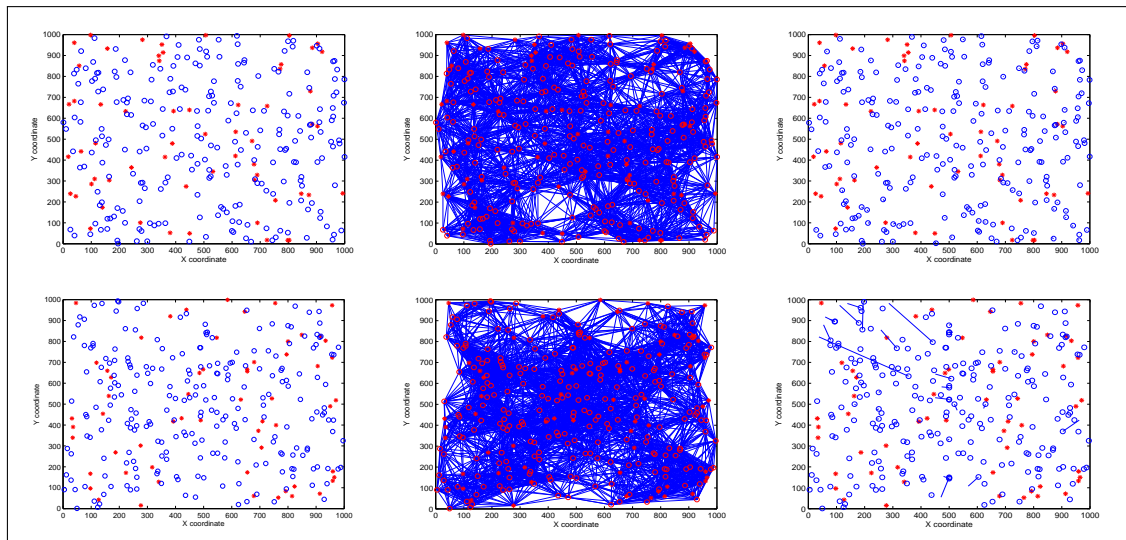


Figure 7: The first row figures are produced by MDS; The second row figures are produced by GAL. The left col figures show the nodes deploy in the open space; The middle col figures show the neighbor relationship have been built; The right col figures show the localization error.

To verify the localization ability of GAL in the environment with obstacles, we put rectangle, Triangle, circle, wall into the region of consideration. When anchor nodes are 20% of all, the localization results are shown in Figure 8.

From the simulate result of Figure 8, it is not difficulty to see that obstacles can't influence localization accuracy much, unknown node also could obtain its position within accepted range.

In the same environment as Figure 8, we use MDS to obtain unknown nodes' position again, the experiment results show in Figure 9.

From Figure 9 it's not hard to find that localization use MDS is much worse than localization using GAL. In the environment like Figure 9 (L) and (O) we can't use MDS to get the position of unknown nodes.

It is well know that MDS has the better localization accuracy than any other classical methods, in order to make this experiment completed we also use DV-hop and centroid algorithm to get the unknown nodes' position in a few environments appeared in Figure 8 and 9. The results

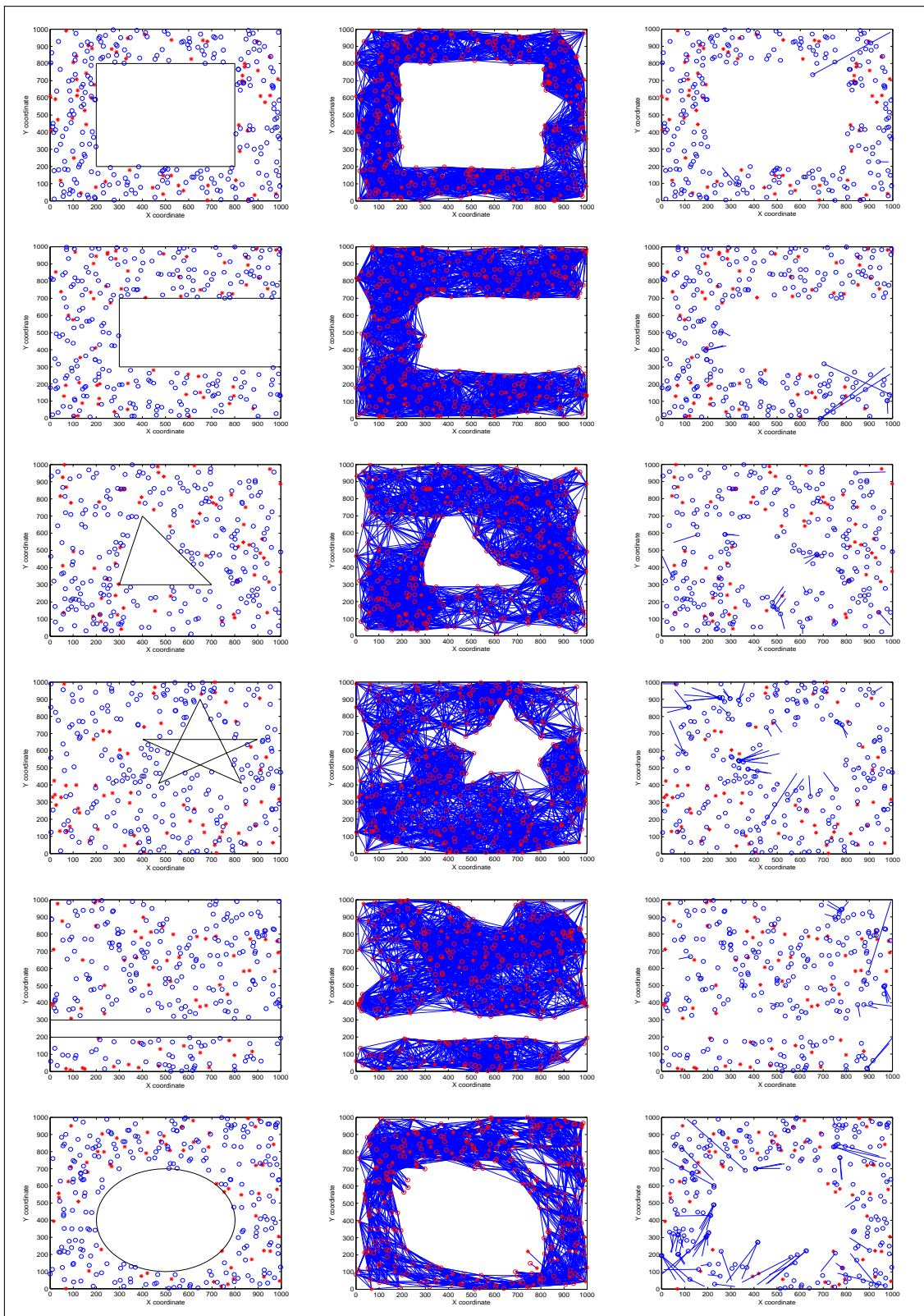


Figure 8: These pictures show the system evaluation by GAL in the environment with obstacles or unconnected. The left col figures show the nodes deploy in different environment; The middle col figures show the neighbor relationship have been built; The right col figures show the localization error of GAL.

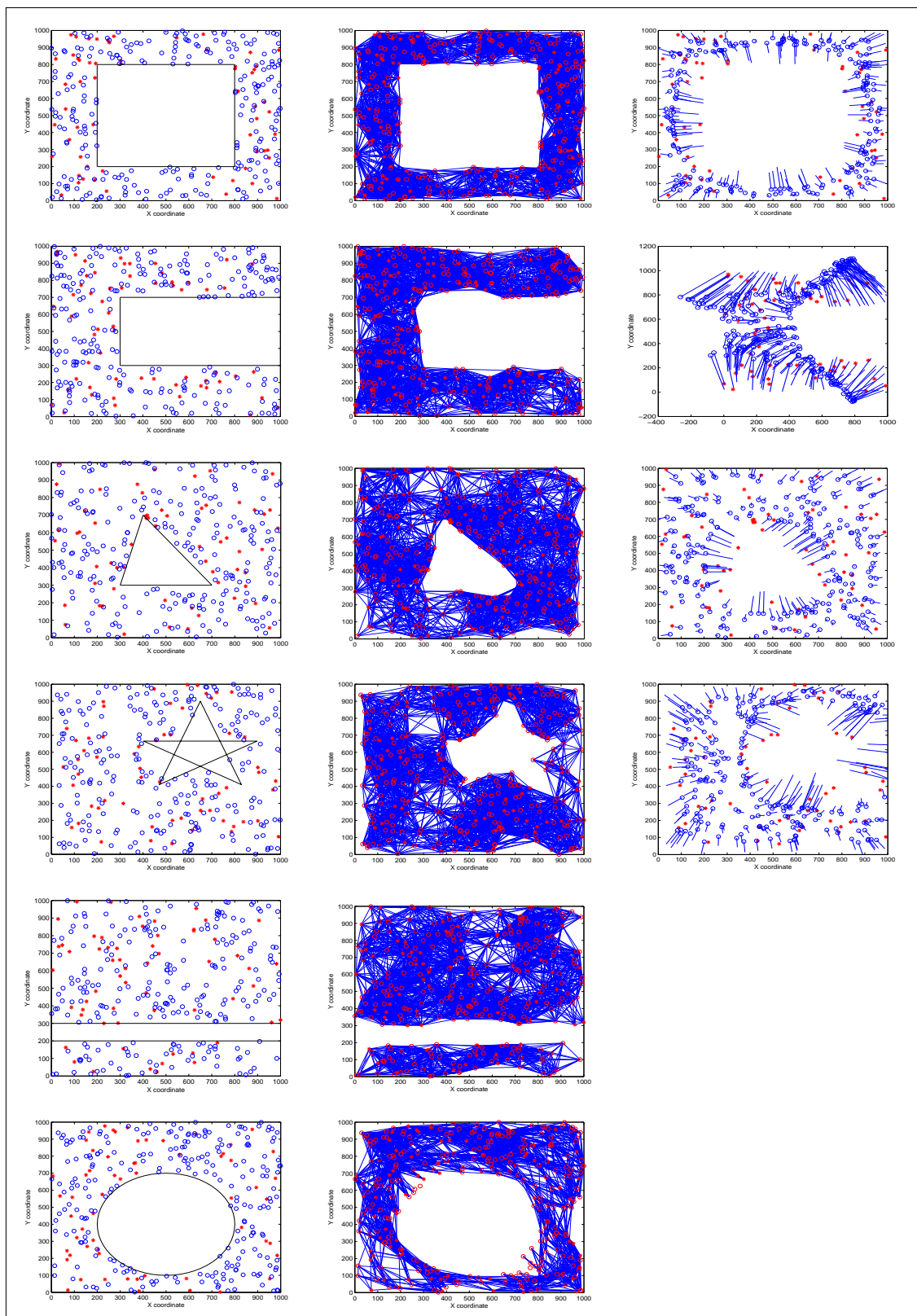


Figure 9: These pictures show the system evaluation by MDS in the environment with obstacles. The left col figures show the nodes deploy in different environment; The middle col figures show the neighbor relationship have been built. The right col figures show the localization error of MDS. But in the last two environment, we can't use MDS to get nodes' position.

of localization error is shown in table 3.

4.2 Contrast GAL with MDS on Characters

In section 4.1 we have presented the result of localization using GAL and MDS in different environment with obstacles. In order to compare with two localization algorithms quantify and objectively, we make experiment in each environment and use different anchor nodes ratio, then analysis the localization error of these results. Figure 10 shows the localization error use different anchor percentage.

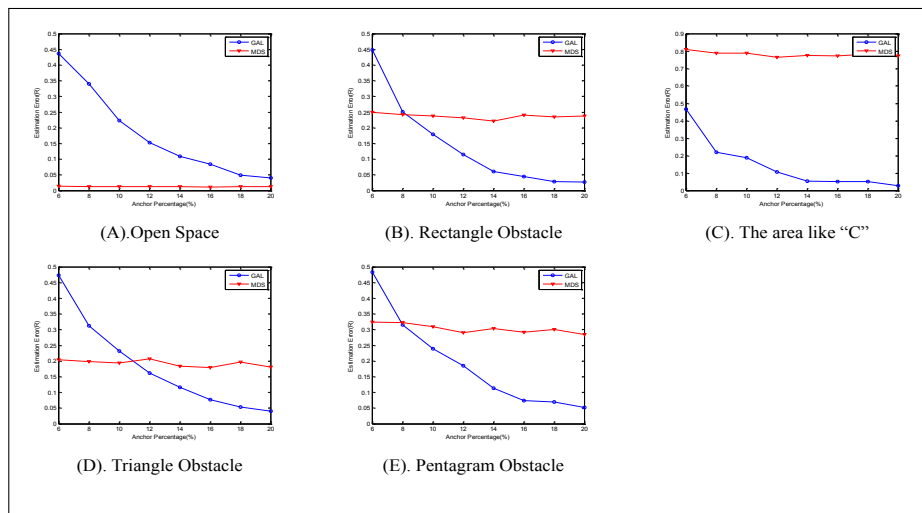


Figure 10: These pictures show the localization error in different anchor percentage by MDS and GAL. In these results, each point is the average result of ten experiments.

Figure 10 (A) shows the localization error in different anchor percentage in the open space. MDS has higher accuracy than GAL, and it just need little number of anchor nodes. The more anchor nodes, the more accuracy of GAL, but the accurate of MDS hardly change. When anchor nodes up to 20%, the localization accuracy of GAL is close to MDS.

Figure 10 (B) to (E) show the localization error comparison of two methods in the environment with obstacles, from these figures we can find that the more anchor nodes percentage, the more localization accuracy of GAL, but the accurate of MDS also hardly to be changed. When percentage of anchor nodes less than 8% in the rectangle or pentagram obstacle environment, or less than 11% in the triangle obstacle environment, MDS could get better localization accuracy than GAL, but the localization error is too high to application. However, GAL could make localization error less than 15%, when percentage of anchor nodes reaches 13%. Therefore, GAL has highly robustness and practicality than any other localization algorithms in the environment not well for communicating.

Table 3: The localization error of DV-hop and Centroid

Method	Without obstruction	Triangle obstruction	Five pointed star obstruction
DV-hop	0.34246	0.375053	0.32259
Centroid	0.318298	NaN	NaN

5 Conclusion

A novel localization approach is proposed in this paper, in this method unknown nodes through their near anchor nodes to obtain their position. In order to reduce error during localization, we use new means to approximate the distance between unknown nodes and anchor nodes when it is larger than node's communication radius. Moreover, we use self-adapting genetic algorithm to calculate the similar position of nodes, it makes the localization error much lower than the common method. Through the contrast experiments in this paper, we find this localization method can work well, whether in the environment with obstruction or not. For the aim of localization in non-planar environment like interior of tall building, we would develop a 3-D localization algorithm in the future study.

Acknowledgment

This work is supported by National Natural Science Foundation of China under Grant No.61063037 and No. 51364001, Key Projects in the Science and Technology Pillar Program of Jiangxi Province of China under Grant No. 20111BBG70031-2 and 20133BBE50033, and Foundation for Young Scientists of Jiangxi Province of China under Grant No. 20133BCB23016.

Bibliography

- [1] Girod, L.; Bychkovskiy, V.; Elson, J.; Estrin, D. (2002); Locating tiny sensors in time and space: A case study, *In Proceedings of IEEE International Conference on Computer Design: VLSI in Computers and Processors*, ISSN 1063-6404, 214-219.
- [2] Girod, L.; Estrin, D.(2001); Robust range estimation using acoustic and multimodal sensing, *In Proceedings of 2001 IEEE/RSJ International Conference on Intelligent Robots and Systems*, ISBN 0-7803-6612-3, 1312-1320.
- [3] A. Savvides; C.C. Han; M.B. Srivastava(2001); Dynamic fine-grained localization in ad-hoc networks of sensors, *In Proceedings of 7th Annual International Conference on Mobile Computing and Networking*, ISBN 1-58113-422-3, 166-199.
- [4] D. Niculescu; B. Nath.(2003); Ad hoc positioning system (APS) using AOA, *In Proceedings of Twenty-Second Annual Joint Conference of the IEEE Computer and Communications*, ISSN 0743-166X, 1734-1743.
- [5] N. Bulusu; J. Heidemann; D. Estrin(2000); GPS-less low-cost outdoor localization for very small devices, *IEEE Personal Communications*, ISSN 1070-9916, 7:28-34.
- [6] D. Niculescu; B. Nath.(2001); Ad hoc positioning system (APS), *In Proceedings of Global Telecommunications Conference*, ISBN 0-7803-7206-9, 2926-2931.
- [7] R. Nagpal(1999); *Organizing a global coordinate system from local information on an amorphous computer*, MIT Artificial Intelligence Laboratory memo no.1666.
- [8] T. He; C. Huang; B.M. Blum; J.A. Stankovic; T. Abdelzaher(2003); Range-free localization schemes for large scale sensor networks, *In Proceedings of 9th Annual International Conference on Mobile Computing and Networking*, ISBN 1-58113-753-2, 81-95.

-
- [9] Y. Shang; W. Ruml; Y. Zhang; M.P.J. Fromherz(2003); Localization from mere connectivity, *In Proceedings of the 4th ACM international Symposium on Mobile Ad Hoc Networking & Computing*, ISBN 1-58113-684-6, 81-95.
- [10] S. Lederer; Y. Wang; J. Gao(2008); Connectivity-based localization of large scale sensor networks with complex shape, *In Proceedings of the 27th Annual IEEE Conference on Computer Communications*, ISSN 0743-166X, 789-797.
- [11] Y. Wang; S. Lederer; J. Gao(2009); Connectivity-based sensor network localization with incremental delaunay refinement method, *In Proceedings of the 28th Annual IEEE Conference on Computer Communications*, ISSN 0743-166X, 2401-2409.
- [12] M. Jin; S. Xia; H. Wu; X. Gu(2011); Scalable and fully distributed localization with mere connectivity, *In Proceedings of the 30th Annual IEEE Conference on Computer Communications*, ISSN 0743-166X, 3164-3172.
- [13] J.H. Holland(1962); Concerning efficient adaptive systems, *Self-Organizing Systems*, ISSN 1069-0948, 215-230.
- [14] J.H. Holland(1992); *Adaptation in Natural and Artificial Systems*, MA: MIT Press.
- [15] M. Srinivas; L.M. Patnaik(1994); Adaptive probabilities of crossover and mutation in genetic algorithms, *IEEE Transactions on Systems, Man and Cybernetics*, ISSN 0018-9472, 24:656-667.
- [16] K. Deep; M. Thakur(2007); A new crossover operator for real coded genetic algorithms, *Applied Mathematics and Computation*, ISSN 0974-4665, 188:895-911.

Active Queue Management of TCP Flows with Self-scheduled Linear Parameter Varying Controllers

C. Kasnakoglu

Cosku Kasnakoglu

TOBB University of Economics and Technology
Department of Electrical and Electronics Engineering
Ankara, Turkey, 06560
kasnakoglu@etu.edu.tr

Abstract:

Control-theoretic approaches to Active Queue Management (AQM) are typically based on linearizations of fluid flow models around design conditions. These conditions depend on the Round Trip Time (RTT), and the AQM performance is known to degrade if RTT values during actual operation depart substantially from design values. To overcome this difficulty a self-scheduled LPV controller for AQM is considered in this paper, where the controller is modified in real-time based on RTT. Simulations show that the self-scheduled LPV controller has good performance for both constant and time-varying RTTs, and outperforms two other common control-theoretic approaches to AQM.

Keywords: AQM, congestion control, control theory, fluid model, LPV systems, self-scheduled controller

1 Introduction

Congestion is one of the most important problems faced in communications networks. Congestion occurs when a link or node is carrying so much data that its quality of service deteriorates. This results in queueing delay, packet loss or the blocking of new connections, leading to low throughput and eventually to congestion collapse. On the other hand, links carrying less data than a certain level are also not desired as this implies that the link capacities are being underutilized [1]. Modern networks try to avoid these situations using congestion control techniques, among which Active Queue Management (AQM) is of particular interest. AQM operates by dropping or ECN-marking packets before the queue is full, according to a probabilistic rule. Earlier AQM disciplines such as RED [2] and REM [3] required careful tuning of parameters in order to provide good performance, while modern AQM disciplines such as ARED [4] and Blue [5] are self-tuning. With the development of dynamical models such as the fluid flow model in [6], control theoretic approaches for AQM have gained interest, including PI/PID controllers [7, 8] and robust \mathcal{H}_∞ controllers [9–11]. These are based on the linearization of the fluid model, which produces a transfer function from probability of package mark p to queue length q . One difficulty is that the transfer function is valid only for a given Round Trip Time (RTT), and a new transfer function must be obtained for a different RTT. In [12] switching between multiple controllers designed for different RTTs is considered and it is seen that higher number of controllers results in improved performance. However, if the number of controllers is too high, the implementation becomes very difficult and complicated. In addition, since RTT takes values on an interval, designing a controller for each RTT value requires an infinite number of controllers. In this paper a self-scheduled control design for AQM is considered to overcome these difficulties. The controller is parameterized in RTT and achieves stability and small tracking error for both fixed RTT as well as time varying RTT within a prescribed range.

2 Mathematical Model of TCP/AQM

A dynamical model for TCP congestion control was developed in [6] using fluid flow approximation. The dynamical model represents a bottleneck with multiple TCP flows sharing the link. The congestion avoidance is modeled as AIMD (additive-increase multiplicative-decrease). The dynamical model consists of the following nonlinear time-delayed differential equations

$$\dot{W}(t) = \frac{1}{\theta(t)} - \frac{W(t)}{2} \frac{W(t - \theta(t))}{\theta(t - \theta(t))} p(t - \theta(t)) \quad (1)$$

$$\dot{q}(t) = \frac{N(t)}{\theta(t)} W(t) - C(t) \quad \text{for } q(t) > 0 \quad (2)$$

where W is the TCP window size, q is the queue length, N is the number of TCP flows, C is the link capacity, p is the probability of packet mark and θ is the RTT. Let C and N be constants for simplicity. The nonlinear system (1)-(2) can be linearized around an operating point to generate a transfer function from p to q . Let $\delta p := p(t) - p_o$ and $\delta q := q(t) - q_o$ where q_o and p_o are the values at the operating point. Then a transfer function from δp to δq can be obtained as [6]

$$\frac{\delta Q(s)}{\delta P(s)} = \frac{(C^3 \theta^3 N) e^{-\theta s}}{2N^2 \theta^3 C s^2 + (2N^2 C \theta^2 + 4N^3 \theta) s + 4N^3} \quad (3)$$

where θ is regarded as a parameter. One then has a different transfer function for each value of the RTT in the range $[\theta_{\min}, \theta_{\max}]$, where θ_{\min} and θ_{\max} are the minimum and maximum values of RTT for the link.

3 Controller Design

The controller must achieve stability of the closed-loop system and the tracking of a desired queue length, under the presence of constant or time-varying RTT within the specified range.¹ A parameter dependent dynamic controller of the following form will be sought

$$\begin{aligned} \dot{\zeta} &= A_K(\theta)\zeta + B_K(\theta)e \\ \delta p &= C_K(\theta)\zeta + D_K(\theta)e \end{aligned} \quad (4)$$

where ζ is the controller's internal state and $e := \delta q_d - \delta q$ is the error between the desired and the actual queue sizes. A systematic method for the selection of the controller matrices A_K , B_K , C_K and D_K in (4) was derived in [13]. For this purpose the system to be controlled, i.e. (3), must be transformed into affine-parameter dependent form. Approximating the time delay $e^{-\theta s}$ in (3) with a second order Pade approximation

$$e^{-\theta s} \approx \frac{12 - 6\theta s + \theta^2 s^2}{12 + 6\theta s + \theta^2 s^2} \quad (5)$$

and transforming into state space form yields

$$\dot{x} = A_P(\theta)x + B_P(\theta)\delta p \quad (6)$$

$$\delta q = C_P(\theta)x \quad (7)$$

¹Under ideal operation, RTT will be small for small queue size. However no such assumption will be made here, since unexpected changes in link conditions (e.g. increased propagation delay) can create longer RTTs even at smaller queue lengths. Thus, the controller is expected to be prepared for any RTT at any queue size (within the limits allowed).

where

$$A_P(\theta) = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -a_3 & -a_2 & -a_1 & -a_0 \end{bmatrix}, \quad (8)$$

$$B_P(\theta) = [0 \ 0 \ 0 \ 1]^T, \quad (9)$$

$$C_P(\theta) = [b_2 \ b_1 \ b_0 \ 0] \quad (10)$$

and

$$\begin{aligned} a_0 &= 7\theta^{-1} + \frac{2N}{C}\theta^{-2}, \quad a_1 = \frac{14N}{C}\theta^{-3} + 18\theta^{-2}, \\ a_2 &= 12\theta^{-3} + \frac{36N}{C}\theta^{-4}, \quad a_3 = \frac{24N}{C}\theta^{-5}, \\ b_0 &= \frac{C^2}{2N}, \quad b_1 = \frac{-3C^2}{N}\theta^{-1}, \quad b_2 = \frac{6C^2}{N}\theta^{-2}. \end{aligned}$$

Letting $\Theta := \theta^{-1}$, and performing a series expansion around point $\Theta = \Theta_o$, one can write A_P , B_P , C_P in terms of Θ as

$$A_P(\Theta) = A_{P0} + A_{P1}\Theta + \mathcal{O}(\Theta^2) \quad (11)$$

$$B_P(\Theta) = B_{P0} + B_{P1}\Theta + \mathcal{O}(\Theta^2) \quad (12)$$

$$C_P(\Theta) = C_{P0} + C_{P1}\Theta + \mathcal{O}(\Theta^2) \quad (13)$$

where A_{P0} has the same structure as in (8) but with

$$\begin{aligned} a_0 &= \frac{-2N\Theta_o^2}{C}, \quad a_1 = \frac{-28N\Theta_o^3}{C} - 18\Theta_o^2, \\ a_2 &= \frac{-108N\Theta_o^4}{C} - 24\Theta_o^3, \quad a_3 = \frac{-96N\Theta_o^5}{C}, \end{aligned}$$

A_{P1} has the same structure as in (8) but with

$$\begin{aligned} a_0 &= \frac{4N\Theta_o}{C} + 7, \quad a_1 = \frac{42N\Theta_o^2}{C} + 36\Theta_o, \\ a_2 &= \frac{144N\Theta_o^3}{C} + 36\Theta_o^2, \quad a_3 = \frac{120N\Theta_o^4}{C}, \end{aligned}$$

B_{P0} is equal to B_P in (9), B_{P1} is zero, C_{P0} has the same structure as in (10) but with

$$b_0 = \frac{C^2}{2N}, \quad b_1 = 0, \quad b_2 = \frac{-6C^2\Theta_o^2}{N}$$

and C_{P1} has the same structure as in (10) but with

$$b_0 = 0, \quad b_1 = \frac{-3C^2}{N}, \quad b_2 = \frac{12C^2\Theta_o}{N}.$$

Approximating A_P , B_P and C_P with the constant and linear terms in (11)-(13) produces an affine parameter dependent state space system (in Θ) as desired. Self-scheduled controller design methods [13] can then be used to design the parameter dependent controller in (4) to meet the desired control objectives. The performance objective is to have a fast and well-damped response over the entire range of parameter values for step-like references. The control law is also required to achieve robustness by avoiding high-gain feedback at high-frequencies. This will prevent the excitation of high frequency modes and nonlinearities that were unmodeled or neglected.

