

ISSN: 2071-2987 (Online version) 2223-523X (Print version)

International Journal of Design, Analysis and Tools for Integrated Circuits and Systems



Preface

Welcome to the Volume 7 Number 1 of the International Journal of Design, Analysis and Tools for Integrated Circuits and Systems (IJDATICS). This volume is comprised of research papers from the International Conference on Recent Advancements in Computing, Internet of Things (IoT) and Computer Engineering Technology (CICET), October 29-31, 2018, Taipei, Taiwan. CICET 2018 is hosted and organized by The Tamkang University amid pleasant surroundings in Taipei, which is a delightful city for the conference and traveling around.

CICET 2018 serves a communication platform for researchers and practitioners both from academia and industry in the areas of Computing, IoT, Integrated Circuits and Systems and Computer Engineering Technology. The main target of CICET 2018 is to bring together software/hardware engineering researchers, computer scientists, practitioners and people from industry and business to exchange theories, ideas, techniques and experiences related to all aspects of CICET.

Recent progress in Deep Learning has unleashed some of the promises of Artificial Intelligence (AI), moving it from the realm of toy applications to a powerful tool that can be leveraged across a wide number of industries. In recognition of this, CICET'18 has selected Artificial Intelligence and Machine Learning as this year's central theme.

The Program Committee of CICET 2018 consists of more than 150 experts in the related fields of CICET both from academia and industry. CICET 2018 is hosted and organized by The Tamkang University, Taipei, Taiwan and supported by:

- Research Institute of Big Data Analytics, Xi'an Jiaotong-Liverpool University, China
- IoT Research Centre, Xi'an Jiaotong-Liverpool University, China
- Swinburne University of Technology Sarawak Campus, Malaysia
- Baltic Institute of Advanced Technology, Lithuania
- Taiwanese Association for Artificial Intelligence, Taiwan
- VersaSense, Belgium
- International Journal of Design, Analysis and Tools for Integrated Circuits and Systems
- International DATICS Research Group

The CICET 2018 Technical Program includes 2 keynotes and 24 oral presentations. We are beholden to all of the authors and speakers for their contributions to CICET 2018. On behalf of the program committee, we would like to welcome the delegates and their guests to CICET 2018. We hope that the delegates and guests will enjoy the conference.

Professor Ka Lok Man, Xi'an Jiaotong-Liverpool University, China and Swinburne University of Technology Sarawak, Malaysia

Dr. Woonkian Chong and Dr. Owen Liu, Xi'an Jiaotong-Liverpool University, China

Chairs of CICET'18

CICET 2018 Organization

Honorary Chairs

- Jian-Nong Cao, Hong Kong Polytechnic University, Hong Kong
- Han-Chieh Chao, National Dong Hwa University, Taiwan

Keynote Speakers

- Steven Guan, Research Institute of Big Data Analytics and Xi'an Jiaotong-Liverpool University, China
- Hui-Huang Hsu, Tamkang University, Taiwan

Advisory Board

- Hui-Huang Hsu, Tamkang University, Taiwan
- Paolo Prinetto, Politecnico di Torino, Italy
- Massimo Poncino, Politecnico di Torino, Italy
- Joongho Choi, University of Seoul, South Korea
- Michel Schellekens, University College Cork, Ireland
- M L Dennis Wong, Heriot-Watt University, Scotland
- Vladimir Hahanov, Kharkov National University of Radio Electronics, Ukraine
- Chun-Cheng Lin, National Chiao Tung University, Taiwan

General Chairs

- Ka Lok Man, Xi'an Jiaotong-Liverpool University, China and Swinburne University of Technology Sarawak, Malaysia
- Woonkian Chong, Xi'an Jiaotong-Liverpool University, China
- Owen Liu, Xi'an Jiaotong-Liverpool University, China