Consider the feedback control structure depicted in Figure 1. The control problem described

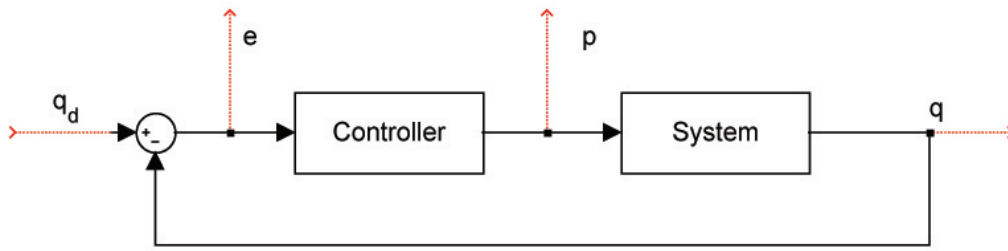


Figure 1: Feedback control structure for AQM

above can be formulated as the minimization of the \mathcal{L}_2 -induced norm of the operator mapping the signal δq_d to signals e and δp . The former map is called the sensitivity function and denoted by S , and the latter map is denoted by KS . Following standard \mathcal{H}_∞ design procedure, the performance objective and robustness objectives are specified through weighting filters $W_1(s)$ and $W_2(s)$, where $W_1(s) = 1/s$ for good tracking of step-like references, and $W_2(s)$ is a third order high pass Butterworth filter. This yields to an \mathcal{H}_∞ optimization problem where the goal is to find a stabilizing controller K for which the inequality

$$\left\| \begin{bmatrix} W_1(s)S(s) \\ W_2(s)K(s)S(s) \end{bmatrix} \right\|_\infty \leq \gamma \quad (14)$$

can be satisfied with γ in \mathbb{R}_+ as small as possible. Let A_{cl} , B_{cl} , C_{cl} and D_{cl} be the state-space matrices corresponding to the one-input two-output system defined by the transfer function matrix $[W_1S, W_2KS]^T$. An LPV controller satisfying (14) can be synthesized and implemented as follows [13]

1. Find a matrix $X_{cl} > 0$, and controller matrices A_{Ki} , B_{Ki} , C_{Ki} , D_{Ki} so that

$$\begin{bmatrix} A_{cl}(\Theta_i)^T X_{cl} + X_{cl} A_{cl}(\Theta_i) & X_{cl} B_{cl}(\Theta_i) & C_{cl}^T(\Theta_i) \\ B_{cl}^T(\Theta_i) X_{cl} & -\gamma I & D_{cl}^T(\Theta_i) \\ C_{cl}(\Theta_i) & D_{cl}(\Theta_i) & -\gamma I \end{bmatrix}$$

is negative-definite for $i = 1, 2$, where $\Theta_1 = \theta_{\min}^{-1}$ and $\Theta_2 = \theta_{\max}^{-1}$.

2. For a given value of $\Theta = \theta^{-1}$, compute the matrices $A_K(\Theta)$, $B_K(\Theta)$, $C_K(\Theta)$ and $D_K(\Theta)$ defining the LPV controller as

$$\begin{bmatrix} A_K(\Theta) & B_K(\Theta) \\ C_K(\Theta) & D_K(\Theta) \end{bmatrix} = \sum_{i=1}^2 \alpha_i(\Theta) \begin{bmatrix} A_{Ki} & B_{Ki} \\ C_{Ki} & D_{Ki} \end{bmatrix}$$

where (α_1, α_2) is a convex decomposition of Θ such that $\Theta = \alpha_1 \Theta_1 + \alpha_2 \Theta_2$ and $\alpha_1 + \alpha_2 = 1$.

4 Simulation Results

The performance of the closed-loop system with the LPV controller is tested using MATLAB/Simulink and discrete event simulations.² The number of TCP flows is taken to be $N = 150$, the link capacity $C = 500$ packets/sec, the desired queue size $q_o = 150$ packets, and the buffer

²The results from both were consistent so only those from the former are shown for better presentation and to save space.

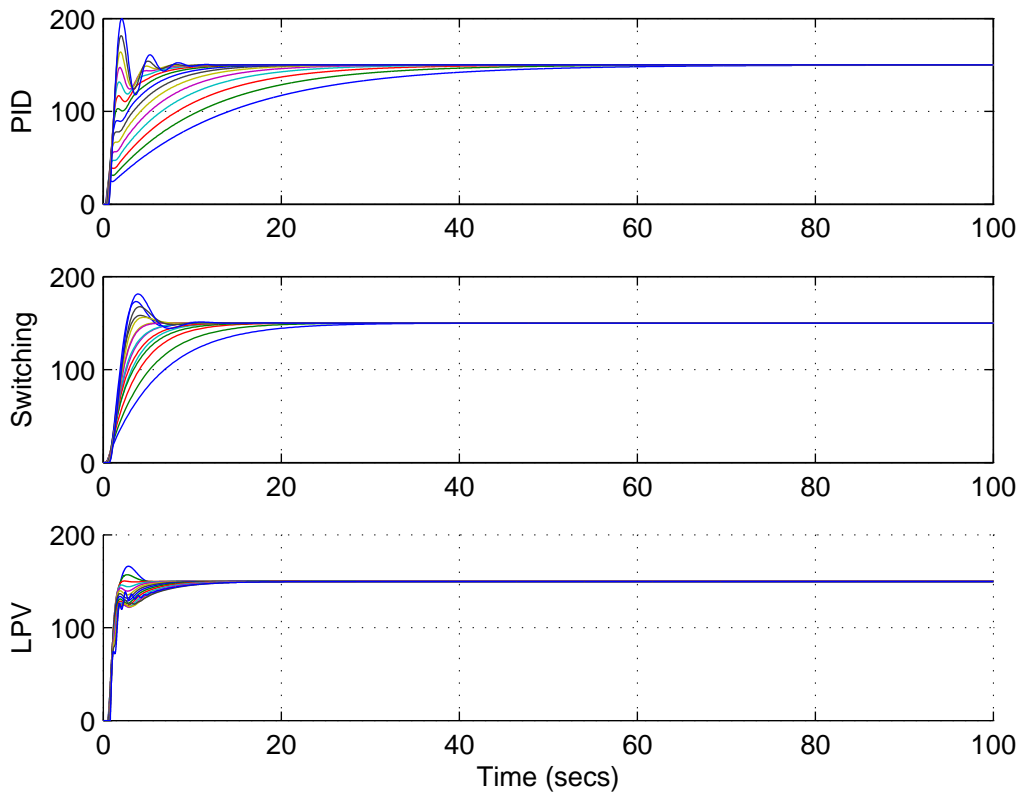


Figure 2: Simulation results for AQM based on PID (top), switching \mathcal{H}_∞ (middle) and LPV controllers (bottom), for 15 fixed values of RTT.

size $q_{\max} = 200$ packets. The RTT is assumed to take values between 0.3 secs and 0.7 secs. The controller design as outlined in Section 3 is carried out with the help of MATLAB Robust Control Toolbox.

Figure 2 shows the simulation results for 15 fixed values of RTT, linearly equally spaced between 0.3 and 0.7. Two other common control-theoretic approaches to AQM design are also implemented for comparison: a PID controller, and a switched \mathcal{H}_∞ controller. After subsequent trials, the best results that could be obtained with the PID controller were those based on a linearization around $\text{RTT} = 0.6$. The switched \mathcal{H}_∞ controller is based on two operating points RTT_1 and RTT_2 , the best results for which were obtained when $\text{RTT}_1 = 0.4320$ and $\text{RTT}_2 = 0.6080$. The simulations illustrate that all controllers eventually succeed in achieving and maintaining the desired queue size of 150 packets. However, the LPV controller response is faster, with less overshoot and better damping over the entire parameter range.

The controllers were also tested for the case when RTT is time varying. The PID controller and the switched controller did not produce a stable closed-loop so these responses are not shown. The result for the LPV controller is shown in Figure 3. The top subfigure shows the variation in RTT versus time, which is a sinusoid with increasing frequency taking values between 0.3 and 0.7.³ The bottom subfigure shows the response of the closed-loop system. It can be observed that the LPV controller produces a stable closed-loop and is capable of maintaining the queue size very close to the desired value, even when RTT is time varying.

³Different cases for the variation of RTT were also tested with successful results, including square, triangle and sawtooth waveforms, as well as the case where RTT varies randomly according to various probability densities. However, only one case is shown in the paper to save space.

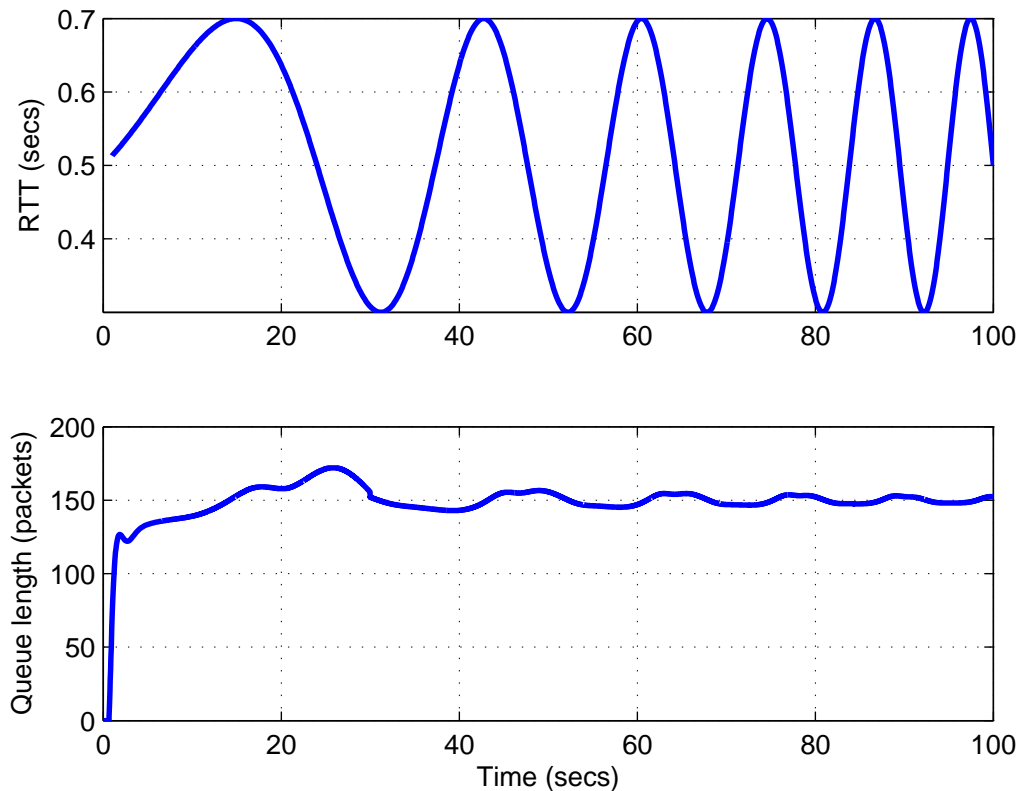


Figure 3: Simulation results for AQM based on LPV, for a time-varying RTT.

5 Conclusions

In this paper self-scheduled LPV control design was implemented for AQM. The controller design is based on a fluid flow approximation of TCP congestion control and utilizes RTT as the scheduling parameter. Simulations were carried out to evaluate the closed-loop system in its ability to keep the queue length at a desired level. Two other control-theoretic approaches to AQM, namely PID control and switching \mathcal{H}_∞ control, were also implemented for comparison. When RTT is constant, it was seen that the LPV controller outperforms the other two, producing a faster and better damped response with lesser overshoot over the entire parameter range. When RTT is time varying, the LPV controller is still capable of maintaining the desired queue size, whereas the PID and switching controllers do not produce a stable closed-loop.

Bibliography

- [1] S. Floyd and K. Fall. Promoting the use of end-to-end congestion control in the Internet. *IEEE/ACM Transactions on Networking*, 7(4):458–472, 1999.
- [2] S. Floyd and V. Jacobson. Random early detection gateways for congestion avoidance. *IEEE/ACM Transactions on Networking*, 1(4):397–413, 1993.
- [3] S. Athuraliya, SH Low, and VH Li. REM: active queue management. *IEEE Network*, 15(3):48–53, 2001.

-
- [4] S. Floyd, R. Gummadi, and S. Shenker. Adaptive RED: An Algorithm for Increasing the Robustness of RED's Active Queue Management. *Preprint*, <http://www.icir.org/floyd/papers.html>, August, 2001.
- [5] W. Feng, K.G. Shin, D.D. Kandlur, and D. Saha. The BLUE active queue management algorithms. *IEEE/ACM Transactions on Networking*, 10(4):513–528, 2002.
- [6] V. Misra, W.B. Gong, and D. Towsley. Fluid-based analysis of a network of AQM routers supporting TCP flows with an application to RED. In *Proceedings of the conference on Applications, Technologies, Architectures, and Protocols for Computer Communication*, pages 151–160. ACM Press New York, NY, USA, 2000.
- [7] CV Hollot, V. Misra, and D. Towsley. Analysis and design of controllers for AQM routers supporting TCPflows. *IEEE Transactions on Automatic Control*, 47(6):945–959, 2002.
- [8] D. Ustebay, H. Ozbay, and N. Gundes. A New PI and PID Control Design Method for Integrating Systems with Time Delays: Applications to AQM of TCP Flows. *Wseas Transactions On Systems And Control*, 2(2):117, 2007.
- [9] P. Yan and H. Ozbay. Robust Controller Design for AQM and \mathcal{H}_∞ -Performance Analysis. *Advances In Communication Control Networks*, 6:7, 2005.
- [10] L. Yu, M. Ma, W. Hu, Z. Shi, and Y. Shu. Design of parameter tunable robust controller for active queue management based on H-infinity control theory. *Journal of network and computer applications*, 34(2):750–764, 2010.
- [11] R. Vilanova and V.M. Alfaro. Robust 2-DoF PID control for Congestion control of TCP/IP Networks. *Int J Comput Commun*, ISSN 1841-9836, 5(5):968–975, 2010.
- [12] D. Ustebay and H. Ozbay. Switching Resilient PI Controllers for Active Queue Management of TCP Flows. In *IEEE International Conference on Networking, Sensing and Control*,, pages 574–578, 2007.
- [13] P. Apkarian, P. Gahinet, and G. Becker. Self-scheduled H-infinity Control of Linear Parameter-varying Systems: a Design Example. *Automatica*, 31(9):1251–1261, 1995.

An Improved Genetic Algorithm for the Multi Level Uncapacitated Facility Location Problem

V. Korać, J. Kratica, A. Savić

Vanja Korać*, Jozef Kratica

Mathematical Institute, Serbian Academy of Sciences and Arts
Kneza Mihaila 36/III, 11 000 Belgrade, Serbia
vanja@mi.sanu.ac.rs, jkratica@mi.sanu.ac.rs

*Corresponding author: vanja@mi.sanu.ac.rs

Aleksandar Savić

University of Belgrade, Faculty of Mathematics
Studentski trg 16, 11000 Belgrade, Serbia
aleks3rd@gmail.com

Abstract: In this paper, an improved genetic algorithm (GA) for solving the multi-level uncapacitated facility location problem (MLUFLP) is presented. First improvement is achieved by better implementation of dynamic programming, which speeds up the running time of the overall GA implementation. Second improvement is hybridization of the genetic algorithm with the fast local search procedure designed specially for MLUFLP. The experiments were carried out on instances proposed in the literature which are modified standard single level facility location problem instances. Improved genetic algorithm reaches all known optimal and the best solutions from literature, but in much shorter time. Hybridization with local search improves several best-known solutions for large-scale MLUFLP instances, in cases when the optimal is not known. Overall running time of both proposed GA methods is significantly shorter compared to previous GA approach.

Keywords: evolutionary approach, metaheuristics, discrete location, combinatorial optimization.

1 Introduction

During the last decades, there was an expansive growth in the ways of solving facility location problems. Research has been concentrated mostly on location problems, which require minimization of total travel time, physical distance, or some other related cost. In most cases it is assumed that facilities are large enough to meet any demand. So far, there are several deterministic uncapacitated models proposed in the literature.

The multi-level uncapacitated facility location problem (MLUFLP) is NP-hard, since it is a generalization of the uncapacitated facility location problem which is proven to be NP-hard in [7].

The multi-level version of uncapacitated facility problem is discussed only in several papers ([1–5, 8, 9]), compared to several hundreds of papers about the basic version of the problem. Moreover, most of the papers that deal with MLUFLP were of theoretical nature, without experimental results. The only exceptions are papers [5] and [8].

In [5] four methods for solving MLUFLP are implemented. Three algorithms are based on linear programming relaxation of the model, described in Section 2, which has an enormous number of variables and moderate number of constrains. Consequently, these three algorithms are capable of solving only small size MLUFLP instances with up to 52 potential facility locations. Running times of these methods are greater than 100 seconds. Gap values are rather satisfying

but the corresponding running times are very long even for these small size instances. Obviously, all three linear programming based methods are unable to solve larger size problem instances. The fourth method produce very quick results (up to 0.07 seconds), but the gaps are very large - up to 732%. From these facts, it can be concluded that none of these algorithms is capable of solving real medium-size and large-scale MLUFLP instances.

The only method capable of reaching optimal solutions in solving small and medium-size MLUFLP instances was an evolutionary approach presented in [8]. It also gives results on large-scale MLUFLP instances in a reasonable running time. A binary encoding scheme with appropriate objective function containing dynamic programming approach for solving subproblem was used, which consisted of finding sequence of located facilities on each level to satisfy clients' demands. Used dynamic programming approach has the polynomial number of states and steps, since the given subproblem is a special case of the shortest path problem. That approach enables the GA to use the standard genetic operators and caching technique, which helps genetic algorithm to reach promising search regions.

2 Previous work

Let us first give a mathematical formulation of the MLUFLP problem as presented in [5]. The input data to the MLUFLP consists of a set of facilities F ($|F| = m$) partitioned into k levels, denoted F_1, \dots, F_k , a set of clients D ($|D| = n$), a fixed cost f_i for establishing facility $i \in F$, and a metric that defines transportation costs c_{ij} for each $i, j \in F \cup D$. A feasible solution assigns each client a sequence of k facilities, one from each level F_k, \dots, F_1 , respectively. A feasible solution is charged the sum of the fixed costs of the facilities used, plus the transportation costs of the clients' assignments. Each client's transportation cost is the sum of transportation cost from itself to the first facility of its sequence, plus the transportation cost between successive facilities. An optimal solution is a feasible solution with a minimum total cost.

The assignment of a client $j \in D$ to a valid sequence of facilities can now be represented as the assignment of j to a path p from j to one of the top level facilities F_1 . The set of all valid sequences of facilities is defined with $P = F_k \times \dots \times F_1$ and the transportation cost of client j 's assignment to sequence $p = (i_k, \dots, i_1)$ by $c_{pj} = c_{ji_k} + c_{i_k i_{k-1}} + \dots + c_{i_2 i_1}$. Variable y_i has value of 1 if facility i is established and 0 otherwise. Similarly, variable x_{pj} represents whether or not a client j is assigned to the path p . Using the notation mentioned above, the problem can be written as:

$$\min \sum_{i \in F} f_i y_i + \sum_{p \in P} \sum_{j \in D} c_{pj} x_{pj} \quad (1)$$

$$\sum_{p \in P} x_{pj} = 1, \quad \text{for each } j \in D, \quad (2)$$

$$\sum_{p \ni i} x_{pj} \leq y_i, \quad \text{for each } i \in F, j \in D, \quad (3)$$

$$x_{pj} \in \{0, 1\}, \quad \text{for each } p \in P, j \in D, \quad (4)$$

$$y_i \in \{0, 1\}, \quad \text{for each } i \in F. \quad (5)$$

Table 1: Fixed costs

Facilities	f1	f2	f3	f4	f5	f6
Fixed cost	70	50	30	20	20	40

Table 2: Distance between facilities on level 1 and level 2

	f1	f2
f3	58	23
f4	44	74
f5	67	15
f6	29	38

The objective function (1) minimizes the sum of overall transportation cost and fixed costs for establishing facilities. Constraint (2) ensures that every client is assigned to a path while constraint (3) guarantees that any facility on a path used by some client is paid for. Constraints (4) and (5) reflect binary nature of variables x_{pj} and y_i .

The problem can be illustrated with one small example.

Example 1 An example of the MLUFLP is shown below. It assumes $k = 2$ levels of $m = 6$ facilities: the first level contains 2 potential facilities and the second 4 potential facilities. in this example there are $n = 5$ clients to be served. The fixed costs of establishing facilities are given in Table 1, the distances between facilities of different levels and the distances between clients and facilities on the second level are given in Table 2 and Table 3, respectively.

The total enumeration technique, described in Section 4, is used to obtain an optimal solution. Established facilities are: f2 on the first and f3, f5 on the second level. The objective function value is 329. The sequences of facilities for each client are shown in Table 4.

3 Improved GA method

Proposed genetic algorithm proposed in this paper is an improved version of GA used in [8], so after the short summary of the overall genetic algorithm implementation we only give the detailed description of the improvement parts. The outline of our GA implementation is given below, where N_{pop} denotes the overall number of individuals in the population, N_{elite} is a number of elite individuals and ind and obj_{ind} mark the individual and its value of objective function.

```

Input_Data();
Population_Init();
while not Stopping_Criterion() do
    for  $ind := (N_{elite} + 1)$  to  $N_{pop}$  do
    
```

Table 3: Distance between clients and facilities

	f3	f4	f5	f6
client 1	38	13	57	41
client 2	47	52	63	15
client 3	42	48	13	54
client 4	9	54	41	43
client 5	15	18	36	22

Table 4: Sequences of facilities for clients

	level 2	level 1
client 1	f3	f2
client 2	f3	f2
client 3	f5	f2
client 4	f3	f2
client 5	f3	f2

```

if (Exist_in_Cache(ind)) then
  objind:= Get_Value_From_Cache(ind);
else
  objind:= Objective_Function(ind);
  Local_Search(ind,objind);
  Put_Into_the_Cache_Memory(ind,objind);
  if (Full_Cache_Memory()) then
    Remove_LRU_Block_From_Cache_Memory();
  endif
endif
endfor
Fitness_Function();
Selection();
Crossover();
Mutation();
endwhile
Output_Data();

```

The encoding of the individuals used in this implementation is binary. The set of facilities F which will be used is naturally represented as the individual represented with a binary string of length m . Digit 1 at the i -th place of the string denotes that $y_i = 1$, while 0 shows the opposite ($y_i = 0$).

Now, for any individual, through its genetic code, string of established facilities is given. Objective function value can now be computed using dynamic programming explained in [8] if the array of established facilities has at least one established facility. Otherwise, solution is not valid and individual is marked as infeasible.

Genetic operators are the same as in [8]. We can briefly summarize them

- Selection operator is fine-grained tournament selection - FGTS [6],
- Crossover operator is a standard one-point crossover operator,
- Mutation operator is a standard simple mutation modified for frozen genes,
- Initialization of first population is random.

The elitist strategy is applied to N_{elite} elite individuals, which are directly passing to the next generation. The genetic operators are applied to the rest of the population ($N_{nnel} = N_{pop} - N_{elite}$ non-elite individuals).

All duplicates of every individual are eliminated. Individuals with the same objective function value, but different genetic codes are limited with the number of their appearance - N_{rv} .

The `Fitness_Function()`, the fitness f_{ind} of individual ind $1, 2, \dots, N_{pop}$ is computed by scaling values of objective function obj_{ind} of all individuals into the interval $[0, 1]$, so that the best individual ind_{min} has fitness 1 and the worst one ind_{max} has fitness 0.

Caching technique applied in this improved version of GA is the same as described in [8].

3.1 Improved objective function

In [8] for feasible individual ind , the `Objective_Function(ind)`, is evaluated in four steps.

1. In the first step, the values of variables y_i are obtained from the genetic code. Let us denote the number of established facilities (number of y_i 's with value 1) with m_1 .
2. In the second step, the array of minimal costs cs is initialized. The array cs carries information about total minimal costs for serving clients and established facilities (except the ones on the first level), regarding the costs for serving facilities on the upper level. The minimal cost values in cs , for the established facilities on the first level are initially set to zero, while the costs for established facilities on remaining levels and clients are set to a large constant $INF = 10^{30}$.
3. The minimal costs for each client and each facility are calculated by dynamic programming. For each facility in each level (except the first one), the array of total minimal costs (initialized in the second step) is updated. The minimal cost value for each facility is the minimum of the sum of the minimal cost for established facility on the upper level and transportation cost between the two facilities. The same procedure is done for each client: Among the established facilities from the last level, the one with the minimal sum of the corresponding minimal cost value and the transportation cost facility-client, is taken.
4. Finally, the objective value is computed by adding all transportation costs facility-client and fixed costs for all established facilities.

As can be seen from the previous algorithm, some of the operation is unnecessary. For example, all pairs of facilities from the two consecutive levels are taken into consideration where as only pairs of established facilities are what matters. In Example 1. there is an optimal solution only for 1 out of 2 facilities established on the first level and only for 2 out 4 facilities established on the second level. Therefore, for calculating minimal cost between the first and the second level, previous algorithm performs $2*4=8$ operations (one operation per facility pair), while as only $1*2=2$ is really needed. A similar situation arises in calculating facility-client costs, which in previous algorithm as $5*4=20$ operation, while only $2*5=10$ is really needed.

In order to correct shortcomings of previous algorithm, the following improvements are proposed. New array of established facility indices is constructed. When calculating minimal costs with dynamic programming approach, instead of using all pairs of facilities on consecutive levels we used only pairs of established facilities obtained from the newly constructed array. Similarly, when calculating minimal facility-client costs, , instead of using of all facilities on the last level, only the established once are taken into consideration.

Therefore, usage of this procedure effectively improves objective function calculation which makes the major part in running time of the overall GA implementation. In presented example we have a significant improvement, but for large scale instances improvement can be even greater.

This can be seen on finding objective function for the largest instance `mt1_5L_60_120_250_500_1070.2000` Best known solutions has 1 established facility on the first 4 levels and 5 established facilities on the 5th level.

Now, instead of calculating

1. 60*120 pairs of facilities between level 1 and 2, there is only one pair of facilities;
2. 120*250 pairs of facilities between level 2 and 3, there is only one pair of facilities;
3. 250*500 pairs of facilities between level 3 and 4, there is only one pair of facilities;
4. 500*1070 pairs of facilities between level 4 and 5, there are five pairs of facilities;
5. 1070*2000 pairs of facilities-clients we have 5*2000 pairs of facilities-clients.

Although objective function calculation has speeded up factor greater than 200, speed up of overall GA implementation is about 6.5 times. This discrepancy can be explained that the running time of calculating objective function was conducted previously, instead of representing the major part, becomes dominated by running time of other parts of GA implementation. It is obvious that this way of calculating objective function is much faster so its running time is significantly shorter than those of genetic operators (selection, crossover and mutation).

3.2 Local search

In order to improve the accuracy of the solutions additionally, the proposed GA approach incorporates a local search procedure. It considers the array of facilities and regards which facilities are established and, which are not. In the array facilities which are established they are marked with 1 and those that are not, are marked with 0. Local search starts from the first facility in the array and tries to change its status. If there is an improvement in the value of objective function, local search starts from the beginning of the array, otherwise it moves to the next facility. If the facility that is momentarily in consideration is the only established facility on that level, local search moves to the next facility in the array. This procedure is repeating until there are no more facilities in the array.

Improvement can be determined in two ways, according to the status of the regarded facility.

1) Let us first consider the case where the facility is non-established, and its status is changed into being established. Then new value of objective function is determined in the following way. All established facilities on the previous levels are unchanged with their fixed transport costs. On the current level costs are added with the fact that running facility is now established. All other transport costs are unchanged. On the next level already found transport costs are compared with transport costs through the new established facilities and changed if the latter is smaller. On the levels following this all transport costs must be evaluated anew. This includes transport costs from the last level of facilities to the clients.

2) In the second case an established facility becomes non-established. Transport costs for that level are then reduced for transport cost of that facility. The array is not reduced in length but all paths through that facility are out of consideration. On the next level transport costs are again calculated only for those that were connected through facility considered. On the following levels, transport costs are determined only for those that were on paths leading through the facility considered.

Local search is applied only on the best individual in population and this only if it stayed unchanged in N_{rep} generations (in this work $N_{rep} = 5$). On the one individual, local search is applied only once, because of its deterministic nature. If the best individual is replaced with another individual with equal value of an objective function but different genetic code and that individual stays unchanged through N_{rep} generations, local search is again applied. In this paper, local search was not applied in the first 50 generations, whatsoever.

Table 5: GA results on the instances with previously known optimal solution

Instance name	<i>Optimal solution</i>	<i>GA</i>		<i>ImprovedGA</i>		<i>GA + LS</i>	
		sol	t_{tot}	sol	t_{tot}	sol	t_{tot}
cap71_2L_6_10.50	1813375.51	<i>opt</i>	0.202	<i>opt</i>	0.195	<i>opt</i>	0.226
cap71_3L_2_5_9.50	4703216.31	<i>opt</i>	0.22	<i>opt</i>	0.195	<i>opt</i>	0.233
cap101_2L_8_17	1581551.39	<i>opt</i>	0.274	<i>opt</i>	0.224	<i>opt</i>	0.269
cap101_3L_3_7_15.50	3227179.81	<i>opt</i>	0.258	<i>opt</i>	0.223	<i>opt</i>	0.265
cap131_2L_13_37.50	1592548.45	<i>opt</i>	0.557	<i>opt</i>	0.358	<i>opt</i>	0.407
cap131_3L_6_14_30.50	3201970.46	<i>opt</i>	0.546	<i>opt</i>	0.355	<i>opt</i>	0.396
cap131_4L_3_7_15_25.50	3630297.67	<i>opt</i>	0.515	<i>opt</i>	0.32	<i>opt</i>	0.352

4 Computational results

All tests were carried out on an Intel 2.5 GHz with 2 GB memory. The algorithms were coded in C programming language.

The GA is tested on the same instances as in [8] to show improvement obtained by improved objective function and hybridization with local search. In this paper, only those instances with multiple levels of facilities are taken into consideration, since for the basic problem with only one level several hundreds of papers are available. Therefore, instances with multiple levels are the main interest of this research, while results on instances with only one level of facilities are omitted. For small size instances, optimal values are known from the literature which is indicated in Table 1.

For fair and direct comparison of the results, we leave same GA parameters as in [8]. The finishing criterion of GA is the maximal number of generations $N_{gen} = 5000$. The algorithm also stops if the best individual or best objective value remains unchanged through $N_{rep} = 2000$ successive generations. Since the results of GA are nondeterministic, the GA was run 20 times on each problem instance.

In [8] experimental results are not totally in accordance with the instances. It was possible to repeat testing and about 1/3 of instances values of objective function cannot be validated as presented in that quoted paper. For example, value of objective function, for instance *cap131_4l_3_7_15_25.50*, is 3630297.67 in the paper and in repeated testing, but for instance *capa_3l_15_30_55.100* value function in the paper is 40725103.254, and 25424361.91 in the repeated testing. Furthermore, for some instances, results obtained in literature and in repeated testing slightly differ, but that can be caused by slightly different random seed. For example, for instance *ms1_4l_64_128_256_552.1000* values of objective function in the literature and in repeated testing are 30936.585 and 31257.27 respectively which means that result in the quoted paper is slightly better. For instance *mt1_4l_120_250_520_1110.2000* respective results are 65044.003 and 64995.454 which means that result in repeated testing is slightly better. Because of these differences, results of repeated testing will be presented in the whole paper as results obtained by original GA implementation.

Table 1 presents the GA result on smaller and medium instances. In the first column names of instances are given. The instance's name carries information about the number of levels, the number of facilities on each level and the number of clients respectively. For example, the instance *capb_3l_12_25_63.1000* is created by modifying ORLIB instance *capb*, which has 3 levels with 12, 25, 63 facilities, respectively and 1000 clients.

The second column contains the optimal solution on the current instance, if it is previously known, otherwise sign as $-$. The best GA value GA_{best} and running time t_{tot} of the original GA is given in the following two columns, with marked optimum in cases when GA reached

Table 6: GA results on the instances with unknown optimal solution

Instance name	<i>GA</i>		<i>Improved GA</i>		<i>GA + LS</i>	
	sol	t_{tot}	sol	t_{tot}	sol	t_{tot}
capa_2L_30_70.1000	14829245.63	32.715	14829245.63	16.589	14829245.63	17.517
capa_3L_15_30_55.1000	25424361.91	17.76	25424361.91	8.679	25424361.91	8.933
capa_4L_6_12_24_58.1000	35421258.15	16.231	35421258.15	7.49	35421258.15	7.636
capb_2L_35_65.1000	14479223.79	27.062	14479223.79	14.118	14479223.79	14.283
capb_3L_12_25_63.1000	25986997.29	28.16	25986997.29	12.147	25986997.29	12.477
capb_4L_6_13_31_50.1000	41787432.24	17.371	41787432.24	8.631	41787432.24	9.083
capc_2L_32_68.1000	14072575.52	29.437	14072575.52	16.218	14072575.52	16.901
capc_3L_13_27_60.1000	26751918.75	26.85	26751918.75	12.932	26751918.75	12.773
capc_4L_4_9_27_60.1000	47109818.66	24.491	47109818.66	11.316	47109818.66	11.184
mq1_2L_100_200.300	8341.287	25.561	8341.287	3.091	8341.287	3.180
mq1_3L_30_80_190.300	12994.871	23.95	12994.871	2.903	12994.871	2.951
mq1_4L_18_39_81_162.300	18048.0305	21.233	18048.0305	3.002	18048.0305	3.000
mq1_4L_20_40_80_160.300	17648.0095	20.759	17648.0095	3.141	17648.0095	3.070
mr1_2L_150_350.500	6733.815	88.939	6733.815	8.979	6733.815	10.277
mr1_2L_160_340.500	6707.505	88.878	6707.505	8.752	6707.505	10.096
mr1_3L_55_120_325.500	10911.319	80.546	10911.319	8.39	10911.319	9.401
mr1_4L_30_65_140_265.500	15237.2605	76.126	15237.2605	7.852	15237.2605	8.297
ms1_2L_320_680.1000	13361.3895	479.724	13361.3895	66.767	13361.3895	98.085
ms1_3L_120_250_630.1000	21923.331	364.135	21881.384	64.965	21881.384	89.360
ms1_4L_64_128_256_552.1000	31257.27	385.73	30902.742	54.715	30902.742	79.729
ms1_5L_25_55_120_250_550.1000	40494.7435	373.781	40249.2415	57.169	40094.335	80.764
mt1_2L_650_1350.2000	27733.057	2410.92	27733.057	457.427	27733.057	685.593
mt1_3L_255_520_1225.2000	46278.719	2398.943	46529.979	426.48	46278.719	622.719
mt1_3L_256_600_1144.2000	46095.09	2403.649	46095.09	406.855	46095.09	585.418
mt1_4L_120_250_520_1110.2000	64995.454	2318.705	64995.454	402.143	64851.16	567.867
mt1_5L_60_120_250_500_1070.2000	83486.9185	2265.78	83363.586	371.867	83363.586	543.857

the optimal solution. In the following two columns are given best values and running times, *Improved GA* and t_{tot} of the improved GA algorithm. Finally, in the last two columns best values and running times are given, *GA + LS* and t_{tot} of the hybridized version of GA algorithm.

As it can be seen from Table 5, there are three variants of GA reached optimal solutions. As expected, running times of the improved GA has been better than for original GA for all instances, and improvement is up to 40 percent. Hybridized GA with LS in some instances also has better running times (4 of 7) than the original GA. It can be noticed that better running times in both variants are accomplished on instances with the greater number of levels and greater number of facilities.

In Table 6 results on larger instances are given. In this case, optimal solutions are not known so that the column is omitted and all other columns have same meaning as in Table 1.

From Table 6 can be concluded that improvements in both new variants of GA were validated. Improved GA has all running times faster than original GA (in some cases like the instances with 500 clients up to 10 times faster) without losing on quality of results. Only in one instance, was result was obtained by improved GA, greater than that obtained by original GA (instance mt1_3L_255_520_1225.2000). In all other instances improved GA has better results. Hybridization with local search also produced improvements. Results obtained from this variant are the best in all instances. This variant of GA produced better results than original GA in 5 instances (mostly in the largest instances), and was better than the improved GA in 3 instances. An interesting fact is that GA with LS has running times shorter almost in all instances (except the smallest instances), and sometimes like in those instances with 300 clients GA with LS run

up to 8 times faster than the original GA. Even in the largest instances hybridized version ran up to 44 times faster and also produced better results. As it could be expected, running times of the hybridized version were always longer than those of Improved GA.

5 Conclusions and Future Works

This paper presented improved genetic algorithm for solving multi-level uncapacitated facility location problem. Improvements are achieved, firstly, in formulating the new version of finding objective function for the problem which resulted in much shorter running times, and secondly, through hybridization with fast and reliable local search procedure, which achieved better results and in almost all cases run faster than original GA approach. These improvements are considerable, since hybrid GA is obtained 5 new best-known solutions and running times are shorter sometimes up to one order of the magnitude compared to original GA version. It can be concluded that both of these versions represent considerable improvements in solving MLUFLP.

Future research can be directed to parallelization of presented GA, hybridization with some other heuristics and its application in solving similar facility location problems.

Bibliography

- [1] K. Aardal, F. Chudak, D.B. Shmoys, A 3-approximation algorithm for the k-level uncapacitated facility location problem, *Information Processing Letters*, 72:161–167, 1999.
- [2] A. Ageev, Improved approximation algorithms for multilevel facility location problems, *Operations Research Letters*, 30: 327–332, 2002.
- [3] A. Ageev, Y. Ye, J. Zhang, Improved combinatorial approximation algorithms for the k-level facility location problem, *SIAM Journal on Discrete Mathematics*, 18: 207–217, 2005.
- [4] A.F. Bumb, W. Kern, A Simple Dual Ascent Algorithm for the Multilevel Facility Location Problem, *Proceedings of the 4th International Workshop on Approximation Algorithms for Combinatorial Optimization Problems and 5th International Workshop on Randomization and Approximation Techniques in Computer Science: Approximation, Randomization and Combinatorial Optimization*, 55-62, August 18–20, 2001.
- [5] N.J. Edwards, Approximation algorithms for the multi-level facility location problem. *Ph.D. Thesis, Cornell University*, 2001.
- [6] V. Filipović, J. Kratica, D. Tošić, D. I. Ljubić, Fine Grained Tournament Selection for the Simple Plant Location Problem. in Proceedings of the 5th Online World Conference on Soft Computing Methods in Industrial Applications - WSC5, 152–158, September 2000.
- [7] J. Krarup, P. M. Pruzan, The simple plant location problem: Survey and synthesis, *European Journal of Operational Research*, 12: 36–81, 1983.
- [8] M. Marić, An efficient genetic algorithm for solving the multi-level uncapacitated facility location problem, *Computing and Informatics*, 29(2): 183–201, 2010.
- [9] J. Zhang, Approximating the two-level facility location problem via a quasi-greedy approach. *Mathematical Programming*, 108: 159–176, 2006.

A Improved EPC Class 1 Gen 2 Protocol with FCFS Feature in the Mobile RFID Systems

X. Li, Y. Quan

Xiaowu Li

1. School of Information Science and Technology
Southwest Jiaotong University
Chengdu 610031, Sichuan, China
NO.111 of the North Second Ring Road
lxw66@126.com

2. School of Information and Technology
Kunming University
Puxin Road 2, Kunming 650214, China

Quanyuan Feng*

School of Information Science and Technology
Southwest Jiaotong University
Chengdu 610031, Sichuan, China
NO.111 of the North Second Ring Road
*Corresponding author: fengquanyuan@163.com

Abstract: In all anti-collision protocols of RFID standards, EPC Global Class 1 Generation 2 (C1G2) protocol has been most widely used in RFID systems since it is simply, efficient and safety. Similar to most existing anti-collision protocols, The C1G2 protocol initially aims at tag identification of static scenarios, where all tags keep still during the tag identification process. However, in many real scenarios, tags generally move along a fixed path in the reader coverage area, which implies that tags stay the coverage area only for a limited time (sojourn time). The scenarios are usually called mobile RFID systems. Because the multiple tag identification based on a shared wireless channel is random, tags entering the reader coverage area earlier may be identified later (random later identification phenomenon). The phenomenon and the limited sojourn time may let some tags lost. In this paper, we propose an improved C1G2 protocol with first come first served feature in mobile RFID systems. The protocol can overcome the RLI phenomenon effectively and retains good initial qualities of C1G2 protocol by modifying it slightly. Simulation results show that the proposed protocol can significantly reduce the numbers of lost tags in mobile RFID systems. The idea of the paper is beneficial for redesigning other existing tag anti-collision protocols so as to make these protocols adapt to mobile RFID systems.

Keywords: RFID, tag anti-collision, mobile RFID systems, EPC C1G2.

1 Introduction

Radio Frequency Identification (RFID) based on wireless radio communication is increasingly used a great deal in many ways. A reader can only interrogate tags within its interrogation region, where the reader's electromagnetic signal is strong enough to energize the tags. This region is also referred to as the reader coverage area. During tag identification process, a tag collision occurs when multiple tags reply simultaneously to a reader. So far, so many tag collision resolutions have been proposed for static scenarios where no tag enters or leaves the reader coverage area during the process of tag identification [1-11].

However, in other many practical scenarios, tags which are usually attached to items and moving along the fixed path in the reader coverage area need to be identified as soon as possible [12-17]. These scenarios are generally called mobile RFID systems and often appear at doorways

of warehouse and highway [13,14]. In these scenarios, tags stay the coverage area only for a time (sojourn time). Since the reader has only limited sojourn time to identify the passing tags, tag loss may sometimes be inevitable. Figure 1 shows the mobile RFID system diagram and illustrates some parameters: the reader coverage area length, the tag moving velocity (assumed constant) and tag sojourn time S (given in slots) which is a quotient between the reader coverage area length and the time interval for one slot [17]. Mobile RFID systems shown in Figure 1 are our research scenarios.

In the mobile RFID systems, some tags may leave the coverage area unidentified under the condition of high tag tag density or is high tag moving speed, etc. We refer to such a tag as a lost tag. How to decrease the number of lost tags is the critical research issue in mobile RFID applications. Therefore, we rate tag loss ratio (TLR) as a critical performance index, which is defined as the quotient between the number of lost tags and the total number of tags entering the reader coverage area [12–15, 17].

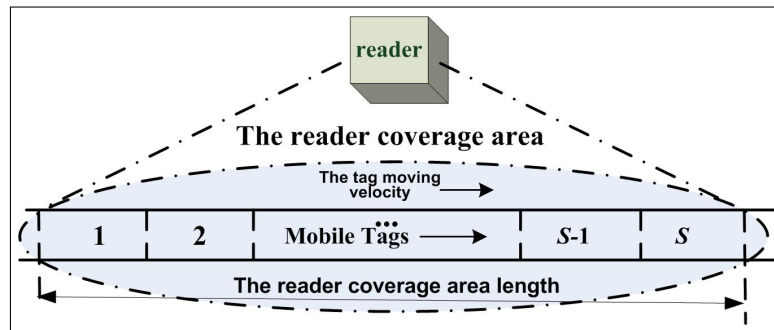


Figure 1: The mobile RFID system diagram.

Besides, since the multiple tag identification based on a shared wireless channel is random, tags entering the coverage area earlier may be identified later, which is named the random later identification (RLI) phenomenon [17]. An example given in Figure 2 illustrates the phenomenon. In this figure, all tags in the coverage area contend the shared wireless channel together under the condition of high tag density. Each tag randomly selects a slot to communicate with the reader. For example, the tag with asterisk, which enters the coverage area earlier than other ones, is closer to the exit of the coverage area than other tags, but selects the latter slot (slot 12) for communication with the reader. So, it may leave the coverage area unidentified because its slot number (12) is relatively greater. According to the analysis, we believe that the RLI phenomenon and limited sojourn time are two main reasons for causing lost tags, especially under high tag density.

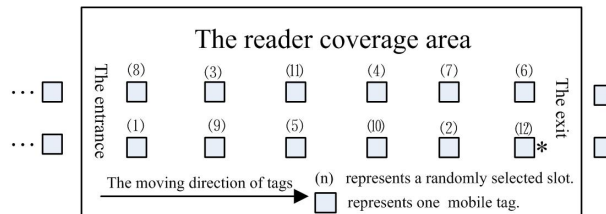


Figure 2: Random later identification (RLI) phenomenon .

In this paper, we will propose an improved EPC C1G2 protocol with first come first served feature in mobile RFID systems, which is an improvement to EPC C1G2 protocol [10]. The improved protocol overcomes RLI problem and converts random access for tags into sequential

access for tag groups. For the sake of brevity, the proposed protocol is named IC1G2MRS protocol.

The remainder of this paper is organized as follows. In Section 2 we briefly review the related work in the area. In Section 3 we propose the principle of tag sequencing. Then, in Section 4 we offer IC1G2MRS protocol. Section 5 provides the simulation results. Finally, section 6 concludes.

2 Related Works

EPC C1G2 protocol is also known as Q algorithm or C1G2 algorithm. In the protocol [10], each tag randomly selects a frame slot and sends a 16-bit random number ($RN16$) to the reader for reserving the remainder of the frame slot at the selected frame slot. There are three kinds of replies (correspond to three kinds of slots): (1) No reply (empty slot): after waiting for a short time, if the reader has not received a $RN16$. The reader terminates the current frame slot. (2) Collision reply (collision slot): if multiple tags transmit $RN16$ in the frame slot, the reader terminates the current frame slot. (3) Success reply (success slot): if only one $RN16$ is sent to the reader in the current frame slot, the reader can successfully get the tag's ID in the rest of the frame slot. One benefit of EPC C1G2 protocol is that empty frame slots and collision frame slots are shorter than success frame slots. Therefore, EPC C1G2 protocol is superior to the FSA very much [8]. Figure 3 shows the tag identification process in EPC C1G2 protocol.

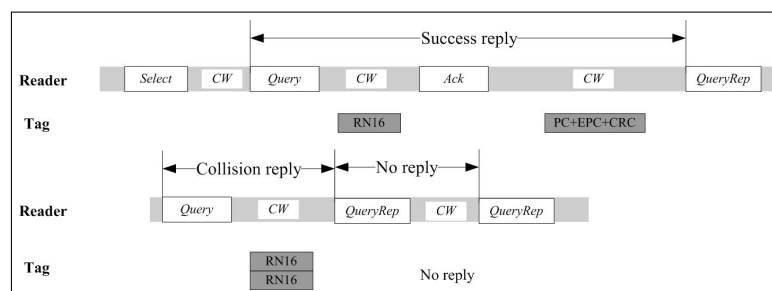


Figure 3: Random later identification phenomenon (RLI).

Until now, some researches on mobile RFID systems have been done [12–14, 17]. Authors in [12, 13] pay attention to TLR computation of frame slotted ALOHA protocol (FSA) and CSMA by Markov model or dynamic systems model.

In [14, 15], authors focused on single tag set passing the reader coverage area in a limited time where no more tag sets enter the coverage area until the previous one has left. Obviously, the scenario is different from our research scenario.

In [16], Sarangan offered a framework which reduces the tag reading time by using bitmaps. The framework can improve to some extent the system performance of mobile RFID systems. However, the method also can not solve RLI problem effectively.

In [17], authors present a grouping based dynamic framed slotted aloha for tag anti-collision protocol in the mobile RFID systems. The protocol has FCFS feature but does not possess the attributes of high system performance of EPC C1G2 protocol in static RFID scenarios.

In this paper, we develop an improved C1G2 protocol with first come first served (FCFS) feature in mobile RFID systems (IC1G2MRS). Some advanced characteristics of the protocol are as follows: (1) IC1G2MRS protocol can be easily implemented without complicated computation. (2) The frequency of Tag sequencing for FCFS mechanism is adjustable. (3) The proposed protocol can retain initial good quantities of C1G2 protocol. (4) The protocol can avoid overflow error effectively when the RFID systems work continuously for an awful long time.

3 The Principle of Tag Sequencing

It is well known that FCFS is based on the sequenced objects. So, tag sequencing is a critical work. In the following, we will give the detailed explanation of tag sequencing, which can refer to [17].

Usually, tags enter the reader coverage area continuously in mobile RFID systems. In this case, the tags in the coverage area can be grouped in their arrival order at regular intervals (assume M slots). Such obtained groups are called time groups and each time group can be assigned a time group sequence number ($TGSN$). We can derive that tags arriving in the same time interval possess the same $TGSN$ and tags arriving in different time interval possess different $TGSNs$. For the special case of static scenarios, all tags' $TGSNs$ are the same because they enter the coverage area in the same time interval. Notice that M can be called the frequency parameter of tag sequencing. The lower M is, the higher the frequency of tag sequencing is.

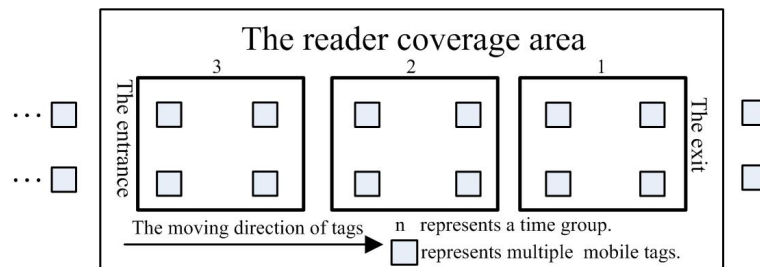


Figure 4: The tags grouped in their arrival time order.

To record a $TGSN$, each tag should possess a memory cell to store $TGSN$. Figure 4 shows that the tags in the coverage area are sequenced to become three time groups. Initial $TGSN$ of new tags are set as $TGSN=null$, which means that the new tags have entered the coverage area but have not received any $TGSNs$ given by the reader. We can derive that if a tag's $TGSN$ equals $null$, the tag has not been sequenced. Notice that new tags may enter the reader's field continuously, which can enlarge the $TGSN$ extremely because every new time group may lead to $TGSN$ plus one. To save the tag's memory space, a circular queue mechanism is used to manage the $TGSN$. Moreover, for some mobile scenarios where RFID systems are required to work continuously for an awful long time, such as several months or years, only using circular queue mechanism to manage $TGSN$ can avoid overflow of $TGSN$. For example, we assume that the length of $TGSN$ is 8 bits. With circular queue mechanism, next time group's $TGSN$ will be set as $TGSN=0$ if ongoing time group's $TGSN==2^8-1$ during the process of tag sequencing. This means that the reader coverage area can hold a maximum of 2^8-1 time groups at any time. The reason for subtracting 1 is that the circular queue has overflowed when there are 2^8-1 time groups in the coverage area. However, with non-circular queue mechanism, the occurrence of next time group will result in an overflow error if ongoing time group's $TGSN==2^8-1$ [17].

We also know that a circular queue mechanism requires the aid of parameters $HEAD$ and $REAR$. In our method, the two parameters are stored in the reader and respectively point to the earliest arrival time group's $TGSN$ and the latest arrival time group's $TGSN$ plus one. Based on the explanation, we can also say that $HEAD$ points the tags in the earliest time group. Figure 5 depicts that a circular queue which is used to manage the $TGSNs$. The circular queue mechanism with 8-bit $TGSN$ can use $TGSN 0 \sim TGSN 2^8-1$ repeatedly. In Figure 5 (a), $HEAD = 1$ and $REAR = 3$ mean that there are two time groups in the coverage area, and time groups 1 and 2 are the earliest and latest time groups respectively. When $HEAD == REAR$, there is no unidentified tag in the reader coverage area, as shown in Figure 5 (b).

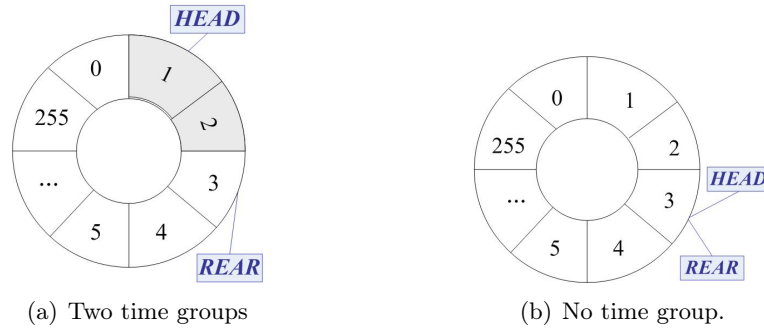


Figure 5: Circular queue mechanism is used to manage the *TGSNs*.

In summary, the tag sequencing is summarized as follows: Upon hearing Sequencing command sent by reader, new tags set their *TGSN* as $TGSN = REAR$ and then the reader set its *REAR* as $REAR = (REAR + 1) \bmod 2^8$ where *TGSN* length equals 8. In this way, tag sequencing can be implemented, that is, tags in the coverage area can be divided into multiple time groups.

4 The Proposed IC1G2MRS Protocol

IC1G2MRS is based on C1G2 algorithm. Compared with C1G2 algorithm, IC1G2MRS has two other features: (1) IC1G2MRS can sequence tags in the reader coverage area to be tag groups (time groups) every M slots (e.g. 20 slots). That is, the reader can group these tags in their arrival order. (2) On this basis, the reader can identify tag groups one by one in order of their *TGSNs* which indicates that IC1G2MRS has a first come first served (FCFS) character. IC1G2MRS protocol is summarized as follows:

Step 1. The reader sequences tags in the coverage area by sending Sequencing command. Upon hearing this command, all tags whose *TGSN* is equal to 'null' will be assigned a true *TGSN* ($0 \sim 255$). More details can refer to section 3. Following operations only aim at the tags in the earliest time group if no special explanation is given.

Step 2. The reader sends a query command (e.g. Query or QueryRep) to tags. If it is currently the time to start a new frame, the reader issues Query command; Otherwise the reader issues QueryRep.

Step 3. Tags receive the query command from the readers. It could be Query or QueryRep. If the received query command is Query, every unidentified tag selects a slot number (SN) between 0 and 2^Q-1 , inclusively. If the received query command is QueryRep, all unidentified tags decrease their SN by 1. After these operations, every unidentified tag randomly selects a slot number between 0 and 2^Q-1 , inclusively. The tags whose SNs equal 0 generate a 16-bit random number (*RN16*) and responds the *RN16* to the reader. Because each tag conducts the operations independently, there are three possible outcomes (correspond to three kinds of slots):

1) Success reply (success slot): Only one tag responds and the reader receives the *RN16* successfully. Then the reader will send out an acknowledgement (ACK) (go to Step 4).

2) Collision reply (collision slot): More than one tag responds simultaneously and collision occurs. Then the reader adjusts Q as $Q = \min(15, Q + c)$ and the reader continues to identify tags by sending QueryRep to tags (go to Step 2).

3) No reply (empty slot): No tag responds. Then the reader updates Q as $Q = \max(0, Q - c)$ and continues to identify tags by sending QueryRep to tags (go to Step 2).

Notice that, a success slot is counted as a standard slot, named slot for short. According to EPC C1G2 standard, a collision and an empty slot respectively equal approximately 0.2 times

and 0.1 times the time cost of standard slots [10]. More details can refer to section 3.

Step 4. After the reader successfully receives a $RN16$ from a tag, it sends an ACK back to all tags. But only the tag that successfully replies in Step 3 recognizes the ACK and then goes to Step 5. The ACK contains the $RN16$ that the reader received in Step 3.

Step 5. All tags may receive the ACK, but only the tag that replies successfully in Step 3 continues to transmit its EPC to the reader. Then the reader continues to identify remaining tags in ongoing earliest time group by sending QueryRep (go to Step 2). Notice that identified tags are silenced. Namely, they will not reply the reader's command in following identification process.

Step 6. The reader sequences tags in the reader coverage area every M (e.g. 20) slots.

Step 7. After one frame ends, if one or more collisions occur in the current frame, the reader sends Query command with updated Q to all unidentified tags in the earliest time group after ongoing frame (go to Step 2). If not, the ongoing earliest time group move next time group. Then go to step 1.

In short, IC1G2MRS protocol is hybrid of C1G2 protocol and first come and first served (FCFS) mechanism. It can group tags in the reader coverage area to be multiple tag groups in their arrival order and then use basic principles of EPC C1G2 protocol to identify each time group in FCFS order.

5 Simulation and Results

Before we evaluate the proposed protocol, we first offer a simplified mobile RFID experiment model because there are no good mathematical methods which can compute the TLR of various anti-collision protocols in the mobile RFID systems so far [17]. More reasons for using the simplified mobile RFID experiment model can refer to [17]. The model is composed of 3 groups of mobile tags, as shown in Figure 6 where $T1$ denotes the time interval between the 1st and 2nd groups of tags, $T2$ denotes the time interval between the 2nd and 3rd groups of tags, $N1$, $N2$ and $N3$ denote the number of tags in the 1st, 2nd and 3rd groups of tags respectively, S denotes the tag sojourn time. Notice that in the paper, all parameters related to time, such as $T1$, $T2$, S and M are measured in slot.

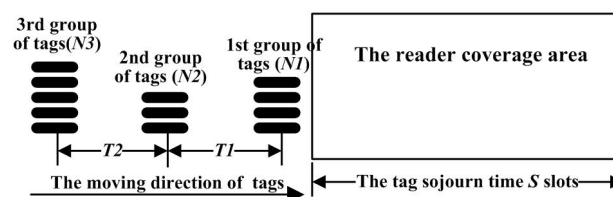


Figure 6: Simplified mobile RFID experiment model.

Now, we comprehensively evaluate IC1G2MRS by comparing it with the typical C1G2 protocol [10]. Our simulation based on Monte Carlo technique and the simplified mobile RFID experiment model. The basic parameters of following simulation experiments are $N1=70$, $N2=70$, $N3=70$, $T1=50$, $T2=80$, $S=200$, $M=20$. The evaluation results are shown in Fig. 6.

Figure 7 (a) depicts the relationship between TLR and the tag density by changing the number of tags $N2$ in the second group. TLR of two protocols increases as the tag density increases. The reason is that the reader has to identify more tags in the same interval. Compared with C1G2, IC1G2MRS has better performance. For example, when the number of tags in the second group $N2$ is less than or equal to 100 tags, TLR of IC1G2MRS is nearly zero while that of C1G2 is 10 percent.