Local Chair

• Chien-Chang Chen, Tamkang University, Taiwan

Industrial Liaison Chair

• Gangming Li, Xi'an Jiaotong-Liverpool University, China

Publicity Chairs

- Vincent Ng, The Hong Kong Polytechnic University, Hong Kong
- Neil Y.(Yuwen) Yen, The University of AIZU, Japan
- Patrick HangHui Then, Swinburne University of Technology Sarawak, Malaysia

Program/Workshop Chairs

- Tomas Krilavičius, Baltic Institute of Advanced Technologies and Vytautas Magnus University, Lithuania
- Seungmin Rho, Sungkyul University, South Korea
- Sheung-Hung Poon, University of Technology Brunei, Brunei Darussalam
- Chuck Fleming, Xi'an Jiaotong-Liverpool University, China
- Yujia Zhai, Xi'an Jiaotong-Liverpool University, China

Program Committee

- Alberto Macii, Politecnico di Torino, Italy
- Wei Li, Fudan University, China
- Emanuel Popovici, University College Cork, Ireland
- Jong-Kug Seon, System LSI Lab., LS Industrial Systems R&D Center, South Korea
- Umberto Rossi, STMicroelectronics, Italy
- Franco Fummi, University of Verona, Italy
- Graziano Pravadelli, University of Verona, Italy
- Yui Fai Lam, Hong Kong University of Science and Technology, Hong Kong
- Jinfeng Huang, Philips & LiteOn Digital Solutions Netherlands, The Netherlands
- Jun-Dong Cho, Sung Kyun Kwan University, South Korea
- Gregory Provan, University College Cork, Ireland
- Miroslav N. Velev, Aries Design Automation, USA
- M. Nasir Uddin, Lakehead University, Canada
- Dragan Bosnacki, Eindhoven University of Technology, The Netherlands
- Milan Pastrnak, Siemens IT Solutions and Services, Slovakia

- John Herbert, University College Cork, Ireland
- Zhe-Ming Lu, Sun Yat-Sen University, China
- Jeng-Shyang Pan, National Kaohsiung University of Applied Sciences, Taiwan
- Chin-Chen Chang, Feng Chia University, Taiwan
- Mong-Fong Horng, Shu-Te University, Taiwan
- Liang Chen, University of Northern British Columbia, Canada
- Chee-Peng Lim, University of Science Malaysia, Malaysia
- Salah Merniz, Mentouri University, Constantine, Algeria
- Oscar Valero, University of Balearic Islands, Spain
- Yang Yi, Sun Yat-Sen University, China
- Damien Woods, University of Seville, Spain
- Franck Vedrine, CEA LIST, France
- Bruno Monsuez, ENSTA, France
- Kang Yen, Florida International University, USA
- Takenobu Matsuura, Tokai University, Japan
- R. Timothy Edwards, MultiGiG, Inc., USA
- Olga Tveretina, Karlsruhe University, Germany
- Maria Helena Fino, Universidade Nova De Lisboa, Portugal
- Adrian Patrick ORiordan, University College Cork, Ireland
- Grzegorz Labiak, University of Zielona Gora, Poland
- Jian Chang, Texas Instruments, Inc, USA
- Yeh-Ching Chung, National Tsing-Hua University, Taiwan
- Anna Derezinska, Warsaw University of Technology, Poland
- Kyoung-Rok Cho, Chungbuk National University, South Korea
- Yuanyuan Zeng, Wuhan university, China
- D.P. Vasudevan, University College Cork, Ireland
- Arkadiusz Bukowiec, University of Zielona Gora, Poland
- Maziar Goudarzi, Sharif University of Technology, Iran
- Jin Song Dong, National University of Singapore, Singapore
- Dhamin Al-Khalili, Royal Military College of Canada, Canada
- Zainalabedin Navabi, University of Tehran, Iran
- Lyudmila Zinchenko, Bauman Moscow State Technical University, Russia
- Muhammad Almas Anjum, National University of Sciences and Technology (NUST), Pakistan
- Deepak Laxmi Narasimha, University of Malaya, Malaysia
- Danny Hughes, Katholieke Universiteit Leuven, Belgium
- Jun Wang, Fujitsu Laboratories of America, Inc., USA
- A.P. Sathish Kumar, PSG Institute of Advanced Studies, India
- N. Jaisankar, VIT University. India
- Atif Mansoor, National University of Sciences and Technology (NUST), Pakistan
- Steven Hollands, Synopsys, Ireland
- Siamak Mohammadi, University of Tehran, Iran
- Felipe Klein, State University of Campinas (UNICAMP), Brazil
- Eng Gee Lim, Xi'an Jiaotong-Liverpool University, China
- Kevin Lee, Murdoch University, Australia
- Prabhat Mahanti, University of New Brunswick, Saint John, Canada
- Kaiyu Wan, Xi'an Jiaotong-Liverpool University, China
- Tammam Tillo, Xi'an Jiaotong-Liverpool University, China
- Yanyan Wu, Xi'an Jiaotong-Liverpool University, China