Figure 7 (b) depicts the relationship between TLR and the tag sojourn time S by changing S (correspond to change of the tag moving speed or the reader coverage area length). TLR of two protocols decreases as S increases. The reason is that the reader has more time to identify tags as S increases. Compared with C1G2, IC1G2MRS has better performance. For example, when tag sojourn time S is larger than or equals 200 slots, TLR of IC1G2MRS is nearly zero while TLR of C1G2 protocol is 5 percent.

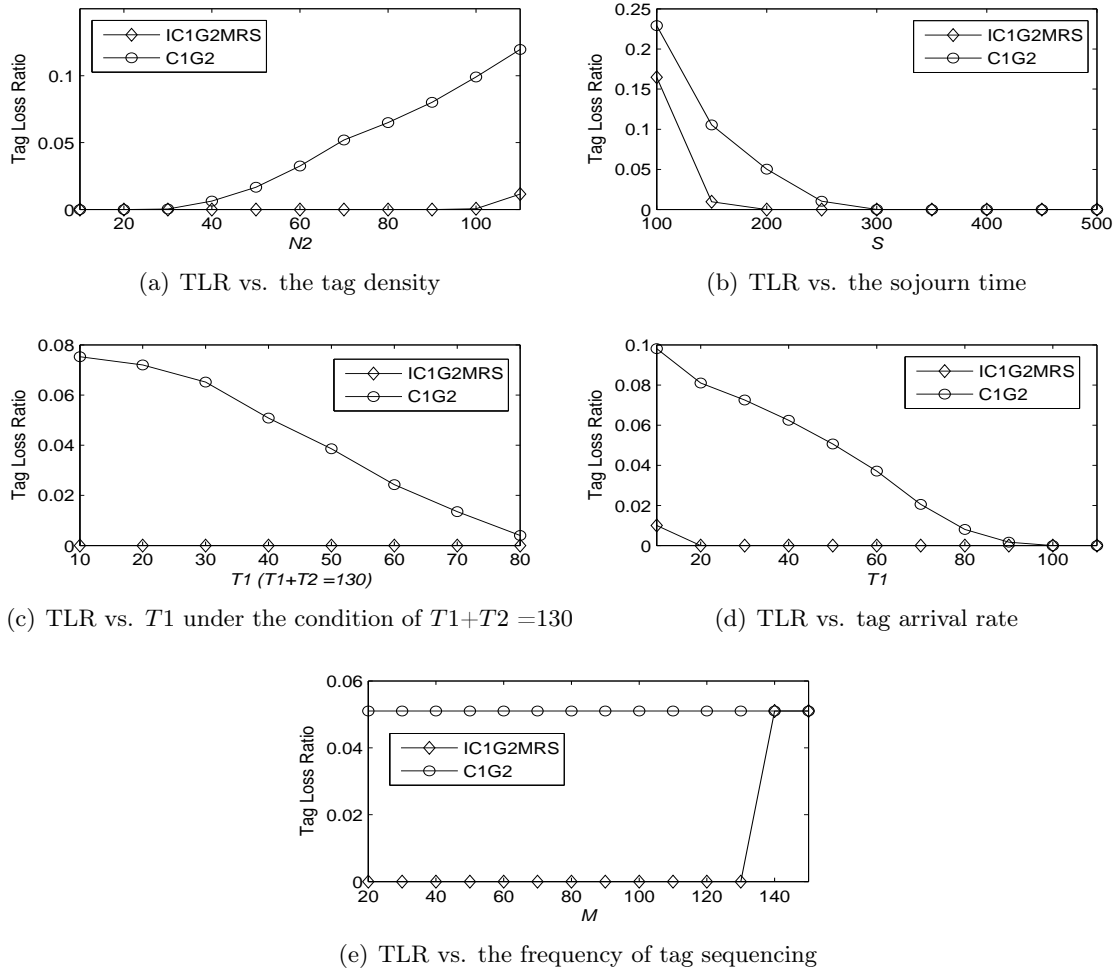


Figure 7: TLR of the proposed IC1G2MRS protocol and C1G2 protocol.

Figure 7 (c) depicts the relationship between TLR and time interval T_1 under the condition that $T_1 + T_2$ equals a constant (e.g. 130), which can to some extent conclude which protocol is sensitive to changing T_1 in the same time span, and which protocol is not. For the experiment, IC1G2MRS can maintain a stable performance whereas C1G2 is unsteady.

Figure 7 (d) depicts the relationship between TLR and tag arrival rate by changing T_1 . From the figure, we can find that TLR of two protocols decreases as the tag arrival rate become slow. The reason is that the reader has more time to identify tags when the rate becomes slow. Compared with C1G2, IC1G2MRS has better performance. For example, when T_1 is equal to 20 slots, TLR of IC1G2MRS is nearly zero percent while TLR of C1G2 protocol is 8 percent.

Figure 7 (e) depicts the relationship between TLR and the frequency of tag sequencing by changing M . From the figure, TLR of IC1G2MRS increases as the frequency of tag sequencing decreases (correspond to increase of M). The reason is that RLI problem may not be resolved

when the frequency of tag sequencing is too low. For example, when $M > 130$, the three groups of tags in the Figure 6 are not sequenced and are seen as a group of tags. Thus TLR of IC1G2MRS significantly increases and is equal to that of C1G2. We can also know easily that if the frequency of tag grouping is too high, great TLR may occur. The reason is that the operation of the tag grouping also needs the time cost. From the experiments, we can derive that the frequency of tag sequencing has a strong influence on TLR.

6 Conclusions

EPC C1G2 protocol is the most popular standard in RFID systems. But it can not be used in mobile RFID systems effectively because the random late identification (RLI) phenomenon can not be overcome. The paper specifies the random late identification (RLI) phenomenon which causes unnecessary tag loss in the mobile RFID systems. Then IC1G2MRS protocol is proposed to resolve the phenomenon and retains initial good quantities of C1G2 protocol, which improves the system performance of mobile RFID systems dramatically.

Acknowledgement

This work is supported by the National Natural Science Foundation of China (NNSF) under Grant 60990320, 60990323, 61271090, and the National 863 Project of China under Grant 2012AA012305, and Sichuan Provincial Science and technology support Project under Grant 2012GZ0101, and Chengdu Science and technology support Project under Grant 12DXYB347JH-002.

Bibliography

- [1] He, M. et al (2011); A fast RFID Tag Identification Algorithm Based on Counter and Stack, *Expert Systems with Applications*, ISSN 0957-4174, 38: 6829-6838.
- [2] Yeh, M. et al. (2009); Adaptive Splitting and Pre-signaling for RFID Tag Anti-collision, *Computer Communications*, ISSN 0140-3664, 32: 1862-1870.
- [3] Finkenzeller, K. (2002); RFID Handbook: Radio-frequency Identification Fundamentals and Applications, *John Wiley Press*.
- [4] Chen, Y. (2013); Multiple-Bits-Slot Reservation Aloha Protocol for Tag Identification, *IEEE Transactions on Consumer Electronics*, ISSN 0098-3063, 59(1): 93-10.
- [5] Alcaraz, J. et al. (2013); A Stochastic Shortest Path Model to Minimize the Reading Time in DFSA-Based RFID Systems, *IEEE Communications Letters*, ISSN 1089-7798, 17(2): 341-344.
- [6] Wong, C. et al. (2007); Grouping Based Bit-slot ALOHA Protocol for Tag Anti-collision in RFID Systems, *IEEE Communication Letters*, ISSN 1089-7798, 11(12): 946-948.
- [7] Jia, X. et al. (2010); An Efficient Anti-collision Protocol for RFID Tag Identification, *IEEE Communication Letters*, ISSN 1089-7798, 14(11): 1014-1016.
- [8] Zhu, L.; Yum, T. (2011); A Critical Survey and Analysis of RFID Anti-Collision Mechanisms, *IEEE Communications Magazine*, ISSN 0163-6804, 5: 214-221.

- [9] Klair, D. et al. (2009); A Survey and Tutorial of RFID Anti-Collision Protocols, *IEEE Communications Surveys and Tutorials*, ISSN 1553-877X, 12(3): 400-421.
- [10] EPCglobal Specification for RFID Air Interface; Radio-frequency identity protocols class-1 generation-2 UHF RFID protocol for communications at 860 MHz - 960 MHz, version 1.0.9, 2005.
- [11] Vogt, H. (2002); Efficient Object Identification with Passive RFID Tag, *Proceedings of 2002 IEEE international conference on systems*, pp. 98-113
- [12] Alcaraz, J. et al (2011); Dynamic System Model for Optimal onfiguration of Mobile RFID Systems, *Computer Networks*, ISSN 1389-1286, 55: 74-83.
- [13] Vales-Alonso, J. et al (2009); Markovian model for Computation of Tag Loss Ratio in Dynamic RFID Systems, *Proceedings of 5th European Workshop on RFID Systems and Technologies, Bremen, Germany*, pp. 16-17.
- [14] Vales-Alonso, J. et al (2011); On the Optimal Identification of Tag Sets in Time-constrained RFID Configurations, *Sensors*, ISSN 1424-8220, 11: 2946-2960.
- [15] Xie, L. (2010); Efficient Tag Identification in Mobile RFID Systems. , *Proceedings of IEEE International Conference INFOCOM*, pp. 15-19.
- [16] Sarangan, V. (2008); A framework for fast RFID tag reading in static and mobile environments, *Computer Networks*, ISSN 1389-1286, 52(5): 1058-1073.
- [17] Li, X.; Feng, Q. (2013); Grouping Based Dynamic Framed Slotted ALOHA for Tag Anti-Collision Protocol in the Mobile RFID Systems, *Appl. Math. Inf. Sci.*, ISSN 1935-0090, (2L): 655-659.

Direct Method for Stability Analysis of Fractional Delay Systems

M.A. Pakzad, M.A. Nekoui

Mohammad Ali Pakzad*

Department of Electrical Engineering
Science and Research Branch, Islamic Azad University
Tehran, Iran

*Corresponding author: m.pakzad@srbiau.ac.ir

Mohammad Ali Nekoui

Faculty of Electrical and Computer Engineering
K.N.Toosi University of Technology, Seyed-Khandan
P.O. Box 16315-1355, Tehran, Iran
manekoui@eetd.kntu.ac.ir

Abstract: In this paper, a direct method is presented to analyze the stability of fractional order systems with single and multiple commensurate time delays. Using the approach presented in this study, first, without using any approximation, the transcendental characteristic equation is converted to an algebraic one with some specific crossing points. Then, an expression in terms of system parameters and imaginary root of the characteristic equation is derived for computing the delay margin. Finally, the concept of stability is expressed as a function of delay. An illustrative example is presented to confirm the proposed method results.

Keywords: Control, fractional delay systems, stability windows, Root-Locus.

1 Introduction

Time delay is an inherent part of many dynamical and physical systems. The device delay often appears in computer based control systems, the wireless and web-based control systems, communication systems and so on. Also, the existence of various sensors and actuators in the feedback loop is the cause of delay in many systems. Moreover, taking an automatic computer-based control system or a process control with networked transmission for example, it is necessary and more reasonable to simultaneously consider the possible transmission delay for the corresponding control law and the device delay due to computing or processing. Recently, much attention has been paid to the subjects of stability, stabilization and control of the time delay systems [1] - [5]. Fractional order models would be more accurate than integer order models. In fact, real world processes generally or most likely are fractional order systems.

In a generic sense, the fractional-order systems are identified by non-integer powers of the Laplace variable s . When time delays and fractional-order derivative are involved in dynamical systems, we have fractional-delay systems. The stability of fractional-delay systems can be determined by the root location of the characteristic equations. The most natural way to find the location of the roots of a linear fractional order system is to solve its corresponding characteristic equation. But in our case this will be a transcendental one, being thus generally impossible to solve it directly. For this reason, most of the existing approaches study stability of such systems by finding the crossings of poles through the imaginary axis [8]. There has been a large effort to deal with this problem, as can be seen by the large quantity of articles dealing with it for the standard case (integer order systems); see [5], and others.

In [6], the necessary and sufficient conditions for the BIBO stability of fractional order delay systems have been introduced. From the numerical analysis point of view, the effective numerical algorithms have been discussed in [7] and [8] for the evaluation of BIBO stability of fractional

order delay systems. In [9] a heavy computation scheme based on the Cauchy's integral has been proposed to test the stability of such systems, and in [10], a technique based on the Lambert W function was used for the same purpose. In this work, an analytical direct approach to determine all possible stability regions in the parametric space of delay is proposed. The original idea in this strategy is derived based on the method reported in [5] to achieve the stability criteria of integer order delay systems. The proposed method is an analytically elegant procedure that first converts the transcendental characteristic equation into an algebraic equation without the transcendentality by eliminating the highest degree of commensurability terms successively. The resulting algebraic equation without the transcendentality also enables us to easily determine the delay dependency of the system stability and the sensitivities of crossing roots (root tendency) with respect to the time delay.

2 Preliminaries and Definitions

A general class of fractional order LTI (linear time invariant) systems with multiple commensurate delays of retarded type is taken into account:

$$D_t^{1/\alpha} x(t) = \frac{d^{1/\alpha} x(t)}{dt^{1/\alpha}} = A_0 + \sum_{\ell=1}^n A_\ell x(t - \ell\tau)$$

The fractional delay characteristic equation of above system can be expressed in the general form of

$$C(\sqrt[\alpha]{s}, \tau) = \det \left(s^{1/\alpha} I - A_0 - \sum_{\ell=1}^n A_\ell e^{-\ell\tau s} \right) = \sum_{\ell=0}^n P_\ell(\sqrt[\alpha]{s}) e^{-\ell\tau s} = 0 \quad (1)$$

where parameter τ is non-negative, such that $\tau \in \mathbb{R}^+$ and $P_\ell(\sqrt[\alpha]{s})$ for $\ell \in \mathbb{N}_N$ is a real polynomial in the complex variable $\sqrt[\alpha]{s}$ (where $\alpha \in \mathbb{N}$). Note that the zeros of characteristic equation (1) are in fact the poles of the system under investigation. We find out from [6] that the transfer function of a system with a characteristic equation in the form of (1) will be H_∞ stable if, and only if, it doesn't have any pole at $\Re(s) \geq 0$. For fractional order systems, if an auxiliary variable of $v = \sqrt[\alpha]{s}$ is used, a practical test for the evaluation of stability can be obtained. By applying this auxiliary variable in characteristic equation (1), the following relation is obtained:

$$C_v(v, \tau) = \sum_{\ell=0}^n P_\ell(v) e^{-\ell\tau v^\alpha} \quad (2)$$

This will transform the domain of the system from a multisheeted Riemann surface into the complex plane, where the poles can be easier calculated. In this new variable, the instability region of the original system is not given by the right half-plane, but in fact by the region described as:

$$|\angle v| \leq \frac{\pi}{2\alpha} \quad (3)$$

with $v \in \mathbb{C}$. Let us assume that $s = \pm j\omega$ or in other words, $s = \omega e^{\pm j\pi/2}$ are the roots of characteristic equation (1) for a $\tau \in \mathbb{R}^+$. Then for the auxiliary variable, the roots are defined as follows:

$$v = \sqrt[\alpha]{s} = \sqrt[\alpha]{\omega} e^{\pm j\pi/2\alpha} \quad (4)$$

3 Crossing position

The main objective of this section is to present a new method for the evaluation of stability and determination of the unstable roots of a fractional delay system. The proposed method eliminates the transcendental term of the characteristic equation without using any approximation and converts it into a equation without the transcendentality such that its real roots coincide with the imaginary roots of the characteristic equation exactly. If the characteristic equation (1) has a solution of $s = j\omega_c$ then $C(\sqrt[\alpha]{-s}, \tau) = 0$ will have the same solution.

$$C(\sqrt[\alpha]{-s}, \tau) = \sum_{\ell=0}^n P_{\ell}(\sqrt[\alpha]{-s}) e^{\ell\tau s} = 0 \quad (5)$$

Characteristic equation (5) can be written in terms of the auxiliary parameter \bar{v} as:

$$C(\bar{v}, \tau) = \sum_{\ell=0}^n P_{\ell}(\bar{v}) e^{\ell\tau v^{\alpha}} = 0 \quad (6)$$

A recursive procedure should be developed to achieve that purpose. Therefore, let us define

$$C^{(1)}(v, \tau) = \sum_{\ell=0}^{n-1} [P_0(\bar{v})P_{\ell}(v) - P_n(v)P_{n-\ell}(\bar{v})] e^{-\ell\tau v^{\alpha}} = \sum_{\ell=0}^{n-1} P_{\ell}^{(1)}(v) e^{-\ell\tau v^{\alpha}} = 0 \quad (7)$$

Then, we have

$$C^{(1)}(\bar{v}, \tau) = \sum_{\ell=0}^{n-1} [P_0(v)P_{\ell}(\bar{v}) - P_n(\bar{v})P_{n-\ell}(v)] e^{\ell\tau v^{\alpha}} = \sum_{\ell=0}^{n-1} P_{\ell}^{(1)}(\bar{v}) e^{\ell\tau v^{\alpha}} = 0 \quad (8)$$

Where

$$P_{\ell}^{(1)}(v) = P_0(\bar{v})P_{\ell}(v) - P_n(v)P_{n-\ell}(\bar{v}) \quad (9)$$

We can easily repeat this procedure to eliminate commensuracy terms successively by defining a new polynomial

$$P_{\ell}^{(r+1)}(v) = P_0^{(r)}(\bar{v})P_{\ell}^{(r)}(v) - P_{n-r}^{(r)}(v)P_{n-r-\ell}^{(r)}(\bar{v}) \quad (10)$$

and an augmented characteristic equation

$$C^{(r)}(v, \tau) = \sum_{\ell=0}^{n-r} P_{\ell}^{(r)}(v) e^{-\ell\tau v^{\alpha}} = 0 \quad (11)$$

By repeating this procedure n times, we eliminate the highest degree of commensuracy terms and obtain the following augmented characteristic equation

$$C^{(n)}(v) = P_0^{(n)}(v) = 0 \quad (12)$$

Where

$$P_0^{(n)}(v) = P_0^{(n-1)}(\bar{v})P_0^{(n-1)}(v) - P_1^{(n-1)}(v)P_1^{(n-1)}(\bar{v}) \quad (13)$$

It should be emphasized that that if $s = j\omega_c$ is the solution of (1) for some τ , then it is also a solution of (12). If we substitute $\bar{v} = (e^{-j\pi/\alpha})v$ and $v = \sqrt[\alpha]{\omega_c} e^{j\pi/2\alpha}$ in (13), we get the following equation in ω

$$D(\omega) = \left(P_0^{(n-1)}(e^{-j\pi/\alpha} v) P_0^{(n-1)}(v) - P_1^{(n-1)}(v) P_1^{(n-1)}(e^{-j\pi/\alpha} v) \right) \Big|_{v=\sqrt[\alpha]{\omega_c} e^{j\pi/2\alpha}} \quad (14)$$

The corresponding value of time delay is then computed by

$$\tau^* = \frac{1}{\omega_c} \tan^{-1} \left(\frac{\Im[(P_0^{(n-1)}(v))/(P_1^{(n-1)}(v))]}{\Re[-(P_0^{(n-1)}(v))/(P_1^{(n-1)}(v))]} \right) \Big|_{v=\sqrt[n]{\omega_c} e^{\frac{j\pi}{2\alpha}}} + \frac{2k\pi}{\omega_c} \quad ; k \in \mathbb{Z}^+ \tag{15}$$

Theorem 1. A system with characteristic equation (1) has finite crossing points for any $\tau \in \mathbb{R}^+$.

Proof: Let assume $s = j\omega_c$ be a pair of roots for $C(\sqrt[n]{s}, \tau)$ then $v = \sqrt[n]{s} = \sqrt[n]{\omega_c} e^{\pm j\pi/2\alpha}$ would be a pair of roots for $D(\omega)$. since $D(\omega)$ is a finite degree polynomial with maximum degree of $(P_0^{(n)}(v))$ then the number of crossing points of (1) that are the real roots of $D(\omega)$ is finite. \square

The whole ω values, for which $s = j\omega$ is a root of equation (1) for some non-negative delays, is defined as the crossing frequency set.

$$\Omega = \left\{ \omega \in \mathbb{R}^+ \mid C(\sqrt[n]{s}, \tau) = 0, \text{ for some } \tau \in \mathbb{R}^+ \right\} \tag{16}$$

Corollary 2. If the system given as (1) is stable for $\tau = 0$ (i.e. system without delay) and $\Omega = \phi$, then the system will be stable for all positive values of $\tau \in \mathbb{R}^+$.

Proof: From the fact that there are no roots crossing the imaginary axis. \square

4 Direction of crossing

After the crossing points of characteristic equation (1) from the imaginary axis are obtained, the goal now is to determine whether each of these root crossings from the imaginary axis is a stabilizing cross or a destabilizing cross. Assume that (s, τ) is a simple root of $C(\sqrt[n]{s}, \tau) = 0$. The root tendency for each ω_{cm} and τ_{mk} is defined as:

$$\text{Root Tendency} = RT|_{s=j\omega_c} = \text{sgn} \left(\Re \left(S_{\tau}^s \Big|_{\substack{s=j\omega_{cm} \\ \tau=\tau_{mk}}} \right) \right) = \text{sgn} \left(\Re \left(- \frac{\partial C / \partial \tau}{\partial C / \partial s} \Big|_{s=j\omega_c} \right) \right) \tag{17}$$

If it is positive, then it is a destabilizing crossing, whereas if it is negative, this means a stabilizing crossing. Notice that root tendency represents the direction of transition of the roots at $j\omega_{cm}$ as τ increases from $\tau_{mk-\epsilon}$ to $\tau_{mk+\epsilon}$, $0 < \epsilon \ll 1$.

Theorem 3. The root tendency at a crossing, $j\omega_c$ is invariant with respect to time delay τ_{mk} .

Proof: One can find $ds/d\tau$ for simple roots of (1) as follows:

$$\frac{ds}{d\tau} = - \frac{\partial C / \partial \tau}{\partial C / \partial s} = \frac{\sum_{\ell=0}^n P_{\ell}(\sqrt[n]{s}) \ell s e^{-\ell\tau s}}{\sum_{\ell=0}^n \frac{P_{\ell}(\sqrt[n]{s})}{ds} e^{-\ell\tau s} - \sum_{\ell=0}^n P_{\ell}(\sqrt[n]{s}) \ell \tau s e^{-\ell\tau s}} \tag{18}$$

Based on (18) and definition given in (17), the root tendency of each time delay τ_{mk} is written as follows:

$$\begin{aligned} RT|_{s=j\omega_c}^{\tau} &= \text{sgn} \left(\Re \left(S_{\tau}^s \Big|_{s=j\omega_c} \right) \right) = \text{sgn} \left(\Re \left(- \frac{\partial C / \partial \tau}{\partial C / \partial s} \Big|_{s=j\omega_c} \right) \right) \\ &= \text{sgn} \left(\Re \left(\frac{\sum_{\ell=0}^n P_{\ell}(\sqrt[n]{s}) \ell s e^{-\ell\tau s}}{\sum_{\ell=0}^n \frac{P_{\ell}(\sqrt[n]{s})}{ds} e^{-\ell\tau s} - \sum_{\ell=0}^n P_{\ell}(\sqrt[n]{s}) \ell \tau s e^{-\ell\tau s}} \right) \right) = \text{sgn} \left(\Re \left(\frac{\sum_{\ell=0}^n \frac{P_{\ell}(\sqrt[n]{s})}{ds} e^{-\ell\tau s}}{\sum_{\ell=0}^n P_{\ell}(\sqrt[n]{s}) \ell s e^{-\ell\tau s}} - \frac{\tau}{s} \right)^{-1} \right) \tag{19} \\ &= \text{sgn} \left(\Re \left(\frac{\sum_{\ell=0}^n \frac{P_{\ell}(\sqrt[n]{s})}{ds} e^{-\ell\tau s}}{\sum_{\ell=0}^n P_{\ell}(\sqrt[n]{s}) \ell s e^{-\ell\tau s}} \right) \right) = \text{sgn} \left(\Im \left(\frac{\sum_{\ell=0}^n \frac{P_{\ell}(\sqrt[n]{s})}{ds} e^{-\ell\tau s}}{\sum_{\ell=0}^n P_{\ell}(\sqrt[n]{s}) \ell e^{-\ell\tau s}} \right) \right) \Big|_{\substack{s=j\omega_c \\ \tau=\tau_{m1} + \frac{2k\pi}{\omega_c}}} \end{aligned}$$

The root tendency in each time delay τ_{mk} is independent from the time delay itself and constant for each crossing frequency, because $e^{-\ell\tau s}$ and $P_{\ell}(\sqrt[n]{s})$ do not depend on τ_{mk} . \square

5 Illustrative Example

We present an example case, which display all the features discussed in the text.

Example 4. *This example has been taken from [3] and [8]. Consider the following linear time-invariant fractional order system with one delay:*

$$C(\sqrt{s}, \tau) = (\sqrt{s})^3 - 1.5(\sqrt{s})^2 + 4(\sqrt{s}) + 8 - 1.5(\sqrt{s})^2 e^{-\tau s} \quad (20)$$

This system has a pair of poles ($s = \pm 8j$) on the imaginary axis for $\tau = 0$. A very involved calculation scheme based on Cauchy's integral has been used in [9] to show that this system is unstable for $\tau = 0.99$ and stable for $\tau = 1$. Our objective in this example is to find all the stability windows based on the method described in this article for this system. Using auxiliary variable $v = \sqrt{s}$ into (20) the characteristic equation is obtained as:

$$C_v(v, \tau) = v^3 - 1.5v^2 + 4v + 8 - 1.5v^2 e^{-\tau v^2} \quad (21)$$

By applying the criterion expressed in the previous section, we can eliminate exponential term from (21) as follows:

$$jv^6 + 1.5(1 - j)v^5 + 14(1 + j)v^3 - j16v^2 + 32(1 - j)v + 64 = 0 \quad (22)$$

By inserting expression $v = \sqrt{\omega} e^{j\frac{\pi}{4}} = a(1 + j)$ in the above equation, we get:

$$8a^6 - 12a^5 - 56a^3 + 32a^2 + 64a + 64 = 0 \quad (23)$$

Where $a = \sqrt{\omega/2}$. The real solutions of (23) for a is

$$a = 1.82, \quad v_2 = 1.82(1 + j), \quad \omega = 6.6248 \quad (24)$$

$$a = 2, \quad v_1 = 2(1 + j), \quad \omega = 8 \quad (25)$$

The corresponding infinite countable time delays of the cross points in (24) and (25) are obtained with regards to relation (15) as $\tau_{1k} = 0.0499 + 0.3019k\pi$ and $\tau_{2k} = 0.25k\pi$, respectively.

By applying the criterion expressed in the previous section, it is easy to find out that a destabilizing crossing of roots ($RT = +1$) has occurred at $\tau = 0.25k\pi$ for $s = \pm j6.6248$ and a stabilizing crossing ($RT = -1$) has taken place at $\tau = 0.0499 + 0.3019k\pi$ for $s = \pm j8$ for all values of $k \in \mathbb{Z}^+$. Therefore, we will have 5 stability windows as follows: $0.0499 < \tau < 0.7854$, $0.9983 < \tau < 1.5708$, $1.9486 < \tau < 2.3562$, $2.8953 < \tau < 3.1416$ and $3.8437 < \tau < 3.9270$ which agree with the results presented by [3] and [8]. Note that at $\tau = 3.9269$, an unstable pair of poles crosses toward the right half-plane, and before this unstable pole pair can turn to the left half-plane at $\tau = 4.7922$, another unstable pair of poles goes toward the right half-plane at $\tau = 4.7123$; and thus, the system can not recover the stability. To get a better understanding of the properties of this system, its Root-Locus curve has been plotted as a function of delay in Fig.1.

6 Conclusions

An efficient method to analyze the BIBO stability of a large class of time-delayed fractional order systems for both single and commensurate-delay cases is proposed. The method introduces an augmented equation whose real roots give the finite values of crossing frequencies at which stability feature of the system change. According to the infinitely countable time delays corresponding to each crossing point, the parametric space of τ is discretized to investigate stability in each interval. Finally, an illustrative example is presented to highlight the proposed approach.

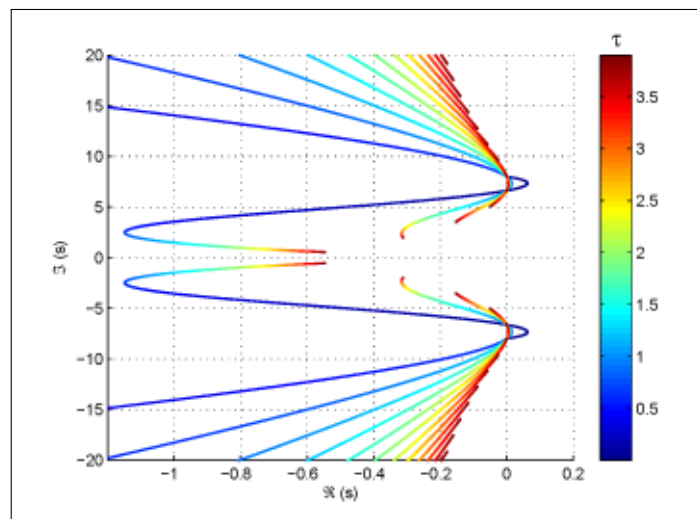


Figure 1: Root-loci for $C(\sqrt{s}, \tau)$ until $\tau=3.9$

Bibliography

- [1] Stojanovic, S.B.; Debeljkovic, D.L.J.; Dimitrijevic, N. (2012); Stability of Discrete-Time Systems with Time-Varying Delay: Delay Decomposition Approach, *Int J Comput Commun*, ISSN 1841-9836, 7(4): 775-783.
- [2] Liu, C.L.; Liu, F. (2010); Consensus Problem of Second-order Dynamic Agents with Heterogeneous Input and Communication Delays, *Int J Comput Commun*, ISSN 1841-9836, V(3):325-335.
- [3] Pakzad, M. A.; Pakzad, S.; Nekoui, M.A. (2013); Stability analysis of time-delayed linear fractional-order systems, *International Journal of Control, Automation, and Systems*, 11(3): 519-525.
- [4] Pakzad, S.; Pakzad, M. A. (2011), Stability condition for discrete systems with multiple state delays, *WSEAS Trans. on Systems and Control*, 6(11): 417-426.
- [5] Walton, J KE.; Marshal, JE. (1987), Direct method for TDS stability analysis, *IEE Proceeding Part D. 134*, 101-107.
- [6] Bonnet, C.; Partington, J.R. (2002), Analysis of fractional delay systems of retarded and neutral type, *Automatica*, 38(7):1133-1138.
- [7] Pakzad, M. A.; Pakzad, S.; Nekoui, M.A. (2013); Stability analysis of multiple time delayed fractional order systems, *American Control Conference, Washington, DC*, 170-175.
- [8] A. R. Fioravanti, C. Bonnet, H. Ozbay and S. I. Niculescu, (2012), A numerical method for stability windows and unstable root-locus calculation for linear fractional time-delay systems, *Automatica*, 48(11):2824-2830.
- [9] Hwang, C.; Cheng, Y.C. (2006), A numerical algorithm for stability testing of fractional delay systems, *Automatica*, 42(5): 825-831.
- [10] Hwang, C.; Cheng, Y.C. (2005), A note on the use of the lambert w function in the stability analysis of time-delay systems, *Automatica*, 41(11):1979-1985.

Dimensionality Reduction and Generation of Human Motion

S. Qu, L.D. Wu, Y.M. Wei, R.H. Yu

Shi Qu*

Air Force Early Warning Academy
No.288, Huangpu Road, Wuhan, 430019, China
*Corresponding author: qushi@nudt.edu.cn

Ling-da Wu, Rong-huan Yu

School of Equipment, Beijing, 101400, China

Ying-mei Wei

National University of Defense Technology
Changsha, 410073, China

Abstract: To reuse existing motion data and generate new motion, a method of human motion nonlinear dimensionality reduction and generation, based on fast adaptive scaled Gaussian process latent variable models, is proposed. Through statistical learning on motion data, the motion data are mapped from high-dimensional observation space to low-dimensional latent space to implement nonlinear dimensionality reduction, and probability distributing of posture space which measures the nature of posture is obtained. The posture which meets constraints and has maximal probability can be computed as the solution of inverse kinematics. This method can avoid cockamamie computation and posture distortion existing in traditional inverse kinematics. The experiments show that our method has higher convergence velocity and precision and extends editing range of motion by adapting motion editing direction.

Keywords: dimensionality reduction, Gaussian process, kernel function, motion generation.

1 Introduction

As the development of motion capture technology, it became reality to generate virtual character motion with motion capture data. This method, easy to implement and generating vivid motion, has been applied to movie, advertisement, game, military affairs and so on. As complexity and multiformity of motion, existing motion capture data are not enough to satisfy application in practice, it is need to reuse existing motion capture data to generate new motion.

Motion capture data, often high dimensional, require to be reduced dimension before analysis and research on them. According to mapping relation between high dimensional data and low dimensional data, dimensionality reduction techniques can be classified to linear or nonlinear. Principal component analysis (PCA) [1] is a linear dimensionality reduction technique in common use, which projects the high dimensional data along principal axis chosen by computing variance. Through improving PCA, many correlative techniques are proposed, which have better capability than PCA in specific application field. The probabilistic principal component analysis (PPCA) [2], introducing probability model into PCA, can complete dimensionality reduction and density estimate by modeling likelihood. In PPCA, which is linear too, high dimensional data and low dimensional data obey probabilistic distribution with noise. Combining PCA and kernel method [3], the kernel principal component analysis (KPCA) [4] is proposed. The basic idea of KPCA is that map the high dimensional data to a feature space by nonlinear method, and do principal component analysis in it. KPCA, as a nonlinear method, has the excellence of kernel method that establishes connection between high dimensional data and low dimensional data with kernel function not requires determining the specific relation between them. Gaussian process

latent variable models (GPLVM) [5] [6] is another effective nonlinear dimensionality reduction technique, which completes dimensionality reduction and density estimate using Gaussian process regress. The GPLVM, different from PPCA, does not estimate the parameters of mapping but margin them. From the viewpoint of Bayesian, GPLVM maximizes the posterior probability to reduce dimensionality given prior probability.

The literature [7] implemented dimensionality reduction and generation of human motion based on scaled Gaussian process latent variable models (SGPLVM), in which there are two shortages, the one is slow convergence, and the other is small editing range of motion. In this paper, we improve on the work to propose fast adaptive scaled Gaussian process latent variable models (FASGPLVM), and implement dimensionality reduction and generation of human motion based on it. Our method has higher convergence velocity and precision and extends editing range of motion by adapting motion editing direction.

2 Overview

2.1 Human Skeleton and Motion Depiction

Human body consists of skeleton, muscle, skin and so on, and the skeleton determines human posture from the angle of motion. In this paper, we construct a simplified human skeleton as shown in Figure 1, which consists of 22 bones. The number in parentheses denotes the degree of freedom of corresponding bone, amounting to 50.

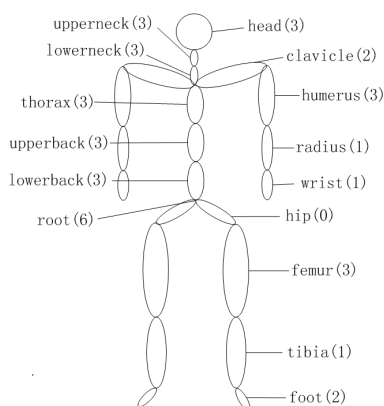


Figure 1: Simplified human skeleton

The degree of freedom denotes translation or rotation direction of bone. In the simplified human skeleton, the root has 3 translation directions and 3 rotation directions, the rest bones have 0 to 3 rotation directions separately. In the process of motion, human posture and location are determined by these degrees of freedom, concretely, the translation of root determines human location and the rotations of all bones determine human posture. So, a feature vector y can be constructed to represent human motion posture, which consists of all rotation degrees of freedom. In this paper, y is a column vector and $y \in R^{47}$. A motion, containing N frames, can be represented as a matrix $Y = [y_1, \dots, y_t, \dots, y_N]^T$, where $t(t = 1, 2, \dots, N)$ is the time index and y_t is the feature vector correspond to the posture at time t .

2.2 Dimensionality Reduction and Generation of Human Motion

The framework of dimensionality reduction and generation of human motion based on FASGPLVM, shown as Figure 2, contains two stages, which are training model stage and generating

posture stage. In training model stage, train FASGPLVM using sample motion data to obtain optimized latent trajectory and model parameters. At the same time, the probability distribution of sample motion is also obtained. In generating posture stage, synthesize the posture which has the maximal probability and satisfies constraints on the end effectors given by user, and update the active set which is a subset of sample motion data. The active set, instead of all sample motion data, are used to train model and synthesize new posture in order to improve computing efficiency.

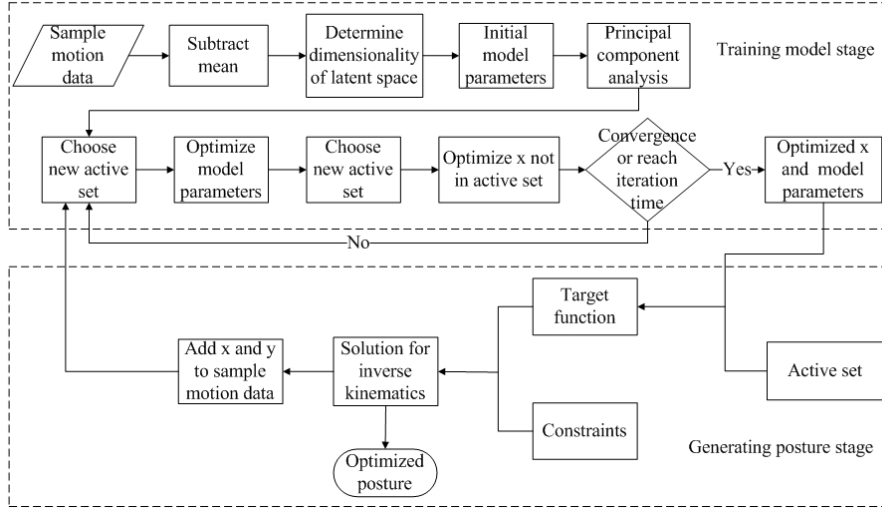


Figure 2: Framework of dimensionality reduction and generation of human motion

3 Fast Adaptive Scaled Gaussian Process Latent Variable Models

Our work is based on the literature [7], in which SGPLVM is proposed by introducing scaled genes to GPLVM. In this paper, FASGPLVM is proposed, based on SGPLVM, to implement dimensionality reduction and generation of human motion. Similar to SGPLVM, scaled genes ω and active set are also included in FASGPLVM. The scaled genes $\omega = \text{diag}(\omega_1, \omega_2, \dots, \omega_D)$ denote relatively important degree of each entry of feature vector y , ωy is used to train model instead of y . The active set is a subset of sample motion data, which is used to improve computing efficiency. Kernel function is the core of model, which has a great influence on training and application of model. With improvement on kernel function of SGPLVM, the summation of radius base function (RBF) and vectorial angle cosine is used as kernel function in FASGPLVM.

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\gamma}{2}(x_i - x_j)^T(x_i - x_j)\right) + \eta \exp\left(\frac{x_i \cdot x_j}{|x_i||x_j|}\right) + \delta_{i,j}\beta^{-1} \quad (1)$$

where x_i is the latent variable correspond to y_i , the parameters α and η mean how correlated pairs of points are in general, γ means the spread of kernel function, β^{-1} is variance of noise and $\delta_{i,j}$ is Kronecker delta function which is 1 when i is equal to j , and 0 otherwise. Kernel function measures how similar two points in latent space are, concretely, a larger $k(x_i, x_j)$ denotes that x_i and x_j are more similar. In kernel function of literature [7], similarity between two points is measured by Euclid distance, and the larger Euclid distance is, the smaller similarity is. In FASGPLVM, vectorial angle cosine is introduced to associate similarity with position of point, and the similarity degrees of two pairs of points who have the same Euclid distance are not always equal. As shown in Figure 3, x_1, x_2, x_3, x_4, x_5 and x_6 are points in latent space which

satisfy $(x_1 - x_2)^T(x_1 - x_2) = (x_3 - x_4)^T(x_3 - x_4) = (x_5 - x_6)^T(x_5 - x_6)$, we have the result $k(x_1, x_2) = k(x_3, x_4) = k(x_5, x_6)$ according to kernel function of literature [7]. According to kernel function of this paper, we have the result $k(x_1, x_2) < k(x_3, x_4), k(x_1, x_2) < k(x_5, x_6)$ because of $\theta > \lambda, \theta > \varphi$, which is more reasonable. x_3x_4 results from rotating x_1x_2 which is perpendicular to center axis (the line connecting origin and center of x_1x_2), and x_1, x_2 have the smallest similarity. When $\lambda = 0$, x_3x_4 and center axis coincide. In this case, $x_3 = cx_4$ (c is a constant) and the similarity between x_3 and x_4 reaches maximum. x_5x_6 results from translating x_1x_2 , the difference between x_5 and x_6 is reduced because of $\|x_5\| > \|x_1\|$ and $\|x_6\| > \|x_2\|$. When x_5 and x_6 apart infinitely from origin, the similarity between x_5 and x_6 reaches maximum since they can be treated as one point.

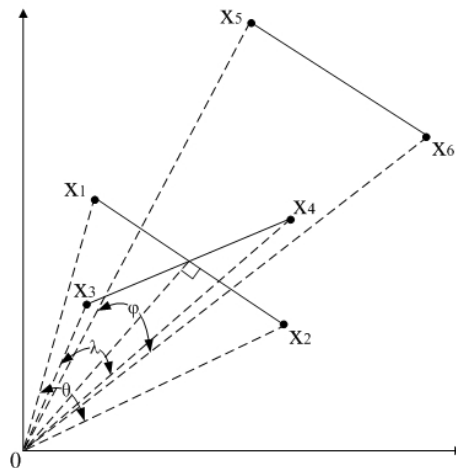


Figure 3: Influence of vectorial angle cosine on similarity measure

Another important improvement on SGPLVM is to introduce dynamic active set. In the motion editing process, the active set of FASGPLVM change, the model parameters $\alpha, \beta, \gamma, \eta, \omega_k$ and K covariance matrix change accordingly, which makes FASGPLVM adaptive to motion editing. In SGPLVM, the new posture is restricted to be similar to sample motion at full steam in order to ensure it vivid. However, the extent of motion editing can not be very large, otherwise, the new posture will be anamorphic. In FASGPLVM, it can enlarge the extent of motion editing and ensure new posture vivid by changing active set. The third improvement on SGPLVM is that convergence judgment, instead of fixed iteration time in SGPLVM, is introduced to training algorithm.

The algorithm of learning sample motion data based on FASGPLVM can be described as follows:

Step 1. Subtract the mean of motion data. The mean vector \bar{y} is computed by

$$\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i \quad (2)$$

Then, subtract the mean from each line of motion data matrix Y .

$$\bar{Y} = [\bar{y}_1, \bar{y}_2, \dots, \bar{y}_N]^T = [y_1 - \bar{y}, y_2 - \bar{y}, \dots, y_N - \bar{y}]^T \quad (3)$$

Step 2. Determine the dimensionality of latent space. FASGPLVM maps the motion data from high-dimensional observation space y^D to low-dimensional latent space x^d . The dimension-

ality of latent space is 2 or 3, and is determined by

$$d = \begin{cases} 2 & D < 60 \\ 3 & D \geq 60 \end{cases} \quad (4)$$

In this paper, $D = 47$, so $d = 2$.

Step 3. Initialize model parameters. In our experiments, $\alpha = 1, \beta = 1, \gamma = 1, \eta = 1, \omega_k = 1, M = 200, T = 100, C = 0.01$, where M is size of active set, T is allowable maximum iteration time, and C is convergence threshold.

Step 4. Initialize the latent variable with PCA. Compute d principal components of motion data and project motion data y_i to these principal components to obtain latent variable x_i . The initialization results are saved as a matrix $X = [x_1, x_2, \dots, x_N]^T$.

Step 5. Update the active set with informative vector machine (IVM) [8].

Step 6. Estimate parameters $\alpha, \beta, \gamma, \eta, \omega_k$ by minimizing (5) with scaled conjugate gradients (SCG) [9].

$$L_{GP} = \frac{D}{2} \ln|K| + \frac{1}{2} \sum_k \omega_k^2 \hat{Y}_k^T K^{-1} \hat{Y}_k + \frac{1}{2} \sum_l x_l^T x_l + \ln(\alpha\beta\gamma\eta / \prod_k \omega_k^N) \quad (5)$$

where parameters D is the dimensionality of feature vector, \hat{X}, \hat{Y} are active set, \hat{Y}_k is the k -th column of \hat{Y} , and K is covariance matrix whose entries are computed by (1), $K(i, j) = k(x_i, x_j)$.

Step 7. Update the active set with IVM according to step 5.

Step 8. Update the latent variable not in active set by minimizing (6) with SCG.

$$L_{IK(x,y)} = \frac{\|\omega y - g(x)\|^2}{2\sigma^2(x)} + \frac{D}{2} \ln\sigma^2(x) + \frac{1}{2} \|x\|^2 \quad (6)$$

Here

$$g(x) = \mu + \hat{Y}^T K^{-1} k(x) \quad (7)$$

$$\sigma^2(x) = k(x, x) - k(x)^T K^{-1} k(x) \quad (8)$$

Where parameters μ is the mean of motion data, $k(x)$ is a column vector in which the i -th entry contains $k(x_i, x)$, and x is latent variable.

Step 9. Stop iteration if $\max\{\|\Delta\alpha\|, \|\Delta\beta\|, \|\Delta\gamma\|, \|\Delta\eta\|, \|\Delta\omega_k\|\} < C$, otherwise go to step 10.

Step 10. Stop iteration if iteration time reaches T , otherwise go to step 5.

After new posture was generated, corresponding (x, y) is added to sample set. The active set is updated using IVM, parameters $\alpha, \beta, \gamma, \eta, \omega_k$ and covariance matrix K is updated according to step 6. These are used in motion editing of next time in order to ensure FASGPLVM adaptive to motion editing and enlarge motion editing extent.

4 Posture and Motion Generation

Posture generation is a problem of inverse kinematics, which generates vivid posture satisfying constraints on the end effectors given by user. Constraints can be satisfied using geometric method easily; the difficulty in inverse kinematics is how to make new posture vivid because it is difficult to model human motion rules. In FASGPLVM, the probability space of motion posture is modeled, in which the posture near sample postures has bigger probability. The sample postures are motion capture data which has the best fidelity, so the probability of posture can be used

as a measure of fidelity and the posture with bigger probability has better fidelity. In inverse kinematics, the posture which satisfies constraints and has the biggest probability is the solution of inverse kinematics. The problem of generating new posture given constraints on end effectors can be described as a nonlinear optimization problem.

$$\begin{aligned} \text{Min}(x, y) &= L_{IK}(x, y) \\ \text{s.t. } f(y) &= c \end{aligned} \quad (9)$$

Where $f(y) = c$ are constraints on end effectors, denoting the destination position for some bones of human skeleton. Generally speaking, not only one but a set of postures, called feasible posture set, can satisfy the constraints on end effectors, but some of them do not satisfy the motion rule and are not vivid. The posture, vivid and satisfying motion rule, is selected from feasible posture set by minimizing the target function $\text{Min}(x, y) = L_{IK}(x, y)$.

In our method, it is not needed to give constraints for all bones of human skeleton. In $f(y) = c$, just a few constraints interested for user are contained, which are called sparse constraints. In the process of optimization, weight gene is introduced to transform this problem to unconstrained nonlinear optimization.

$$\text{Min}(x, y) = (1 - \lambda)L_{IK}(x, y) + \lambda(f(y) - c) \quad (10)$$

where λ is the weight gene ($0 < \lambda < 1$), it denotes intensity of constraints. In this paper, $\lambda = 0.999$.

Motion generation is based on posture generation, which is implemented by editing motion trajectory. Given a motion, we edit the motion trajectory of one joint and optimize (9) constrained by new motion trajectory to generate new motion.

5 Results and Analysis

In this section, we design experiments of dimensionality reduction and generation of human motion to prove the algorithms in this paper correct and effective. In our experiments, the human skeleton consists of 22 bones with 50 degrees of freedom. The observation space is 47 dimensional and the latent space is 2 dimensional. The sample motion data come from human motion capture database of Carnegie Mellon University¹.

5.1 Posture Generation

Jump and kick motion data are used to train FASGPLVM, which consist of 180 frames shown as figure 4. To represent motion process clearly, the left arm and left leg of human skeleton are rendered with gray color.

After trained with jump and kick motion data, the parameters of FASGPLVM are estimated as $\alpha = 456.912$, $\beta = 0.116$, $\gamma = 4.745$, $\eta = 326.582$. Then, we can generate new posture satisfying constraints given by the user. To compare with SGPLVM, train SGPLVM with the same motion data and generate new posture satisfying the same constraints. The experiment results are shown as figure 5, in which (a) is the origin posture, coming from sample motion data; (b), (c) and (d) are new postures generated by FASGPLVM; (e), (f) and (g) are new postures generated by SGPLVM. The dot on the right foot denotes the position of the right foot constrained by the user, (b) and (e) have the same constraint, and so do (c) and (f), (d) and (g). The process of editing posture (a) is to raise the right foot to generate new postures (b), (c), (d), (e), (f)

¹<http://mocap.cs.cmu.edu/>

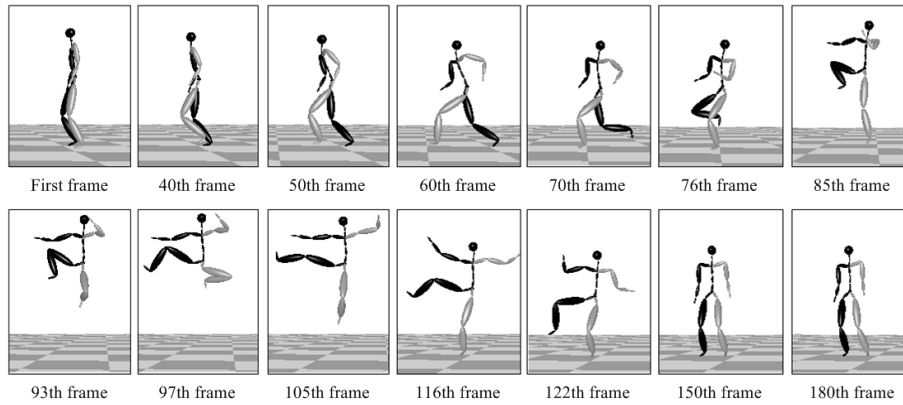


Figure 4: Motion capture data(jump and kick)

and (g). It can be seen from figure 5 that (d) is a vivid posture but the posture (g) distorts. This is because that the height of the right foot in the posture (g) oversteps the limitation of editing range of SGPLVM; but the posture (d) is generated based on FASGPLVM which has larger editing range.

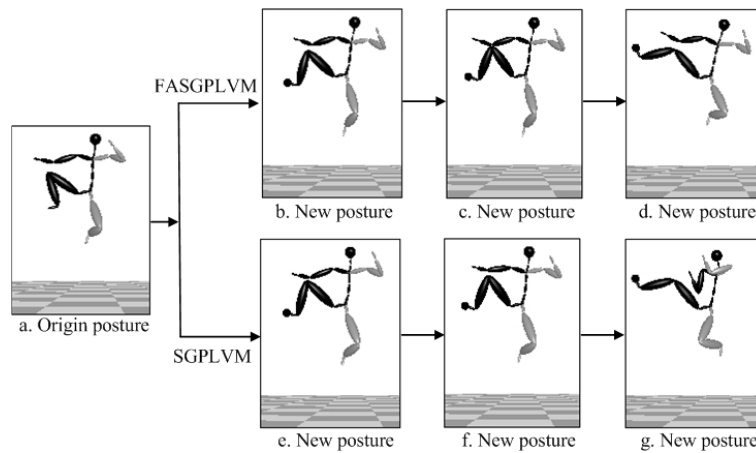


Figure 5: Posture generation (jump and kick)

5.2 Motion Generation

In this section, we generate new motion based on FAGPLVM. Editing the jump and kick motion shown in figure 4 to generate new motion. The 103th frame and the 117th frame are edited, raising the right foot. To generate new motion, the right foot positions of the 89th frame, the 103th frame, the 117th frame and the 131th frame are used as controlling points for Hermite interpolation to generate new motion trajectory. The new motion trajectory constrains the positions of the right foot in new motion, optimizing posture of every frame satisfying this constraint to generate new motion. As shown in figure 6, the black curve denotes the motion trajectory of the right foot.

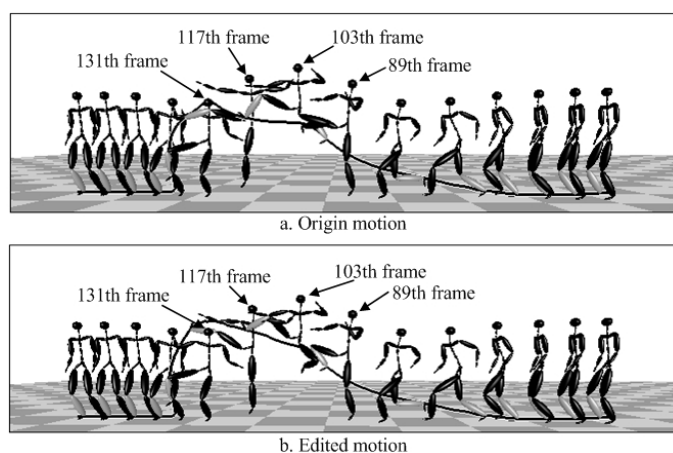


Figure 6: Motion generation (jump and kick)

5.3 Comparison and Analysis

Table 1 is the convergence precision comparison of FASGPLVM and SGPLVM after the same iteration time. The convergence precision is defined as the change of model parameters. So, the convergence precision of FASGPLVM is $\max\{\|\Delta\alpha\|, \|\Delta\beta\|, \|\Delta\gamma\|, \|\Delta\eta\|, \|\Delta\omega_k\|\}$, and the convergence precision of SGPLVM is $\max\{\|\Delta\alpha\|, \|\Delta\beta\|, \|\Delta\gamma\|, \|\Delta\omega_k\|\}$. It can be seen from table 1 that the convergence precision of FASGPLVM is higher than that of SGPLVM, this is because of improving on the kernel function in FASGPLVM.

Tab.1 The convergence precision comparison of FASGPLVM and SGPLVM

Sample motion	Frame amount	Iteration time	convergence precision (FASGPLVM)	convergence precision (SGPLVM)
Kick	341	52	0.004	0.021
Run	264	46	0.002	0.014
Run and jump	536	71	0.005	0.017
Box	452	63	0.008	0.023
Long jump	367	58	0.006	0.026

Table 2 is the editing range comparison of FASGPLVM and SGPLVM retaining the new posture vivid. It can be seen from table 2 that the editing range of FASGPLVM is larger than that of SGPLVM, this is because of introducing dynamic active set in FASGPLVM.

Tab.2 The editing range comparison of FASGPLVM and SGPLVM

Sample motion	Edited frame number	Editing way	Editing range/cm (FASGPLVM)	Editing range/cm (SGPLVM)
Kick	206	Raise fight foot	43	21
Run	124	Raise left hand	26	16
Run and jump	385	Raise fight foot	37	14
Box	162	Lower right hand	28	17
Long jump	256	Lower head	22	12

6 Conclusions

By improving on SGPLVM, FASGPLVM is proposed. A method of dimensionality reduction and generation of human motion based on FASGPLVM is also proposed. Compared with SGPLVM, FASGPLVM has some excellences, such as higher convergence precision and larger editing range.

Acknowledgments

Thanks to the anonymous reviewers for their helpful comments, Neil Lawrence for his on-line source code, and Carnegie Mellon University for human motion capture data.

Bibliography

- [1] Sun Hongwei, Gu Ming and Sun Jiaguang, A coding algorithm using PCA-based correlation vector quantization, *Journal of Computer-Aided Design & Computer Graphics*,17(8):1662-1666,2005.
- [2] Gift N, Lorraine B and Isobel C G, Probabilistic principal component analysis for metabolomic data, *BMC Bioinformatics*,11(1):571-582,2010.
- [3] Nisbet R, Elder J and Miner G, *Statistical analysis and data mining*, New York:Academic Press,2009.
- [4] Roman Rosipal et al, Kernel PCA for feature extraction and de-noising in non-linear regression, *Neural Computing & Applications*,10(3):231- 243,2001.
- [5] Carl H, Philip H S and Neil D L, Gaussian process latent variable models for human pose estimation, *Proc. of Machine Learning for Multimodal Interaction*. Brno: Springer-Verlag Press:132-143,2007.
- [6] QU Shi et al, Pose Synthesis of Virtual Character Based on Statistical Learning, *The International Symposium on Computer Network and Multimedia*, Wuhan, China:36-39,2009.
- [7] Keith Grochow et al, Style-based inverse kinematics, *ACM Transactions on Graphics*,23(3): 522-531,2004.
- [8] Neil D. Lawrence, Matthias Seeger and Ralf Herbrich, Fast sparse Gaussian process methods: the informative vector machine,*Proceedings of Neural Information Processing Systems 15*. MIT Press:609-616,2003.
- [9] Martin F. Muller, A scaled conjugate gradient algorithm for fast supervised learning, *Neural Networks*,6(4):525-533,1993.