- Wen Chang Huang, Kun Shan University, Taiwan
- Masahiro Sasaki, The University of Tokyo, Japan
- Shishir K. Shandilya, NRI Institute of Information Science & Technology, India
- J.P.M. Voeten, Eindhoven University of Technology, The Netherlands
- Wichian Sittiprapaporn, Mahasarakham University, Thailand
- Aseem Gupta, Freescale Semiconductor Inc., Austin, TX, USA
- Kevin Marquet, Verimag Laboratory, France
- Matthieu Moy, Verimag Laboratory, France
- RamyIskander, LIP6 Laboratory, France
- Chung-Ho Chen, National Cheng-Kung University, Taiwan
- Kyung Ki Kim, Daegu University, Korea
- Shiho Kim, Chungbuk National University, Korea
- Hi Seok Kim, Cheongju University, Korea
- Brian Logan, University of Nottingham, UK
- AsokeNath, St. Xavier's College (Autonomous), India
- Tharwon Arunuphaptrairong, Chulalongkorn University, Thailand
- Shin-Ya Takahasi, Fukuoka University, Japan
- Cheng C. Liu, University of Wisconsin at Stout, USA
- Farhan Siddiqui, Walden University, Minneapolis, USA
- Katsumi Wasaki, Shinshu University, Japan
- Pankaj Gupta, Microsoft Corporation, USA
- Masoud Daneshtalab, University of Turku, Finland
- Boguslaw Cyganek, AGH University of Science and Technology, Poland
- Yeo Kiat Seng, Nanyang Technological University, Singapore
- Tom English, Xlinx, Ireland
- Nicolas Vallee, RATP, France
- Rajeev Narayanan, Cadence Design Systems, Austin, TX, USA
- Xuan Guan, Freescale Semiconductor, Austin, TX, USA
- Pradip Kumar Sadhu, Indian School of Mines, India
- Fei Qiao, Tsinghua University, China
- Chao Lu, Purdue University, USA
- Ding-Yuan Cheng, National Chiao Tung University, Taiwan
- Pradeep Sharma, IEC College of Engineering & Technology, Greater Noida, GB Nagar UP, India
- Ausra Vidugiriene, Vytautas Magnus University, Lithuania
- Lixin Cheng, Suzhou Institute of Nano-Tech and Nano-Bionics (SINANO), Chinese Academy of Sciences, China
- Yue Yang, Suzhou Institute of Nano-Tech and Nano-Bionics (SINANO), Chinese Academy of Sciences, China
- Yo-Sub Han, Yonsei University, South Korea
- Hwann-Tzong Chen, National Tsing Hua University, Taiwan
- Michele Mercaldi, EnvEve, Switzerland

Table of Contents

Vol. 7, No. 1, October 2018

Preface	i
Table of Contents	vi

1.	Rico Thomanek, Christian Roschke, Benny Platte, Robert Manthey, Maik	1
	Benndorf and Marc Ritter, Audio- and Location-Based Interface Generator	
	Using Laravel and iBeacon (ALBIGen)	