Security Ontology for Adaptive Mapping of Security Standards

S. Ramanauskaitė, D. Olifer, N. Goranin, A. Čenys

Simona Ramanauskaitė*, Dmitrij Olifer,
Nikolaj Goranin, Antanas Čenys

Vilnius Gediminas Technical University

simona.ramanauskaite,dmitrij.olifer, nikolaj.goranin, antanas.cenys@vgtu.lt

Lithuania, LT-10223 Vilnius, Sauletekio al. 11

*Corresponding author: simona.ramanauskaite@vgtu.lt

Abstract: Adoption of security standards has the capability of improving the security level in an organization as well as to provide additional benefits and possibilities to the organization. However mapping of used standards has to be done when more than one security standard is employed in order to prevent redundant activities, not optimal resource management and unnecessary outlays. Employment of security ontology to map different standards can reduce the mapping complexity however the choice of security ontology is of high importance and there are no analyses on security ontology suitability for adaptive standards mapping.

In this paper we analyze existing security ontologies by comparing their general properties, OntoMetric factors and ability to cover different security standards. As none of the analysed security ontologies were able to cover more than 1/3 of security standards, we proposed a new security ontology, which increased coverage of security standards compared to the existing ontologies and has a better branching and depth properties for ontology visualization purposes. During this research we mapped 4 security standards (ISO 27001, PCI DSS, ISSA 5173 and NISTIR 7621) to the new security ontology, therefore this ontology and mapping data can be used for adaptive mapping of any set of these security standards to optimize usage of multiple security standards in an organization.

Keywords: security ontology, security standards, adaptive mapping.

1 Introduction

An ontology defines the basic terms and relations comprising the vocabulary of a topic area as well as the rules for combining terms and relations to define extensions to the vocabulary [9]. The ontology provides a better communication, reusability and organization of knowledge by decreasing language ambiguity and structuring transferred data [1], [2], [3], [4].

Security becomes fundamental in our society and the survival of organizations depends on the correct management of up-to-date security elements [6]. As the security area is very broad and has many relations between its concepts, usage of security ontology could improve security knowledge description unambiguity in information systems. The necessity of security ontology can be noticed in various security communities and considered as an important challenge and a research branch [5], [7], [8].

In small and medium enterprises the knowledge database of security area and its unambiguity is very important in formal/legal activities, such as certification, standard compliance, etc. In many cases organizations have to meet certain security requirements from different sources, which may be redundant or overlapping by simultaneous usage. Therefore standard mapping should be put into practice in cases where more than one security standard has to be met. The mapping of security standards allows the optimization of resources by indicating matching elements of standards and by eliminating duplicated activities and security measures to meet it. However mapping of security standards can be complicated if more than two standards have to be mapped.

Security ontology can be used to simplify the mapping of more than two security standards [10]. This solution suggests mapping all security standards under one security ontology. This ontology would act as a basis for standard knowledge formalization and would allow adaptive mapping of any standards, mapped to the ontology. Therefore the security ontology plays a large role in adaptive mapping and covers a wide area and should be detailed to meet all security standards.

The aim of this paper is to analyse suitability of existing security ontologies to be used for adaptive mapping of security standards and to propose a new one, more suitable for the purpose.

2 Security Ontology

Existing security ontologies vary according to described area and level of detail. One of the first works mentioning information system knowledge concepts concerning security was published in 1990 by J. Mylopoulos et al. The paper "Telos: Representing Knowledge about Information Systems" [12] describes a Telos language to describe the knowledge about information systems and suggests it can be employed for security specification as well. C. E. Landwehr et al. on 1994 published a paper called "A taxonomy of computer program security flaws" [13] where types of computer program security flaws were summarized and claimed it can be used for introduction to the characteristics of security flaws and their origins. A. Avizienis et al. also proposed a taxonomy, concerning security concepts [14]. This taxonomy describes more abstract and wide concepts than C. E. Landwehr et al. provided however clear relationships between categories of taxonomy are missing too.

The need of ontology rather than taxonomy was indicated in a paper "Toward a Security Ontology" by M. Donner on 2003 [7]. On the same year G. Denker et. al. presented security related ontologies for web services and published it in paper "Security in the Semantic Web using OWL" [15] while H. Mouratidis et al. published work "An Ontology for Modelling Security: The Tropos Approach" [16] where presented ontology for security modelling in agent-based information systems. H. Mouratidis provided more works concerning security ontologies [17], [5] where clear orientation to usage of security ontologies in software developments is noticed, therefore these ontologies are meant more for system requirement representation rather than for basic security concepts.

There are ontologies concentrated specifically on security requirements only. One of such ontologies is presented by F. Massacci [18]. Other specific security ontologies are proposed by D. Geneiatakis et al. [19] (designed for describing Session Initiation Protocol security flaws), by M. Karyda et al. [20] (dedicated for describing applications of e-government), by J. Undercoffer et al. [21] (designed for describing computer attacks), by A. Souag [22] (designed for requirements engineering process) and by other authors. A. Kim extended specific ontologies and created one which can be applied to any electronic resource [23]. However this ontology does not overlay all the concepts of information security. More detailed general security ontologies were proposed by A. Herzog et al. [24] and S. Fenz et al. [25].

Security ontology, proposed by Herzog et al. represents information security domain that includes both general concepts and specific vocabulary of the domain. The proposed ontology has 4 top level concepts: assets, threats, vulnerabilities and countermeasures. The ontology overviews the information security domain in a context-independent and application neutral manner. Similar properties apply to security ontology proposed by S. Fenz et al. however it has more concepts in it including non-core concepts such as the infrastructure of organizations. The main top level concepts in this ontology are: asset, control, organization, threat and vulnerability.

Table 1: Data of general comparison of security ontologies

Property	Ontology			
	G. Denker	A. Herzog	S. Fenz	S. Fenz (raw)
Total number of classes	39	460	641	311
Total number of data types properties	0	7	16	14
Total number of object properties	12	30	58	58
Total number of annotation properties	2	4	10	10
Total number of individuals	117	211	486	478
Number of sub-classes	11	571	1051	409
Max. depth of class tree	4	8	6	6
Min. depth of class tree	1	1	1	1
Avg. depth of class tree	1,4	4,1	3,0	3,2
Max. branching factor of class tree	27	83	199	114
Min. branching factor of class tree	1	1	1	1
Avg. branching factor of class tree	7,6	3,2	3,9	14,5

3 Analysis of Security Ontologies

Three security ontologies were chosen for deeper analysis because of its particularity: security ontology created by G. Denker; Security ontology, created by A. Herzog et al.; Security ontology, created by S. Fenz.

While S. Fenz security ontology includes concepts of several security standards (ISO 27001, Grundschutz) in it, one more version of S. Fenz's security ontology will be analyzed in this study (hereinafter S. Fenz (raw)). All classes and elements of security standards will be excluded from S. Fenz's ontology, relying solely on raw concepts of ontology security.

3.1 General Comparison

In general comparison of security ontologies the total number of different ontology elements, the depth and branching metric of the ontology tree are put into comparison. To get these metrics an OWL ontology editor SWOOP was used. Data obtained by this tool are presented in Table 1.

As results of general ontology comparison reveal, G. Denker's ontologies have the least number of concepts, while security ontologies, created by S. Fenz and A. Herzog have the largest number of concepts. G. Denker's ontology is intended to interface between various notations of security standards while ontologies of S. Fenz and A. Herzog represent the whole area of security, therefore have more concepts.

The purpose of ontology usage inflicts on the number of individuals as well - wider range ontologies have more individuals to allow user to chose from; specific purpose ontologies have less or no individuals as all individuals should be known or unnecessary to the user.

Another important metric is the depth and branching factor of ontology class tree. It defines the main properties of tree structure of the ontology and can be exercised to define how intuitive the ontology should be for individual users. Our analysis displays the security ontology of A. Herzog has the deepest class structure and has the most substantial detailing level. However the maximum branching factor of class tree is equal to 83, which may result in human users facing difficulties while viewing the ontology. Ontology of S. Fenz should be difficult to visualize as well, because of its branching factor.

Table 2: OntoMetric analysis data of the ontologies content and the contents organization

Characteristic	Ontology			
	G. Denker	A. Herzog	S. Fenz	S. Fenz (raw)
Concepts (factor)	2	4	4	4
Relations (factor)	3	3	3	3
Taxonomy (factor)	2	3	3	3
Axioms (factor)	2	4	4	4

3.2 OntoMetric Analysis of Security Ontologies

General comparison of security ontologies gives just a few main quantitative metrics, while the quality of ontology is not taken into account. OntoMetric [26] is a method for ontology quality measurement. This method compares ontologies in five dimensions (the ontologies content and the contents organization; the language in which it is implemented; the methodology that has been followed to develop it; the software tools used to build and edit the ontology; the costs that the ontology will require in a certain project) and measures all the characteristics from 1 to 5 according to their low or high degree of accomplishment.

While all ontologies we are analyzing are written in the same file format, we are analyzing the content metrics alone (metric of language, tools and costs should be equal, because all analyzed ontologies are written in OWL files, while the development process of ontology do not have significant influence on its usage and are unknown to us). According to OntoMetric, the content of ontology can be defined by 4 factors: concepts, relations, taxonomy and axioms.

As evaluation of OntoMetric is qualitative, we will evaluate all of them as all security ontologies are meant for presenting the broadest security area possible and should be able to present any situation in area of information security. The imagination of ideal security ontology is important in order to evaluate the concept factor in OntoMetric analysis as this measurement should provide information on how well the security area is covered by the ontology.

Other factors in OntoMetric analysis are more relative and describes how well the relations, taxonomy and axioms are described in the ontology, not the whole security area.

All data of our OntoMetric analysis are presented in Table 2.

The OntoMetric analysis shows the G. Danker ontology has the lowest scores, while S. Fenz and A. Herzog ontologies have similar scores, the level of detail and provides wide range of security concepts. However the data of OntoMetric analysis does not show differences between S. Fenz and A. Herzog.

While comparing the differences in S. Fenz's and A. Herzog's ontologies, it can be noticed that ontology, created by A. Herzog has more of a theoretical approach rather than the ontology of S. Fenz and describes more definitions, formal concepts of information security area. S. Fenz's ontology provides more information on practical side of information security, by listing basic controls as a guide for security administrators for system security assurance however does not mention concepts, related to organizational security.

3.3 Research of Security Ontology Usage for Mapping of Security Standards

As security ontologies, proposed by A. Herzog and S. Fenz have similar ontology comparison results, a deeper analysis has to be performed to select the best one for usage in adaptive mapping of security standards.

Adaptive Mapping of Security Standards

To ensure security in an organization, security standards or best practices can be employed. In some cases compliance to a certain security standard is even required to obtain privileges to supply or to get different services. However when organization uses more than one security standard, mapping or integration of security standard usage should be done in order to avoid redundant activities, not optimal resource management, unnecessary outlays etc. Integration or direct mapping of security standards are time and knowledge consuming as well as very static (everything has to be redone when a standard has to be removed or added), while adaptive mapping of security standards provides more flexibility to change the list of used standards as well as requires less work to map a larger number of standards as each standard have to be mapped to ontology only. Therefore n mapping activities have to be done to map n standards in stead of $n*(n-1)$ mappings for direct mapping. The process of adaptive mapping and integrated standard generation is presented in Fig. 1

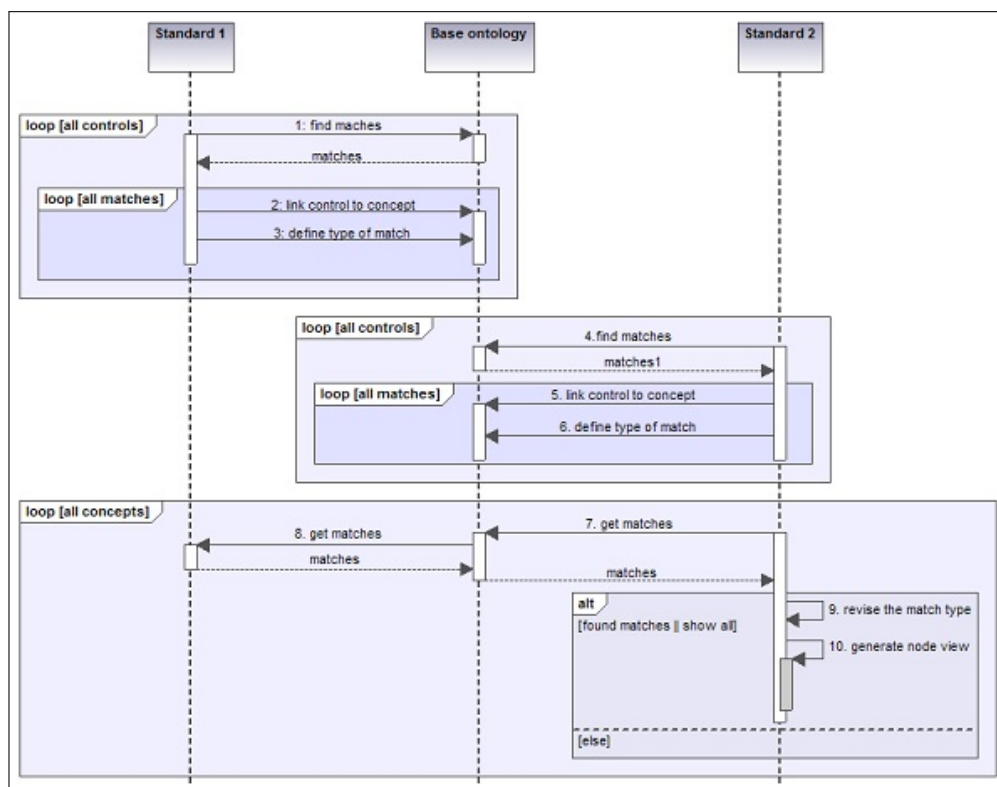


Figure 1: Sequence of mapping two standards and generating the mapped standard in relation to the structure of ontology

When a standard is mapped to the base ontology, all matching controls and concepts between ontology and standards have to be linked. This has to be done once for all standards which have to be mapped together. The generation of standard maps or integrated standards is dynamic and can be done on demand by changing standards which have to be mapped or integrated, properties for relation type estimation etc. The map generation process finds similar controls in selected standards by comparing its linking to the base ontology. An example of relation type estimation in adaptive mapping is provided in Fig. 2.

In this example a control in ISO 27001 (A.8.3.3_Removal_of_access_righ...) and control in PCI DSS (PCI_DSS_8_5_4) standards are mapped with the same links to the security ontology (control in one standard has the same relations to concepts of security standard as control in

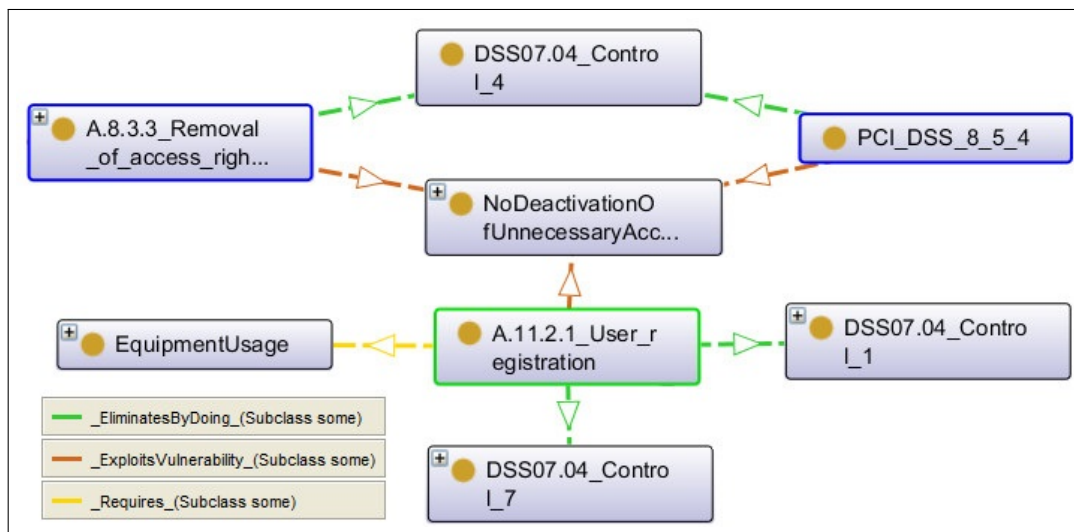


Figure 2: Example of standard mapping through ontology

another standard). As these two controls have no differences in mapping, the full match relation between these two controls of different standards can be generated. One more control of ISO 27001 standard (A.11.2.1_User_registration) is presented in this example to illustrate relevant (not matching) controls. These two controls of ISO 27001 security standard define situations, where vulnerability of nonblocked unnecessary accounts or terminals can be exploited. However both of ISO 27001 controls have more links to different concepts of security ontology, therefore these two ISO 27001 controls can not be treated as equal, however are relevant on certain levels. This kind of information can be used to analyse security standards and to optimize the resource usage when multiple security standards have to be met in an organization.

For visualization and analysis of overlapping of multiple security standards a tool for adaptive mapping of security standards was created (see Fig. 3). This tool uses security ontology and maps security standards data to generated tree structure hierarchies, representing a chosen security standard or integration of few security standards as well as providing additional data through node notation and explanation boxes on similarities of controls in chosen security standards.

Research of Security Standards Coverage by Security Ontologies

To compare which security ontology is more suitable for adaptive security standards mapping and adaptive mapping, ontology and standard concept coverage were analyzed. A. Herzog's and S. Fenz's ontologies were mapped with:

- ISO27001 - the most popular security standard, which was created according to British Security standard BS7799. This standard covers practically all security areas, provides certification opportunity and is widely recognized.
- PCI DSS - security standard developed by such worldwide organizations as Visa, MasterCard, American Express, Discover and JCB. Standard developed to ensure cardholder information protection. This standard is a "Must-have" for all organizations who handle debit, credit, prepaid and other cards. Otherwise these organizations are forbidden to use Visa, MasterCard, American Express and other cards.
- ISSA 5173 - security standard for SME (Small Medium Enterprise). This standard has not been approved or officially recognized, however describes main security requirements,

which need to be implemented in any organization.

- NISTIR 7621 - security standard, developed by national institute of Standards and technology. Document clearly defines which actions are "absolutely necessary" for information, systems and networks protection. It also provides best practices on needed security level implementation.

Data on links between these security standards are presented in static form for 2 specific standards (as a table with matching controls between two security standards [28]) mostly. S. Fenz was the first who mapped ISO 27001 and Grundschutz security standards to his ontology. He used this mapping for purposes of automated risk and utility management [11], however this information can be also used for adaptive standard mapping. S. Fenz mapped two standards only, therefore links can only be generated between ISO 27001 and Grundschutz security.

We analyzed all controls in all 4 chosen standards and mapped them to related concepts in S. Fenz and A. Herzog security ontologies. The mapping of security standards was performed by mapping the lowest level concepts (usually it's a certain control, requirement for the organization), while the classes in security standards, used for presentation of class hierarchy were not accounted as mapping objects.

The process of security standard mapping to security ontologies revealed differences between analyzed ontologies as well. Biggest part of mapping links in S. Fenz's ontology are very direct - one requirement of the standard has an equal or very similar control in S. Fenz's ontology. This type of mapping links are very direct, easy to understand for individual users, however the controls have to be detailed by other links between different concepts of the ontology, otherwise it will be difficult to define relations between standards controls, clustering, etc. Meanwhile mapping security standards according to A. Herzog's ontology was done from logical structure standpoint - one requirement of security standard is to have several links to ontology, by describing which concepts of ontology are related to this requirement (by defining what and how one has to do or use to protect against certain threat or vulnerability). This type of mapping requires more mapping links and has a potential to be easier to cluster controls of security standards into relevant groups. This type of mapping would be more understandable to information systems however would require more analysis or visualizing tools for people to understand links between two security standards, mapped through ontology this way.

Summarizing the security standard mapping process to security ontologies - S. Fenz's ontology can be used to simplify the mapping of security standards because all the most important concepts for mapping are described as list of classes, while in ontology of A. Herzog's mapped classes have more links to ontology and provide more analysis and application possibilities after the mapping is done.

In table 3 data on ontology coverage by standard (covered) and standard coverage by ontology (covers) are provided. Column "covered" defines what part of security ontology was used to map certain standard while column "covers" defines what percentage of security standard was mapped to the security ontology. The property "covers" is more important in this research as it provides information on how well the ontology is capable to present certain security standards in the knowledge database.

The analysis of security ontology and standard coverage revealed that ontologies of A. Herzog and S. Fenz are not capable to fully cover none of analyzed security standards: only security standards with small number of controls or requirements can be mapped with security ontology to cover more than 50% of standard controls; security standards with more than 100 controls or requirements can not be mapped to A. Herzog's and S. Fenz's security ontologies to cover more then 30% of standard controls or requirements. This shows the fact that these two security ontologies do not have all necessary concepts to be fully mapped to security standards.

Analysis of concepts of security ontologies to be employed to map security standard revealed that just a small part (5-18%) of classes from A. Herzog's and S. Fenz's ontologies are mapped directly to security standards. This number could be improved by providing more detailed concept of relationship, however it allows defining what part of ontology is directly related to concepts, mentioned in security standards.

Security ontology, created by S. Fenz was able to cover a larger part of analyzed security standards than A. Harz's ontology. The biggest difference (29% and 19%) was noticed in PCI DSS standard. This could be an argument to chose S. Fenz's security ontology if a company is working with PSI DSS standards, while coverage differences for other analyzed standards are minor. However to cover 29% of PSI DSS standard is not enough to represent it. A new security ontology with more security concepts could help to improve the situation and would allow mapping of bigger parts of security standards.

4 New Security Ontology

As S. Fenz's and A. Herzog's ontologies have low security standard coverage and are not the excellent choice for adaptive mapping of security standards we have created a new general purpose security ontology, which would extend these two ontologies and would be more suitable for adaptive mapping of security standards.

Our ontology has 5 top level classes (see. Fig. 3): asset, countermeasure, organization, threat and vulnerability. These 5 classes are the most basic in security area and are detailed in lower levels.

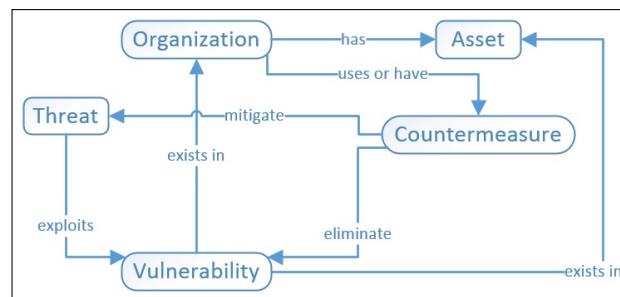


Figure 3: Top level structure of proposed Security ontology

Asset class describes both tangible and intangible asset an organization can have. We describe this class more appropriately than other security ontologies do with the addition of more knowledge on used data by the organization, location and other equipment, owned or used in the company and related to organization's security. The intangible asset is divided into Data and Software (see Fig. 4), while inner structure describes various types of data and software. The tangible assets are structured to subclasses of Movable and Unmovable asset (see Fig. 4). Immovable asset describes location and building concepts and main elements which can be found in it. We structured movable assets into 4 subclasses: alarm systems and detectors; furniture; IT components; utilities. These 4 categories allows the creation of more links to security standards by defining what kind of assets are involved into certain controls (who is at risk, who has a vulnerability etc.).

Countermeasure and Threat classes are described pretty well in A. Herzog's ontology, therefore we made minor changes to it and use similar structure and components as A. Herzog did.

The need for more organizational concepts arose during the mapping of security standards to A. Herzog's and S. Fenz's security ontologies. These two ontologies have a poor description

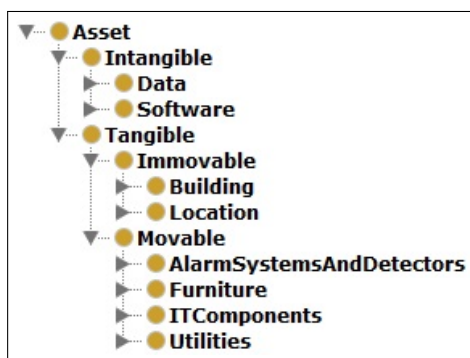


Figure 4: Basic hierarchy of asset concept

of organization structure and policies, while the companies' information security policy is the most important to ensure its security. As subclasses of organization Department, Personnel and Policy concepts were distinguished (see Fig. 5). Precise control description of security standards can be achieved if links to certain executor can be made. Therefore Department and Personnel classes were added and detailed to distinguish plausible types of departments and positions in it.

COBIT is a framework [27], which describes the best ideas for information technology management, quality, evaluation and improvement. Therefore we adopted COBIT 5 framework into our ontology, by defining IT policy class as organization policy subclass, where all COBIT 5 ideas are detailed (see Fig. 5). The COBIT 5 framework was exercised as it is in IT policies class. This guaranties the intuit of ontology usage to those, who is familiar with COBIT framework. Meanwhile in order to propose multiple views and ways to find necessary concepts in the ontology, more subclasses were added to Policy class (see Fig. 5). These classes should present more general policies of the organization however most of them have relations to classes of IT policies class (COBIT 5 framework).



Figure 5: Basic hierarchy of organization concept

Vulnerability class was not detailed properly in A. Herzog's and S. Fenz's ontologies as well. S. Fenz provides a list of vulnerabilities describing individuals with no structure, while A. Herzog describes simply basic types of vulnerabilities. Therefore we extended vulnerability class by dividing it into Code vulnerabilities, Configuration vulnerabilities, Design vulnerabilities, Policy vulnerabilities and Transfer vulnerabilities (see Fig. 6). Those classes are detailed to reflect the basic security vulnerabilities, however they are more structured than in S. Fenz's ontology, to make it more intuitive and simpler to visualize.

To use this ontology as a base for adaptive mapping of security standards a clear and intuitive

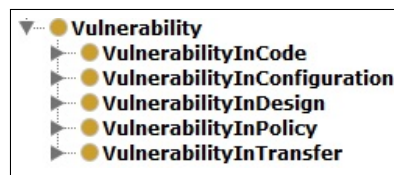


Figure 6: Basic hierarchy of vulnerability concept

ontology structure has to be maintained. We optimized the tree structure of the ontology, therefore now it has 1795 classes, average depth of class tree is 6,5 (has maximum up to 9 depth of class tree) and average branching factor of class tree is 4,8 (has from 1 to 18 subclasses). Such a structure is more viewable in tree structure and should be more intuitive for ontology users (see Fig. 7).

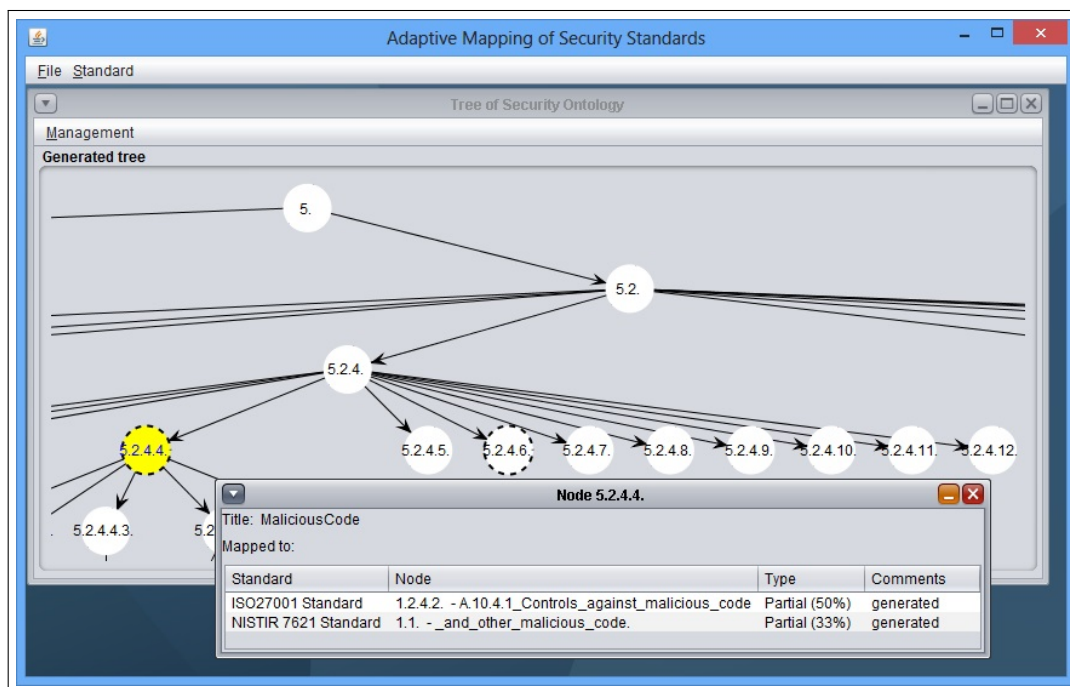


Figure 7: Structure fragment of proposed security ontology and adaptive mapping data in AMSS (created tool for adaptive mapping of security standard)

While structure optimization of new security ontology is more important to ensure user friendly usage and understanding, new concepts allowed a better coverage of security standards. We do not provide direct list of controls and use similar ontology structure for standard mapping as A. Herzog therefore security standard mapping to this ontology has to be done by defining more than one relation to ontology. This mapping property is useful to analyze and to map security concepts in different standards.

ISO27001, PCI DSS, ISSA 5173 and NISTIR 7621 security standards were specified and mapped to it in order to evaluate its suitability to map security standards. Using this ontology as a base for adaptive mapping, 80% of ISO27001, 100% of PCI DSS, ISSA 5173 and NISTIR 7621 standards were mapped to the ontology (see Table 3).

The 100% mapping of ISO27001 standard was not achieved because we did not mapped very specific requirements in security standard (like security properties of used operating system etc.) to more abstract in our ontology.

Table 3: Coverage of ontology to standard and standard to ontology

Standard	Ontology/Standard coverage					
	S. Fenz		A. Herzog		Proposed ontology	
	Covered	Covers	Covered	Covers	Covered	Covers
ISO27001	35/311 (11%)	23/133 (17%)	26/460 (6%)	19/133 (14%)	130/1795 (7%)	107/133 (80%)
PCI DSS	42/311 (14%)	48/165 (29%)	25/460 (5%)	32/165 (19%)	132/1795 (7%)	165/165 (100%)
ISSA 5173	31/311 (10%)	7/12 (58%)	29/460 (6%)	6/12 (50%)	15/1795 (1%)	12/12 (100%)
NISTIR 7621	14/311 (5%)	8/10 (80%)	21/460 (5%)	8/10 (80%)	19/1795 (1%)	10/10 (100%)

This ontology and mapping of these 4 security standards to it can be used to generate adaptive maps between any of the two mapped security standards or integrated standard can be created with the usage of any set of mapped security standards without the necessity to map two security standards directly. As our proposed security ontology can cover a larger part of concepts in analysed security standards (the average coverage of these 4 security standards is 92%, while S. Fenz's average coverage of these security standards is 27%, A. Herzog - 20%) the adaptive mapping of security standards will be more precise by applying it as a base ontology. However the ontology does not cover all standards by 100%, therefore should be improved to get even bigger precision of adaptive mapping.

5 Conclusions and Future Works

General comparison of G. Denker's, A. Herzog's and S. Fenz's security ontologies has shown the necessity of user friendly ontology structure - all three ontologies have classes, with more than 25 subclasses in them. Such ontology could be difficult to use for visual presentation or quick knowledge search.

OntoMetric methodology allows a more precise judgment on security ontologies rather than general comparison, because it enables an evaluation of the content of compared ontologies. However the evaluation marks are very dependable on the evaluator's opinion and requirements for the ontology. Evaluation of ontologies' ability to be mapped to security standards is a more suitable measurement to choose the base ontology for adaptive mapping of security standards comparing to OntoMetric.

In order to evaluate ontologies' suitability to map different security standards we compared percentage of concepts in security standard (ISO 27001, PCI DSS, ISSA 5173 and NISTR 7621) which can be mapped to security ontology. This research revealed there are no security ontologies, that would be able to map at least 50% of any security standards we have analyzed. This fact implies the necessity of new or modified ontology, which could be used to present larger parts of knowledge, used in security standards.

We proposed a new security ontology, by integrating concepts of COBIT framework, part of classes of A. Herzog's and S. Fenz's ontologies. This new ontology increased the coverage of security standards. Using this security ontology, from 80% to 100% of analyzed security standards (ISO 27001, PCI DSS, ISSA 5173 and NISTR 7621) can be mapped to it. This percentage can be increased even more with the addition of more specific (related to payment cards, law and standard requirements etc.) concepts to this ontology. The proposed security ontology has a more balanced tree structure as well, which increases its visualization possibilities.

Acknowledgements

The study was carried out within the framework of the National Project No.VP1-3.1-MM-08-K-01-012: "Virtualisation, visualization and e-services security technologies and research", supported by the EU Social Fund.

Bibliography

- [1] Gruber, T (1995). Towards Principles for the Design of Ontologies used for Knowledge Sharing, *International Journal of Human-Computer Studies*, ISSN 1071-5819, 43(5-6): 907-928.
- [2] Dobson, G.; Sawyer P. (2006). Revisiting Ontology-Based Requirements Engineering in the age of the Semantic Web, *In: Dependable Requirements Engineering of Computerised Systems at NPPs*, Institute for Energy Technology (IFE), Halden, 2006.
- [3] Fernandez-Breis, J. T.; Martiinez-Bejar R (2002). A cooperative framework for integrating ontologies, *International Journal of Human-Computer Studies*, ISSN 1071-5819, 56(6): 665-720.
- [4] Gruninger, M.; Lee J. (2002). Ontology Applications and Design, *Communications of the ACM*, ISSN 0001-0782, 45(2): 39- 41.
- [5] Mouratidis, H.; Giorgini P. (2006). *Integrating Security and Software Engineering: Advances and Future Visions*, IGI Global.
- [6] Dhillon, G.; Backhouse J. (2000). Information system security management in the new millennium, *Communications of the ACM*, ISSN 0001-078, 43(7): 125-128.
- [7] Donner, M. (2003). Toward a Security Ontology, *IEEE Security and Privacy*, ISSN 1540-7993, 1(3): 6-7.
- [8] Tsoumas, B.; Gritzalis D. (2006). Towards an Ontology-based Security Management, *Advanced Information Networking and Applications*, ISSN 1550-445X, 1: 985 - 992.
- [9] Gomez-Perez A.; Fernandez-Lopez M.; Corcho O. (2004). *Ontological Engineering*, Springer.
- [10] Ramanauskaite, S.; Goranin, N.; Cenys, A.; Olifer, D. (2013) Ontology-based security standards mapping poptimization by the means of Graph theory, *Proceedings of International congress on engineering and technology ICET 2013*, ISBN 978-80-87670-08-8: 74-83.
- [11] Fenz S. (2010). Ontology-based Generation of IT-Security Metrics, *Proceedings of the 2010 ACM Symposium on Applied Computing*, ISBN 978-1-60558-639-7: 1833-1839.
- [12] Mylopoulos J.; Borgida A.; Jarke M.; Koubarakis M. (1990). Telos: Representing Knowledge About Information Systems, *ACM Transactions on Information Systems*, ISSN 1046-8188: 325-362.
- [13] Landwehr C. E.; Bull A. R.; McDermott J. P.; Choi W. S. (1994). A taxonomy of computer program security flaws, *ACM Computing Surveys*, ISSN 0360-0300, 26(3): 211-254.
- [14] Avizienis A.; Laprie J. C.; Randell B.; Landwehr C. (2004). Basic concepts and taxonomy of dependable and secure computing, *IEEE Transactions on Dependable and Secure Computing*, ISSN 1545-5971, 1(1): 11-33.

- [15] Denker G.; Kagalb L.; Finin T. (2005). Security in the Semantic Web using OWL, *Information Security Technical Report*, ISSN 2214-2126, 10(1): 51-58.
- [16] Mouratidis H.; Giorgini P.; Manson G. (2003). An Ontology for Modelling Security: The Tropos Approach, *Proceedings of the KES 2003 Invited Session Ontology and Multiagent Systems Desing.*
- [17] Giorgini P.; Manson G.; Mouratidis H. (2004). Towards the Development of Secure Information Systems: Security Reference Diagrams and Security Attack Scenarios, *Proceeding of 16th Conference On Advanced Information Systems Engineering.*
- [18] Massacci F.; Mylopoulos J.; Paci F.; Tun T. T.; Yu Y. (2011). An Extended Ontology for Security Requirements, *Advanced Information Systems Engineering Workshops*, ISSN 1865-1348, 83: 622-636.
- [19] Geneiatakis D.; Lambrinouidakis C. (2007). An ontology description for SIP security flaw, *Computer Communications*, ISSN 0140-3664, 30(6): 1367-1374.
- [20] Karyda M.; Balopoulos T.; Gymnopoulos L.; Kokolakis S.; Lambrinouidakis C.; Gritzalis S.; Dritsas S. (2006). An ontology for secure e-government applications, *Proceedings of the The First International Conference on Availability, Reliability and Security, ARES 2006.*
- [21] Undercoffer J.; Joshi A.; Pinkston J. (2003). Modeling Computer Attacks: An Ontology for Intrusion Detection, *The Sixth International Symposium on Recent Advances in Intrusion Detection.*
- [22] Souag A. (2012). Towards a new generation of security requirements definition methodology using ontologies, *Proceedings of 24th International Conference on Advanced Information Systems Engineering*: 1-8.
- [23] Kim A.; Lou J.; Kang M. H. (2005). Security Ontology for Annotating Resources, *On the Move to Meaningful Internet Systems 2005: CoopIS, DOA, and ODBASE* ISSN 0302-9743, 3761: 1483-1499.
- [24] Herzog A.; Shahmehri N.; Duma C. (2007). An Ontology of Information Security, *International Journal of Information Security and Privacy*, ISSN 1930-1650, 1(4): 1-23.
- [25] Fenz S.; Ekelhart A. (2009). Formalizing information security knowledge, *Proceedings of the 4th International Symposium on Information, Computer, and Communications Security*, ISBN 978-1-60558-394-5: 183-194.
- [26] Lozano-Tello A; Gomez-Perez A. (2004). ONTOMETRIC: A method to choose the appropriate ontology, *Journal of database management*, ISSN 1063-8016, 15(2): 1-18.
- [27] ISACA (2013). COBIT 5: A Business Framework for the Governance and Management of Enterprise IT.
- [28] Hofherr M. (2011). Mapping ISO27001 <>PCI DSS 2.0, *ForInSecT*, http://www.forinsect.com/downloads/Mapping-ISO27001-PCI_public.pdf.

RSEA-AODV: Route Stability and Energy Aware Routing for Mobile Ad Hoc Networks

P. Srinivasan, P. Kamalakkannan

P. Srinivasan*

Mahendra Institute of Technology, Namakkal
Tamil Nadu, India

*Corresponding author: salemsrini4u@gmail.com

Dr. P. Kamalakkannan

Department of Computer Science,
Government Arts College, Salem
Tamil Nadu, India

Abstract: Frequent changes in network topology, confined battery capacity of nodes and unreliable nature of wireless channels are the main challenges for reliable routing in Mobile Ad hoc Network (MANET). Selecting a long lasting route is a critical task in MANET. In this paper, we propose a new protocol, Route Stability and Energy Aware Ad hoc On-demand Distance Vector (RSEA-AODV) protocol which is an enhancement of Ad hoc On-demand Distance Vector (AODV) protocol. It designs a bi-objective optimization formulation to compute the reliability factor based on stability and residual energy of nodes. The route with the highest reliability factor value is selected for data transmission. This protocol is compared with other similar routing protocols: PERRA and AODV. We use ns-2 for simulation. Our simulation results show that, the proposed protocol increases the node expiration time by 12 - 32 % and accomplishes 7 - 13 % higher packet delivery ratio compared to PERRA or AODV. The packet delay and control overhead of the proposed protocol is comparable to that of AODV.

Keywords: MANET, reliable routing, route stability, biobjective optimization.

1 Introduction

Mobile ad hoc networks (MANET) are groups of wireless mobile devices, which can communicate with each other without any infrastructure support. It is a self-configured and self-maintained network with no central authority. Every node in MANET acts as both a host and a router. Dynamic topology, limited bandwidth, battery, CPU resources and multi-hop communication are the characteristics that put special challenges in routing protocol design.

Several routing protocols have been proposed for MANETs. Based on the route discovery principle, we can classify them into either proactive or reactive. Proactive routing protocols update routes for every pair of nodes at regular intervals irrespective of their requirement. The reactive or on-demand routing protocols, determine route only when there is a need to transmit a data packet, using a broadcasting query-reply (RREQ-RREP) procedure. Most of these protocols use min-hop as the route selection metric and found that the routes discovered by these protocols are not stable.

The stability based routing protocols [1]- [4] are designed to choose stable route passing through stable links. These protocols improve route lifetimes and packet delivery ratio compare to the shortest path routing protocols. The energy-aware routing protocols [5, 6] are designed considering factors like residual energy, total transmission power or both. These protocols avoid over using of certain nodes and reduce total energy consumptions. But, there exist a very few protocols in the literature [6] - [8] that consider both stability and energy metric during route discovery and maintenance.

In this paper, our objective is to design a routing scheme based on route stability and residual energy metrics during route discovery and maintenance. This scheme compute the link stability based on measurement of received signal strength of successive packets and route stability is computed as the product of link stability of all links that make up the route. This scheme allows nodes that satisfy required energy metric to act as intermediate nodes.

2 Related Work

AODV [9] a reactive routing protocol establish a route to a destination only on demand.. The advantage of AODV is reduced control overhead. But, multiple RREP packets in response to a single RREQ packet can lead to heavy control overhead. Presence of stale routes and unnecessary bandwidth consumption due to periodic beaconing are the drawbacks of AODV.

Link stability is a measure of how stable the link is and how long it will last. Signal strength, pilot signals, link duration distributions, residual energy of the nodes and relative speed between nodes are the parameters used for the computation of link stability. Stability based routing protocols use link stability factor and path stability factor calculated using the above specified parameters to select stable path for data transmission.

The lifetime of a network is one of the important factors to be considered in designing a MANET routing protocol. Maximizing the network lifetime by minimizing the power consumption for the data transfer is the main aim of Energy aware routing protocols [5]. Optimizations carried out in these protocols are classified in the following schemes. (i) Minimize the total energy consumed along the route (ii) Avoid using node with minimum residual energy and (iii) Minimize the total battery cost along the route. They reduce the energy wastage ensuing from retransmission due to bit error rate, frame error rate and link failures due to energy depletion.

In [4], the authors propose Route Stability based QoS Routing (RSQR) where they use the received signal strength of the packets to calculate link stability and route stability. It considerably reduces the number of route breakages by restricting traffic admission. However, the above described protocols do not consider energy metrics during route discovery.

In [7], the authors propose Link-Stability and Energy Aware Routing protocol (LAER). It considers joint metric of link stability and energy drain rate into route discovery. It balances the traffic load on the nodes and considerably decreases the control overhead. However, LAER does not able to discriminate between links of the same age.

In [8], the authors propose Power Efficient Reliable Routing Protocol for Ad Hoc networks (PERRA), a reactive routing protocol, which accounts for both link stability and power efficiency on route discovery. It considerably reduces control overhead by constrained flooding of route requests. However, PERRA does not able to discriminate longer link from lifetime point of view as it does not consider link lifetime and its statistical behavior.

3 Route Stability and Energy Aware (RSEA) Model

3.1 Stability Aware Metric

RSEA model considers signal strengths and mobility for computing the probability of link failures. It computes link stability (LS) using signal strength values received from the MAC layer. Any link e has an associated link stability $LS(e)$ [4] and it is given by

$$LS_{i,j} = \frac{u_2 - DSS_{i,j}}{u_2 - u_1} \quad (1)$$

where DSS is the differentiated signal strength to decide whether the signals are getting stronger or weaker. It is computed as follows.

$$DSS_{i,j} = SS_{cur_{i,j}} - SS_{new_{i,j}} \quad (2)$$

A path between source s and destination d is given as $P(s, d) = (s, e(s, x), x, e(x, y), y, \dots, e(z, d), d)$. Formally, a path between two nodes s and d is a set of all feasible path between them and can be represented as $P(s, d) = P_0, P_1, \dots, P_n$, where each P_i is a feasible path between s and d .

We define the stability of the path P , by the product of link stability of its edges as follows

$$Stability(P) = \prod_{e \in P} LS(e) \quad (3)$$

The path with higher path stability value contains more stable links and choosing it will considerably reduce the probability of link failure.

3.2 Energy Aware Metric

It is assumed that all wireless nodes come with residual power detection device. The energy required to transmit a packet (E_{tx}) [8] can be computed as

$$E_{tx} = \frac{P_{size} * P_{tx}}{BW} \quad (4)$$

where P_{size} is the packet size, P_{tx} is the packet transmitting power and BW is the bandwidth of the link. The transmitting energy is directly proportional to the distance between nodes.

The source application layer communicates n value to the network layer, for selecting nodes that meet the energy requirement. It avoids link breakages due to energy depletion. The total energy required (REQ_e) for data packet transmission is given by

$$REQ_e = n * (E_{tx} + E_{proc}) \quad (5)$$

where E_{proc} is the energy required for packet processing and n is the number of packet .

The energy metric (EM) of the path is given by

$$EM(P) = \prod_{i=1}^n \left(\frac{R_i}{F_i} \right) \quad (6)$$

where $R_i(t)$ is remaining battery capacity and F_i is full battery capacity of intermediate node i , at time t . The goal of this metric is to maximize EM. It takes the product of the residual battery of the intermediate nodes to select a path that has nodes with maximum residual energy among the path that just meet the basic energy requirement REQ_e .

3.3 Problem Formulation

The problem can be stated as, ‘*To find a reliable path for data communication based on route stability and residual energy metrics*’. The above bi-objective optimization problem can be transformed into a single objective problem, by providing importance factor (i.e. W_1 and W_2) for each criterion of the objective. We combine the objectives into a single objective function to calculate the Reliability Factor (RF) of the path P , can be mathematically stated as

$$RF(P) = W_1 \cdot Stability(P) + W_2 \cdot EM(P) \quad (7)$$

where the parameters w_1 and w_2 are chosen based on the network dynamics and application requirements. In this study, in order to give equal importance to both stability and energy metrics, we assign 0.5 to both W_1 and W_2 , such that $W_1 + W_2 = 1$ condition is satisfied.

Consequently, the sum of the objectives has to be maximized and Maximum Reliability Factor (MRF) can be computed by

$$MRF = \max(RF(P_1), RF(P_2), \dots, RF(P_n)) \quad (8)$$

The path with MRF value is selected as a reliable path for data transmission.

4 Route Stability and Energy Aware (RSEA) routing in MANET

4.1 Route Discovery

When a node S needs to send packets to a destination D , it searches for the route in its route table. If a route to the destination D is not available, then the source S will broadcast a route request (RREQ) message to its neighbors. The RREQ of RSEA-AODV is an extension of a RREQ packet of AODV routing protocol. Three new fields Accumulated Path Stability (APS), Accumulated Energy Metric (AEM) and required energy (REQ_e) are added to the RREQ Packet. It initializes the values to the added fields as follows: APS, AEM with 1. The required energy is calculated using equation 5 and is initialized to REQ_e field.

4.2 Route Discovery at Intermediate Nodes

If the strength of the RREQ packet is poor, then it drops the RREQ. Then node i checks whether its residual energy will meet the required energy REQ_e specified in the RREQ packet. If the above conditions are satisfied, then node i make a reverse route entry in the Routing Table (RT). Then it calculates LS. If the signal strength is above $SThr_1$, then it assigns 1 to LS. It implies that nodes are close and link is sufficiently stable. Otherwise, it calculates LS using equation 1. Energy metric is calculated using equation 6.

After these steps, it updates the APS and AEM fields such that the updated values contain the route stability and energy metric of the explored route up to the current node. It enters the relevant information from RREQ into Route Request Forward Table (RFT). Then, Node i broadcast the RREQ to its neighbor. In case of duplicate packet, if it contains better values for APS or AEM, then it makes an entry in RFT and discards the packet. On receiving a RREP packet, node i measures the strength of RREP. If its strength is poor, then it will drop the RREP packet. It looks up RFT for the corresponding RREQ entries, to select the node with the highest RF value. It forwards the RREP packet to the node with the highest RF value. It is shown in Algorithm 2. It makes an entry in the RT.

4.3 Route Selection at Destination

Destination node will receive RREQ packets from different possible routes. On receiving the first RREQ packet, the node D starts a timer t_1 for the duration of Route Reply Latency (RRL) time. It stores all the RREQ that arrives, in its routing table. It computes RF value for the path explored by the RREQ. If destination node receive more than one RREQ before the timer t_1 expires, then it forwards the RREP packet to the node with the highest RF value. It is shown in Algorithm 3. This considerably reduces the amount of control overhead incurred during the route establishment due to multiple RREP for single RREQ as in AODV. The Destination node makes an FT entry for the flow. In case, if it does not receive any data packets within the timeout period, then it will delete the respective entry from the FT table.

Alg. 1 Implemented in Intermediate nodes.

Input: A RREQ packet P from neighbor node.

```

1:   if (Node Battery < EThr) or (SSnew < SThr2) then
2:     Drop Packet P
3:   end if
4:   if ((RREQ not already forwarded) or (RREQ has better APS or AEM value))
5:     then
6:       if (SSnew > SThr1) then LS = 1
7:       DSS = SScur - SSnew
8:       if (SSnew < SThr1) and (SSnew > SThr2) then
9:         if (DSS < u1)
10:          then
11:            LS=1
12:          else
13:            LS = (u2 - DSS)/ (u2 - u1)
14:          end if
15:        end if
16:        APS = APS * LS
17:        AEM = AEM * ((RBCi/FBCi))
18:      Update RREQ with APS and AEM
19:      Broadcast RREQ packet to the next hop
20:    else
21:      Drop Packet P
22:    end if

```

Alg. 2 Implemented in Intermediate nodes

Input: A RREP packet P from neighbor node.

```

1:   Lines 1 - 3 of Algorithm 1
2:   N = number of entry in RFT for the corresponding RREQ.
3:   F = index of the first entry in RFT for the corresponding RREQ
4:   i = F
5:   Count = N
6:   while (Count > 1) do
7:     j = Next entry in RFT for that RREQ
8:     if (w1 * RFT[j].APS + w2 * RFT[j].AEM) > RelFact)
9:       then
10:        i = j
11:       end if
12:     Count = Count - 1
13:   end while
14:   S = RFT[i].PrevHop
15:   Make an entry in RT table
16:   Delete RFT entries for the corresponding RREQ
17:   Forward RREP packet to S
18:   end if

```

Alg. 3 Implemented in Destination Node

Input : A RREQ packet from node N with APS and AEM

```

1: Lines 4- 18 of Algorithm 1
2:   N* = N
3:   MRF: = W1*APS + W2 *AEM
4:   start timer (t1) *for 1st copy of RREQ only *
5:   while (! time-out) do
6:     if (a RREQ arrives at D) then
7:       NEWRF: = w1*APS + w2 *AEM
8:       if ( NEWRF > MRF) then
9:         MRF := NEWRF
10:      N* = N
11:    end if
12:  end if
13: end while
14: Make an entry in RT and FT
15: Send a RREP packet to N*
```

4.4 Route Maintenance

In dynamic mobile ad-hoc networks, link-breaks occur frequently. RSEA-AODV comes with *make-before-break* route maintenance mechanism. This mechanism quickly adapt to the link breakage likely to occur due to the mobility and energy drain. It is depicted in Algorithm 4. It executes the mechanism after every t_2 seconds, to monitor the status of the route established. If an intermediate node is in the critical battery status or it is receiving weaker signal packets, then it creates an HLP message with Time To Live (TTL) set to 1 and broadcasts it to its neighbors. Neighbors on receiving the HLP packet, check for route availability in its routing table for the destination specified in the HLP packet. If the route is available, then it returns the route to the downstream node of the node broadcasting the HLP packet. The downstream node on receiving the route it updates its routing table. It routes the data packets in the new route available, preventing packet losses due to link breakage.

Alg. 4 Route Maintenance by make-before-break mechanism

Input : A packet P from neighbor node

```

1: Executed periodically after timer  $t_2$  expires
2: if ( $SS_{new} < SThr_2$ ) or ( $NodeBattery < Ethr$ ) then
3:   if (intermediate node)
4:     then
5:       Send HLP packet to all its 1 hop neighbors for alternate path
6:     else//if destination node
7:       send Stop – Traffic intimation to the source node
8:     endif
9:   endif
10:  if (timeout)//Alternate route not found
11:    Send RCR to the source node
12:  endif
```

If the node with critical battery is the destination node, then it will send the stop traffic intimation to the source node, to avoid future packet drops and wastage of resources. If there is no alternate route available with the one-hop neighbors, then after the expiry of timer i.e. timeout, the node will send the Route Change Request (RCR) to the source node. The source node on receiving the RCR will go for the re-route discovery to the destination.

4.5 Simulation Parameters

We have simulated a scenario of 50 mobile nodes in a rectangular topology of 1000m * 500m, with each node having a transmission range of 250m. The values of SThr1 and SThr2 are set to $1.5 * RxThr$ and $1.2 * RxThr$, respectively. We simulate it in NS2 2.31. The simulation has been run for 600 seconds. The results show 95% confidence interval on all observed metrics.

The first experiment evaluated the overall performance of RSEA-AODV, PERRA and AODV. During simulation, we generated 12 CBR connections producing 3 packets /second.

Table 1: Simulation Results

Parameters	RSEA-AODV	PERRA	AODV
PDR	96.42	90.54	86.29
COH	0.8342	1.0871	1.1956
Delay(sec)	0.0698	0.0732	0.0659

Table 1 presents the results of the simulations. We observe that the PDR of RSEA-AODV is improved by 6.5% and 13.2 % relative to PERRA and AODV, respectively. This is because of node selection and, route maintenance strategy carried on by the RSEA-AODV. It chooses the intermediate nodes considering the stability and residual energy issues. In case of AODV, it chooses a path with the shortest route. It may contain low energy nodes which lead to disconnections of sessions. The control overhead of RSEA-AODV is reduced by 22 % and 30.2% relative to PERRA and AODV, respectively. This is due to reduction in path reconstruction and constrained flooding of control packets. As RSEA-AODV selects the most reliable path, it considerably reduces the total number of required messages for route reconstruction.

The packet delay for RSEA-AODV is comparable to AODV and PERRA. This is due to following opposing factors. (i) Finding a reliable route, considering the stability and energy metrics, increases the delay on the one hand. (ii) On the other hand, route selection and route maintenance procedure of RSEA-AODV increases the lifetime of the routes and the bottle-neck nodes. It significantly reduces the need for packet retransmissions.

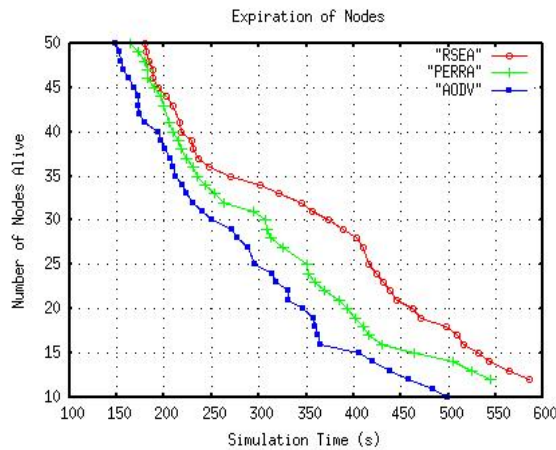


Figure 1: Node ExpirationTime

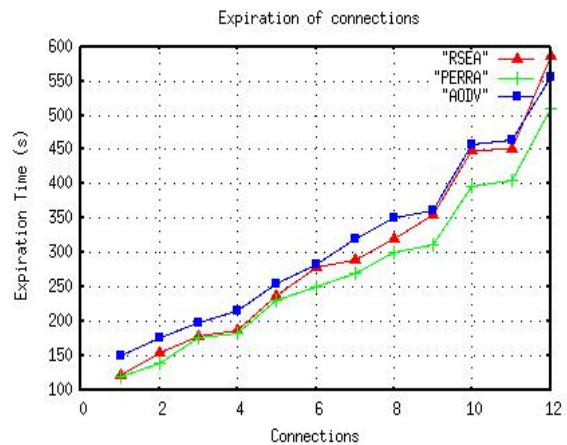


Figure 2: Connection Expiration Time

The plots of NET and CET are presented in Figures 1 and 2. It is observed that RSEA-AODV extends the node lifetime by between 12%-24% over PERRA and by between 20%-32% over AODV. This is because of RSEA-AODV's route maintenance mechanism, in which nodes stop transmitting data traffic if they are about to drain and find an alternate route. AODV exhibits

the worst performance in terms of the nodes expiration time because intermediate nodes continue packet forwarding regardless of the remaining battery until battery depletion.