- 2. Takeshi Nagata, A Multi-agent System for Voltage Control of Power Distribution 7 System Corresponding to Large-scale Renewable Energy Sources
- Difeng Yu, Juntao Zhu, Wenge Xu, Hai-Ning Liang, Charles Fleming and Yong Yue, An Investigation of Micro- and Macro-Interaction for 3D Manipulation using Dual-Hand Controller in Virtual Reality Environments
- 4. Lim Kuoy Suong, Yun Borin, Lee Sunwoo and Kwon Jangwoo, Abnormal 14 Pedestrian Behavior Detection System Using Multiple Deep Convolutional Neural Networks
- Benny Platte, Rico Thomanek, Tony Rolletschke, Christian Roschke and Marc 20 Ritter, Person Tracking and Statistical Representation of Person Movements in Surveillance Areas
- 6. Weihao Liu, Haoyang Wang, Limin Yu, Mark Leach and Fei Ma, *Wireless* 26 Sensor Network Traffic Modeling and Anomaly Simulation based on OPNET
- 7. M L Dennis Wong and Jia Jun Tay, A Low Multiplicative S-Box for a Stochastic 30 Random Number Generator
- 8. P. Chomtip, V. Adhisaya, L. Kanhokthorn and L. Phoptorn, *Banana (Musa 34 acuminata Triploid AAA, Cavendish) Sweetness Measurement by Digital Image Processing Technique*
- Dingkun Li, Keun Ho Ryu, Erdenebileg Batbaatar, Hyun Woo Park, Seon Phil 38 Jeone and Zhou Ye, An Effective Feature Selection and Classification Model for High Dimensional Big Data Sets
- I. Bumbuliene, J. Mandravickaite, A. Bielinskiene, L Boizou, J. Kovalevskaite, 44
 E. Rimkute, L. Vilkaite-Lozdiene and T. Krilavicius, *RNNs for Lithuanian Multiword Expressions Identification*
- 11. Qi Chen, Wei Wang and Xin Huang, Long Short-Term Memory Encoder- 48 Decoder for Traffic Flow Prediction
- 12. Yujia Zhai, Yuanye Fang, Zhejian Zhang, Sanghyuk Lee and Kejun Qian, *Design* 52 of an Intelligent Temperature Control for the MIMO Thermal System
- 13. Chengkai Yu, Charles Fleming and Hai-Ning Liang, *Scale Invariant Privacy* 56 *Preserving Video via Wavelet Decomposition*
- 14. Yanda Zhu and Yuxuan Zhao, Comparative Studies of Segmentation Algorithms 59
- 15. Gangmin Li and Bei Yao, *Classification of Genetic Mutations for Cancer* 63 *Treatment with Machine Learning Approaches*

- 17. Yi-Jen Su, Chao-Ho Chen, Tsong-Yi Chen and Cheng-Chan Cheng, *Chinese* 71 *Microblog Sentiment Analysis by Adding Emoticons to Attention-Based CNN*
- Shih-Hao Chang, Chih-Chieh Hung and Mao-Sheng Hung, Visible Light Optical 76 for Indoors Toxic Gas Detection and Positioning System
- Tengfei Qian and Ou Liu, Detection of Intentional and Unintentional Financial 81 Restatements using Data Mining Techniques
- 20. S. -I. Kang, S. B. Kim and S. M. Lee, *Improved GMDA based DOA Technique* 85 using Pre-training Phase Unwrapping for Source Localization
- 21. Jongmin Lee, Kwangho Lee, Kidong Yun, Mucheol Kim, Geuchul Park and Chanyeol Park, *HyDM: Data Migration Methodology for Hybrid Memories targeting Intel Knights Landing Processor*
- 22. C. C. Chen, J. M. Yang and H. Q. Liu, *An Index Table Based Optimally Matched* 94 *Reversible Image Watermarking Scheme*
- 23. Mucheol Kim, Junho Kim, Jongmin Lee, Geunchul Park and Chanyeol Park, 98
 PCA based Performance Analysis with System Profiling Data in Many-core system
- 24. Jieming Ma, Blended Learning Design for Computer Programming Courses 103