AODV exhibits a longer lifetime of connections despite a shorter lifetime of nodes. It is noted that, AODVs connection expiration times being anywhere between 25%-5% better than RSEA-AODV or PERRA. This is because, in RSEA-AODV: (i) a source node tries for route reestablishment only fixed number of times after route change request or link failures, and (ii) Nodes reject request for new session once their residual energy is below energy threshold. But AODV keeps retrying and so has a higher chance of finding an alternate route and keeps connection alive for a longer time.

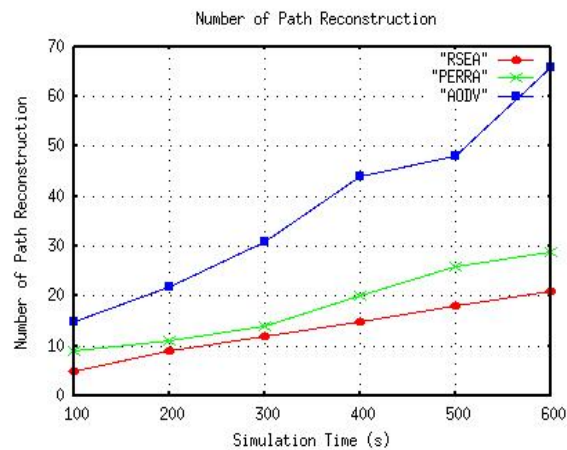


Figure 3: Average number of path reconstructions

Figure 3 shows the number of path reconstructions due to node mobility and energy shortages. RSEA-AODV chooses the more reliable route and thus, reduces the need for route reconstructions. In addition, it decreases the control overhead and offers some energy benefits. The route discovery and maintenance mechanism carried on by PERRA considerably reduced the number of route reconstructions compared to AODV.

In the second experiment, the stability weight (w_1) was provided with the values 0.1, 0.3, 0.5, 0.7 and 0.9 (lowest to highest) and performance of the protocols are observed. The nodes moved at the maximum speed of 5 meters / second. From Figure 4, it is observed the packet delivery ratio of RSEA-AODV was higher than that of PERRA and AODV. It is because it uses more stable routes and presence of enhanced route maintenance mechanism. PERRA showed comparatively better performance than AODV, due to its route discovery mechanism and maintenance of alternative route. As AODV does not discriminate the link in stability issue, it does not show any deviation.

The variance of residual node energy is depicted in Figures 5. This parameter shows the load balancing capability of the protocols. It is noted that as stability weight factor increases, there is an increase in energy variance. As AODV overloads certain nodes and shows a big variance between remaining energies of the node.

Figure 6 shows that there is a gradual increase in the hop count as stability weight increases in both RSEA-AODV and PERRA. The increase in hop count is due to the selection of short and stable link as the stability weight increases. The average physical length of the hops chosen by RSEA-AODV is 45-55% of the transmission range. It is 50-60% in the case of PERRA and 55-65% in the case of AODV.

The third experiment measured the average residual energy of the nodes on the chosen path.

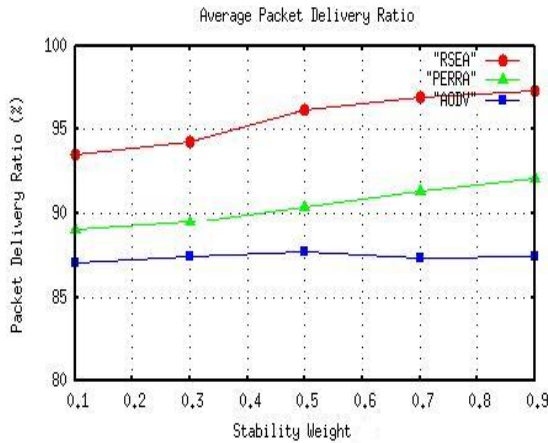


Figure 4: Packet delivery ratio vs. Stability

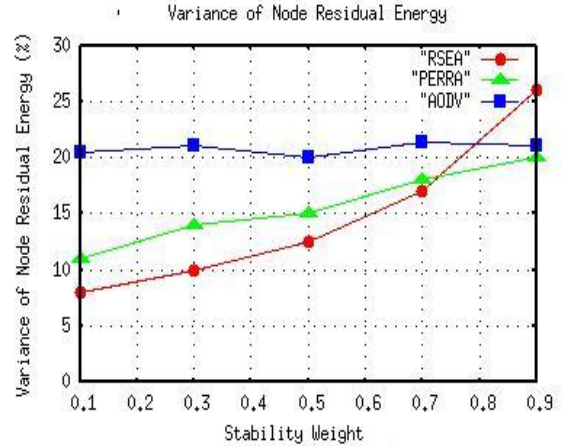


Figure 5: Residual energy vs Stability

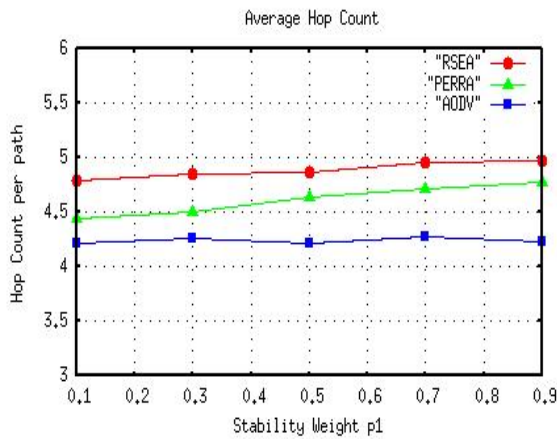


Figure 6: Hop count vs. Stability

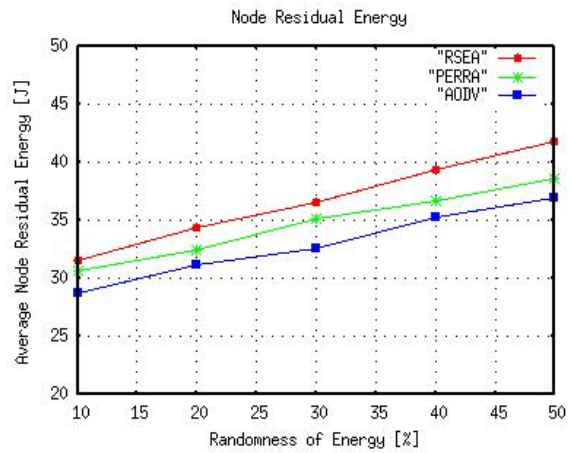


Figure 7: Residual energy vs. Energy

The average residual energy of nodes was set at 30 [J] and with some randomness as in [8]

$$RE_e = E[RE_e] \pm \gamma \quad (9)$$

where $E[RE_e]$ is the residual energy. 10 - 50 % randomness in the residual energy for the nodes was provided in this simulation. For example, in the case of 10 % randomness, each node had 27 - 33 [J] residual energy. The sources energy requirement was 30J. Figures 7 depicts the average energy of the nodes on the chosen path by the protocols. It is noted that RSEA always had higher residual energy compared to PERRA and AODV protocols. This was due to the consideration of nodes residual energy during route selection. But, AODV showed lower residual energy as it goes for the shortest path selection and does not consider residual energy during route discovery.

5 Conclusions and Future Work

In this paper, we proposed a new Route Stability and Energy Aware routing scheme. It establishes route based on the joint metric of link stability and residual energy. During the route discovery process, it constructs the route with nodes that satisfy the sources energy requirement. Among the possible routes, RSEA selects the route with the highest reliability factor for data

transmission. This scheme can be incorporated in any of the existing MANET routing protocols. Using the ns2 simulator, we compared the performance of RSEA-AODV with PERRA and AODV. The simulation results highlight that RSEA-AODV was outstanding in terms of node expiration time, delivery ratio and path reconstruction overhead.

We intend to evaluate the performance of the protocol under different node densities and traffic loads, as part of our future work.

Bibliography

- [1] Abid, M.; Belghith, A. (2011); Stability routing with constrained path length for improved routability in dynamic MANETs, *Personal and Ubiquitous Computing*, 15: 799-810.
- [2] Hui, Z.; Ning, D.Y. (2007); A novel path stability computation model for wireless ad hoc networks, *IEEE Signal Process. Letter*, 14: 799-810.
- [3] Sarma, N.; Sukumar Nandi. (2010); Multipath QoS Routing with Route Stability for Mobile Ad Hoc Networks, *IETE Technical Review*, 27: 380-397.
- [4] Nandi, S.; Sarma, N. (2009); Route Stability Based QoS Routing in Mobile Ad Hoc Networks. *Wireless Pers. Commun, Springer*, 54: 203-224.
- [5] Tan, W.C.W.; Bose, S.K. (2012); Power and mobility aware routing in wireless ad hoc network. *Commun, IET*, 6: 1425-1437.
- [6] Lee, S.; Park. (2008); A routing protocol for Extent Network Lifetime through the Residual Battery and Link Stability in MANET, in *Proc. ACC 08*, 199-204.
- [7] Rango, F. De.; Fazio, P.; (2012); Link Stability and Energy Aware Routing Protocol in Distributed Wireless Networks, *IEEE Trans. on Parallel and Distributed systems*, 23: 713-726.
- [8] Kim, K.J.; Sang, Y. (2005); Power-Efficient Reliable Routing Protocol for Mobile Ad hoc Networks, *IEICE Trans. on Commun.*, 88: 4588-4597.
- [9] Perkins, C.E.; Royer, E.M.; and Das, S.R. (2003); Ad hoc on-demand distance vector (AODV) routing, *IETF RFC 3561*.

Design of Congestion Control Scheme for Uncertain Discrete Network Systems

H. Wang, C. Yu

Hongwei Wang, Chi Yu

Qinhuangdao Branch, Northeastern University
Qinhuangdao, China, 066004
wanghw0819@163.com; yuchi0319@163.com

Abstract: For a class of uncertain discrete network systems, a sliding mode control algorithm is presented for active queue management (AQM) in order to solve the problem of congestion control in transmission control protocol (TCP) communication. First, the sliding surface is designed based on linear matrix inequality (LMI) technique. Then, we analyze the mechanism of chattering for the discrete-time exponential approximation law, a modified one is presented and applied to the network systems. Simulation results demonstrate that the proposed controller has good stability and robustness with respect to the uncertainties of the number of active TCP sessions, link capacity and the round-trip time.

Keywords: sliding mode control, network systems, linear matrix inequality(LMI)

1 Introduction

As the rapid expansion of network scale, congestion control has become an important issue. AQM is a router-based control mechanism, which can implement the end system to demand quality of service. So the combination of TCP and AQM is the main ways to solve the problems of current network congestion control.

Random early detection (RED) as the earliest well-known AQM algorithm is sensitive to parameter variations [2]. So some improved RED methods are presented [3-5]. However, these algorithms can not guarantee high network utilization and low packet loss[6-7]. Recently, some AQM algorithms have been proposed based on mathematical models, which give the basis for control theory research. In [6], a fluid-flow model for TCP/AQM networks has been introduced. The proportional-integral (PI) controller is designed in [8], also some robust control schemes [9-12], such as intelligent PID, variable structure sliding mode controller, H-infinity controller, and so on. These methods can obtain well performance for practical network systems.

With the rapid development of computer technology and digital signal processing chips, the study of discrete-time control theory is rather important. Sliding mode control (SMC) is a robust technique for its unique ability to withstand external disturbance, it has achieved fruitful for continuous system. However, sliding mode control schemes are relatively small for discrete-time system [13-15], a few discrete algorithms are applied to the network control. In this paper, a robust discrete-time sliding mode controller is designed for TCP network model with uncertain disturbance. The aim is to avoid network congestion.

2 Problem Statement and Preliminaries

In [6], a model of TCP connection through a congested AQM router is developed.

$$\begin{cases} \dot{W}(t) = \frac{1}{R(t)} - \frac{W(t)W(t-R(t))}{2R(t)}p(t-R(t)) \\ \dot{q}(t) = \frac{N(t)}{R(t)}W(t) - C(t) \end{cases} \quad (1)$$

where $C(t)$ is the capacity of link, $W(t)$ is the size of TCP congestion window, $q(t)$ is the length of queue in buffer, $p(t)$ is packet-dropping probability function ($0 \leq p(t) \leq 1$), $N(t)$ is the number of active TCP link, $R(t)$ is the transfer delay and $R(t) = T_p + q(t)/C(t)$.

To linearize(1), we first assume $R(t) = R_0$, $N(t) = N$ and $C(t) = C$ is normal value of $R(t)$, $N(t)$ and $C(t)$, the equilibrium point (W_0, q_d, p_0) is defined by $\dot{W} = 0$ and $\dot{q} = 0$. Let $\delta W(t) = W(t) - W_0$, $\delta q(t) = q(t) - q_d$, $\delta p(t) = p(t) - p_0$. A linearized model is given.

$$\begin{aligned} \delta \dot{W}(t) &= -\frac{2N}{R_0^2 C} \delta W(t) - \frac{R_0 C^2}{2N^2} \delta p(t - R_0) \\ \delta \dot{q}(t) &= \frac{N}{R_0} \delta W(t) - \frac{1}{R_0} \delta q(t) \end{aligned} \tag{2}$$

Let $x(t) = \left(\delta q(t) \quad \delta \dot{q}(t) \right)^T = \left(x_1 \quad x_2 \right)^T$, $u(t) = \delta q(t)$, $(-p_0 \leq u(t) \leq 1 - p_0)$. We have

$$\dot{x}(t) = \bar{A}x(t) + \bar{B}u(t) \tag{3}$$

where

$$x(t) = \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix}, \bar{A} = \begin{bmatrix} 0 & 1 \\ -\frac{2N}{R_0^2 C} & -\left(\frac{1}{R_0} + \frac{2N}{R_0^2 C}\right) \end{bmatrix}, \bar{B} = \begin{bmatrix} 0 \\ -\frac{C^2}{2N} \end{bmatrix}.$$

Then the discrete-time uncertain system can be expressed as the sampling period T .

$$x(k+1) = (\tilde{A} + \Delta \tilde{A})x(k) + (\tilde{B} + \Delta \tilde{B})u(k) \tag{4}$$

where $\tilde{A} = e^{\bar{A}T}$, $\tilde{B} = \left(\int_0^T e^{\bar{A}(T-\tau)} \bar{B} d\tau\right)$, $\Delta \tilde{A}$ and $\Delta \tilde{B}$ are depending on network parameters.

In the process of designing controller, the following assumptions are taken.

A₀. The pair (\tilde{A}, \tilde{B}) is controllable, and $\tilde{B} = \begin{pmatrix} \tilde{B}_1 & \tilde{B}_2 \end{pmatrix}$, $\det(\tilde{B}_2) \neq 0$.

A₁. The matrix $\Delta \tilde{A}(k)$ satisfies mismatch condition and $\Delta \tilde{B}$ satisfies $\Delta \tilde{B} = \tilde{B} \times \Delta \hat{B}$.

We do a linear transformation as follows:

$$z = Tx = \begin{bmatrix} I_{n-m} & -\tilde{B}_1 \tilde{B}_2^{-1} \\ 0 & \tilde{B}_2^{-1} \end{bmatrix} x \tag{5}$$

The system (4) is written as follows.

$$z(k+1) = (A + \Delta A)z(k) + B(I + \Delta \hat{B})u(k) \tag{6}$$

where $z(k) = \begin{bmatrix} z_1(k) \\ z_2(k) \end{bmatrix}$, $A = T\hat{A}T^{-1} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$, $\Delta A = T\Delta \tilde{A}T^{-1} = \begin{bmatrix} \Delta A_{11} & \Delta A_{12} \\ \Delta A_{21} & \Delta A_{22} \end{bmatrix}$,

$B = T\tilde{B} = \begin{bmatrix} 0 \\ I_m \end{bmatrix}$, $z_1(k) \in R^{n-m}$, $z_2(k) \in R^m$.

We have

$$z_1(k+1) = A_{11}z_1(k) + A_{12}z_2(k) + \Delta A_{11}z_1(k) + \Delta A_{12}z_2(k) \tag{7}$$

$$z_2(k+1) = A_{21}z_1(k) + A_{22}z_2(k) + \Delta A_{21}z_1(k) + \Delta A_{22}z_2(k) + I_m(I + \Delta \hat{B})u(k) \tag{8}$$

3 Design of Controller for Discrete-time Network Systems

3.1 Designing Sliding Mode Surface

Without loss of generality, we suppose that the sliding surface is

$$s(k) = \bar{M}z(k) = \begin{bmatrix} -M & I_m \end{bmatrix} \begin{bmatrix} z_1(k) \\ z_2(k) \end{bmatrix} = 0 \tag{9}$$

where $\bar{M} \in R^{m \times n}$, $M \in R^{m \times (n-m)}$. Substituting (9) into (7) gives the sliding motion

$$z_1(k+1) = (A_{11} + \Delta A_{11} + A_{12}M + \Delta A_{12}M) z_1(k) \tag{10}$$

A₂. The $\Delta A_{11}(k)$ and $\Delta A_{12}(k)$ satisfy $\Delta A_{11}(k) = DF(k)E_1$, $\Delta A_{12}(k) = DF(k)E_2$, D and $E_i(i=1, 2)$ are constant matrices of appropriate dimensions, $F(k)$ satisfies $F^T(k)F(k) \leq I$.

Lemma 1^[16]. *Given constant matrices D, E and symmetric matrix Y of appropriate dimensions, the following inequality holds $Y + DFE + E^T F^T D^T < 0$. where F satisfy $F(k)^T F(k) \leq I$, if and only if for some constant $\varepsilon > 0$, we have $Y + \varepsilon DD^T + \varepsilon^{-1} E^T E < 0$.*

Theorem 1. *If there exists a symmetric and positive definite matrix P , some matrix W and some scalar ε such that the following LMI(11) is satisfied, then the reduced-order discrete-time system(6) is asymptotically stable by the sliding mode surface(9).*

$$\begin{bmatrix} -X & * & * \\ A_{11}X + A_{12}W & -X + \varepsilon DD^T & * \\ E_1X + E_2W & 0 & -\varepsilon I \end{bmatrix} < 0 \tag{11}$$

where $X = P^{-1}$, $W = MP^{-1}$ and * denotes the transposed elements in the symmetric positions.

Proof: For system (6), we choose the following Lyapunov- Krasovskii function

$$v(k) = z_1^T(k) P z_1(k) \tag{12}$$

Differential equation along the trajectory of the system in (10) is given by

$$\begin{aligned} \Delta v(k) &= v(k+1) - v(k) \\ &= z_1^T(k) (Q^T P Q - P) z_1(k) \end{aligned}$$

where $Q = A_{11} + A_{12}M + DFE_1 + DFE_2M$. If the $Q^T P Q - P < 0$, we can the Theorem 1.

3.2 Design of discrete-time sliding mode controller

Discrete-time approximate law is given in [17].

$$s(k+1) = (1 - \eta T) s(k) - \lambda T \operatorname{sgn} s(k) \tag{13}$$

where $\lambda > 0$, $\eta > 0$, $0 < \eta T < 1$, T is sampling period.

The (13) can not guarantee the system reaches to point. The modified reaching law is given.

$$s(k+1) = (1 - \eta T) s(k) - \left(1 - e^{-|s(k)|} - \eta\right) |s(k)| T \operatorname{sgn}(s(k)) \tag{14}$$

According to the equations (6), (9) and (14), we get the control law as follows.

$$u(k) = (\bar{M}B)^{-1} [\bar{M}Ax(k) + (1 - \eta T)s(k) - (1 - e^{-|s(k)|} - \eta)|s(k)|T \operatorname{sgn}(s(k)) - \bar{M}f(k)] \tag{15}$$

where $f(k) = \Delta Az(k) + \Delta Bu(k)$. it is an unknown number, the controller can not be achieved.

The system existing unknown disturbances of the dynamics is much slower compared with the sampling frequency, so we have

$$f(k-1) = z(k) - Az(k-1) - Bu(k-1) \tag{16}$$

Considering the practical network system, we can obtain the following controller.

$$\begin{aligned} u(k) &= -p_0, u(k) < -p_0 \\ u(k) &= (\bar{M}B)^{-1} [-\bar{M}Az(k) + (1 - \eta T) s(k) - (1 - e^{-|s(k)|} - \eta) |s(k)| T \text{sgn}(s(k)) \\ &\quad - \bar{M}(z(k) - Az(k-1) - Bu(k-1))] , -p_0 \leq u(k) \leq 1 - p_0 \\ u(k) &= 1 - p_0, u(k) > 1 - p_0 \end{aligned} \tag{17}$$

Remark: 1) when $s(k) \geq 0$, the modified reaching law (14) is rewritten as follows.

$$s(k+1) - s(k) = -(1 - e^{-|s(k)|})s(k)T \leq 0$$

2) when $s(k) \leq 0$, the modified reaching law (14) is rewritten as follows.

$$s(k+1) - s(k) = (1 - e^{-|s(k)|})s(k)T \leq 0$$

We can see that the design control law satisfying the sliding mode reaching condition.

4 Simulation Results

Let we choose parameters. $N = 50, C = 300$ packets/s, $R_0 = 0.5$ s, $W_0 = 3, p_0 = 0.22$. $\tilde{A} = \begin{bmatrix} 0.9999 & 0.0098 \\ -0.0263 & 0.9671 \end{bmatrix}, \tilde{B} = \begin{bmatrix} -0.0445 \\ -8.8514 \end{bmatrix}, \Delta\tilde{A} = \begin{bmatrix} 0.0001\sin(0.1\pi k) & 0 \\ 0 & 0.0005\sin(0.1\pi k) \end{bmatrix}, \Delta\tilde{B} = \begin{bmatrix} -0.0013 & -0.2655 \end{bmatrix}, \varepsilon = 0.01$. Using LMI toolbox in the matlab, we can get $M = -8.3$. Then choosing $\lambda = 0.1, T = 0.01$ s, $\eta = 5, x(0) = [30 \ 2]^T$. Fig.1 is the control law response curve based on the proposed control law, which can have much lower packet-dropping probability, satisfying the demand of system response.

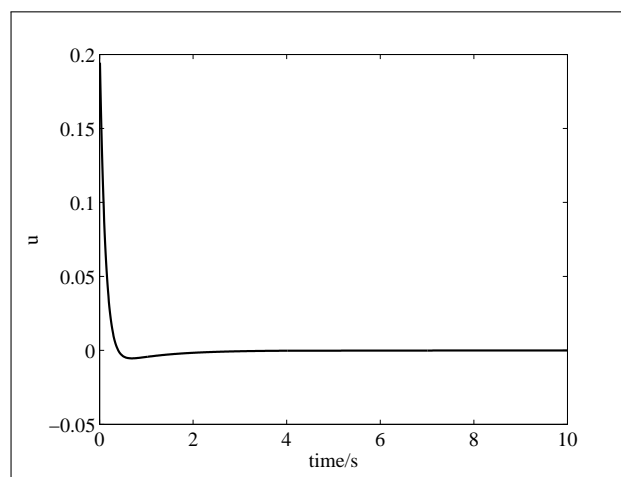


Figure 1: The motion curve with designed control law $u(k)$ in this paper

We give the system state responses curve based on reaching law (13)and (14). The proposed control scheme can obtain much better performance both response time and chattering from Fig.2 and Fig.3.

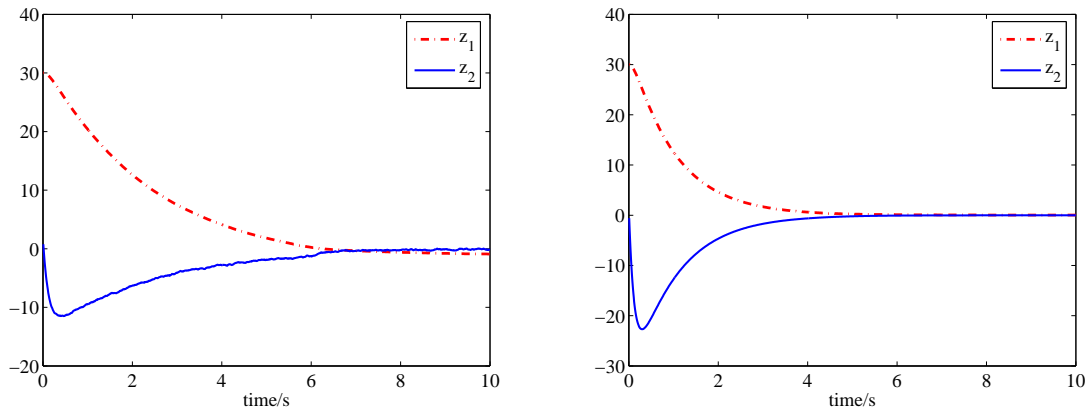


Figure 2: State responses with control law (13) Figure 3: State responses with control law (14)

5 Conclusion

This paper gives a discrete sliding mode control algorithm for network systems. A modified reaching law is presented and applied to the system. Simulation results show that the controller has better stability and robustness, which can get a faster transient response and smaller steady state error. The scheme can effectively avoid network congestion.

Bibliography

- [1] S. Floyd and V. Jacobson, Random early detection gateway for congestion avoidance, *IEEE/ACM Transaction on Networking*, 1(4): 397-413, 1993.
- [2] B. Braden, D. Clark and J. Crowcroft, Recommendations on queue management and congestion avoidance in the Internet, *IETF Request for Comments, RFC 2309*, April 1998.
- [3] D. Lin and R. Morris, Dynamics of random early detection, *Proc. of the ACM SIGCOM'97*, Cannes, 127-137, 1997.
- [4] W. Feng, D. Kandlur and D. Saha, A self-configuring RED gateway, *Proc. of the IEEE INFOCOM*, New York, 1320-1328, 1999.
- [5] T. J. Ott, T. V. Lakshman and L. H. Wong, SRED: stabilized RED, *Proc. of the IEEE INFOCOM*, New York, 1346-1355, 1999.
- [6] V. Misra, W.B.Gong, D.Towsley, Fluid-based analysis of a network of AQM routers supporting TCP flows with an application to RED, *Proc. of the ACM/SIGCOM*, Stockholm, 152-160, 2000.
- [7] C. Hollot, V. Misra, D. Towsley, and W.-B. Gong, On designing improved controllers for AQM routers supporting TCP flows, *Proc. of the IEEE INFOCOM'01*, Anchorage, Alaska, USA, 1726-1734, 2001.
- [8] S. K. Nguang and P. Shi, Fuzzy H-infinity output feedback control of nonlinear systems under sampled measurements, *Automatica*, 39(12): 2169-2174, 2003.

- [9] C. K. Chen, Y.C. Hung, Design of robust active queue management controllers for a class of TCP communication networks, *Information Sciences*, 177: 4059-4071, 2007.
- [10] F.Y. Ren, C. Lin , A robust active queue management algorithm in large delay networks, *Computer Communications*, vol. 28, pp.485-493, 2005.
- [11] W. J. Chang, Robust fuzzy control for uncertain stochastic time-delay Takagi-Sugeno fuzzy models for achieving passivity, *Fuzzy Sets and Systems*, 161: 2012-2032, 2010.
- [12] D. Y. Gu and W. D. Zhang, Design of an H_AT based on PI controller for AQM routers supporting TCP flows, *Proc. Of the 48th IEEE Conference on Decision and Control and 28th Chinese Control Conference*, Shanghai, PR.China, 2009.
- [13] Z.Y. Xi and T Hesketh, On Discrete Time Terminal Sliding Mode Control for Nonlinear Systems with Uncertainty, *Proc. of American Control Conference Marriott Waterfront*, Baltimore, MD, USA, 980-984, 2010.
- [14] K. Abidi, A discrete-time terminal sliding mode control approach applied to a motion control problem, *IEEE Trans. on Industrial Electronics*, 56(9):3619-3627, 2009.
- [15] Y. Q. Xia, M . Y. Fu, Robust sliding-mode control for uncertain time-delay systems based on delta operator, *IEEE Trans. on Industrial Electronics*, 56(9): 3646-3655, 2009.
- [16] L. H. Xie, Output feedback H_AT control of system with parameter uncertainty, *Int J of Control*, 63(4): 741-750, 1996.
- [17] W. B. Gao, Theory and design of nonlinear systems, *Nanjing: Southeast University Press*, 1990.

Author index

Čenys A., 878

Arumugam S., 791

Chan Yi-Cheng, 800

Goranin N., 878

Guan X., 812

Guo Y., 825

He Y., 825

Hu Ya-Yi, 800

Jiang N., 825

Jin S., 825

Kamalakkannan P., 891

Kasnakoglu C., 838

Korać V., 845

Kratica J., 845

Li X., 854

Nekoui M.A., 863

Olifer D., 878

Pakzad M.A., 863

Qu S., 869

Quan Y., 854

Ramanauskaitė S., 878

Sabeen S., 791

Savić A., 845

Srinivasan P., 891

Wang H., 901

Wei Y.M., 869

Wu L.D., 869

Yu C., 901

Yu R.H., 869

Zhang S., 812

Zhang Z., 812