Audio- and Location-Based Interface Generator Using Laravel and iBeacon (ALBIGen)

Rico Thomanek, Christian Roschke, Benny Platte, Robert Manthey, Maik Benndorf and Marc Ritter

Abstract-Digital signage systems are an important human machine interface in public institutions. By distributing services in IoT environments, dynamically generated and personalized content can be displayed. Currently existing digital signage systems are mostly proprietary and difficult to maintain and administer. In addition, such systems require the installation and configuration of an application on the client side, which can lead to high deployment and management costs. This paper describes a framework architecture for creating such systems, based on free and standardized technologies. The focus of the framework is on distributed and web-based communication of all components. The use of the MVC-Framework Laravel to create a basic system and its extension by plugin functionalities allows a platform independence, maintainability and extensibility. In the proposed architecture, external interaction mechanisms such as voice control, location-based services via iBeacons and bidirectional server-client connections are added to the basic system. The framework enables the development of personalized digital signage systems and connects them with Internet of Things components, such as sensors and actuators. The implementation of a prototype based on the framework and the results show the advantages of the proposed architecture.

Index Terms—Digital Signage System, Internet of Things, location-based services, iBeacon, Laravel

I. INTRODUCTION

THIS paper describes the framework "audio- and locationbased interface generator" (ALBIGen) for creating and managing digital information and guidance systems. The framework is based on Laravel and current web technologies and enables the visualization of personalized information via a smart user interface using iBeacon technology in combination with a smartphone. In addition, it is possible to connect Internet of Things technologies, especially in the smart home sector. This allows information to be created and displayed dynamically as soon as an event is triggered by a sensor or an actuator. Newer architectures for information and control systems often exist as proprietary solutions based on predefined hardware and software. Due to missing and non-standardized

R. Thomanek is with the Hochschule Mittweida, Mittweida, Germany, phone: +49 3491 58 1407; e-mail: rthomane@hs-mittweida.de.

C. Roschke is with the Hochschule Mittweida, Mittweida, Germany, phone: +49 3491 58 1146; e-mail: roschke@hs-mittweida.de).

B. Platte is with the Hochschule Mittweida, Mittweida, Germany.

R. Manthey is with the TU Chemnitz, Chemnitz, Germany.

M. Benndorf is with the Hochschule Mittweida, Mittweida, Germany.

M. Ritter is with the Hochschule Mittweida, Mittweida, Germany.

interfaces, extensions and user-specific adaptations are only inadequately realizable. The development of web-based systems for displaying specific information based on an MVC framework has already been discussed in several papers. Recent work in [1], [3], [5], [9] describes such architectures for displaying information using distributed systems. None of these developments focuses on linking such systems with voice control and location-based services using iBeacons to display personalized content. Furthermore, neither the administration nor the development effort was considered in the context of minimization. The architecture of the framework described in this paper allows a platform-independent development and easy maintenance of a digital signage system. By using standardized technologies, dynamic content can be automatically generated and displayed in a user-oriented manner. In addition, content can be imported from the system via standardized interfaces (XML, JSON). Moreover, iBeacons enable the framework to visualize personalized information. In order to guarantee the platform independence of the management and playout devices, a web-based architecture concept was implemented. The purpose of this architectural concept is to minimize the administrative effort when adding new interfaces. This is realized by a dynamic RESTful web service. As a result of permanent web socket connections between interface and web server, each interface can display individual content via HTTP requests. These requests can be made manually or automatically by IoT components.

II. SYSTEM ARCHITECTURE

We selected a web-based architecture as the basic system because it offers several advantages over native applications. On the client side, a browser installed on most operating systems can be used and there is no need to install additional tools. Another advantage is the automatic and persistent backup of data in a central database. Each transaction of a user is timeexactly related to the current data and is stored time exactly to these. In the domain of web development there are various frameworks to support the development process. As a basis we chose the MVC Framework Laravel for ALBIGen, because it offers the following advantages: [4]

1. Logic can be distributed as required so that calculations can be performed on the server, client and database side and a simple load distribution is possible.

- 2. All developed systems can easily expand and maintain due to the clear separation of logics. In addition, changes can be made and displayed directly at runtime.
- 3. Created modules can be easily reused, parameterized and updated.
- 4. Many required functions are already integrated and simplify the development process.

In order to display personalized content on a display, we use iBeacons to identify people. The user's smartphone detects the presence of an iBeacon and sends control commands to the associated display. This device is connected to the server via a unidirectional connection and can receive and evaluate events and display personalized content. In addition to this approach, we use speech recognition to customize the content. Voice commands can be used to trigger actions such as highlighting a defined element. The Framework is also coupled with sensors and actuators from the smart home domain by a management system. This allows commands to be sent to a defined client after an actuator has been triggered or a value read from a sensor exceeds or falls below a predefined threshold value.

A. Laravel as a basic framework

The opensource framework Laravel follows the development pattern Model View Controller (MVC) and is licensed under MIT. Laravel offers multiple components in the standard installation, which can be used by developers for larger and more complex developments. In addition to the integration of Laravel modules, Symfony components can also be integrated. [1] Laravel supports true code modularization through a combination of drivers and bundles. Drivers are interfaces to functional components such as cache, session, database and authentication. Bundles allow code sections to be packed and reused. In addition, Laravel offers possibilities to design, extend and modify databases. Queries can be described in general using the Fluent Query Builder. The system then translates the data and forwards it to the linked database. MySQL, PostgreSOL, MSSQL and SOLite are supported. Communication with the database is made possible by the Eloquent module, which is an implementation of the Active Record Pattern. Eloquent allows database entries to be created, queried, updated and deleted without creating SQL queries. To simplify the creation of required components, the framework has the command line tool Artisan. This allows developers to initialize database migrations, run unit tests and create new components. Artisan can be extended with own functions. [4]

1) The architecture of a Laravel application

Laravel has a routing module that allows to forward HTTP requests to a controller. It can be differentiated which request method will be used and which functions are to be called. This allows different views to be delivered on the base of the URL and to transmit parameters for requests. The controller processes the transferred parameters and initiates the use of a model to perform operations on the database. The controller is used to pass content from the database to a view. Data entered in a form can be checked, database entries queried, views loaded, and files uploaded into the system. Laravel View component uses the blade template engine to create views. Developers can write blade templates and integrate PHP output into HTML code. The Laravel Model allows you to represent a database connection as a PHP object. [11]

As shown in figure 1, when interacting with a Laravel-based application, an HTTP client transmits a request that is received by an HTTP server and forwarded to the Laravel routing engine. The router sends this request to a controller. The route used must be defined beforehand. The controller interacts with models and prepares data for the view. The view delivers the formatted data to the HTTP client using any markup language.



2) Route Parameters

We added route parameters to the system and linked them to the creation of an interface instance. Each time an interface view is called, a persistent entry is created in the database and a fixed ID is assigned to the requesting device. This allows any number of devices to start an individual view. After the view has been successfully created, the corresponding route is blocked for the requesting system on the basis of the IP. Due to the fixed IP, specific control commands can be sent to individual devices later on. The advantage of this method is the minimization of the administration effort, since not every device has to be explicitly entered in a configuration file and a direct assignment is made after a successful request.



Figure 2 illustrates the sample route structure that the client calls to display a view. The client sends any ID to the server using a given GET parameter. The server reads this parameter and generates an individual interface and transmits it to the client. This can lead to several devices requesting the same ID and therefore no unambiguous assignment can be carried out. Therefore, after a successful request, the system blocks the ID and additionally secures the routes with a user authentication provided by Laravel. Only requests for which the correct authentication key is transmitted in the POST body are accepted and forwarded.

3) Events

For the communication between client and server after the view was transmitted, we decided to use Server-sent events. With this

technology it is possible to establish and maintain a persistent unidirectional connection between the client and the server. The EventSource API used is a standardized part of HTML5, so to open a connection to the server to receive events from it, an EventSource object with the reference to a server script must be created. When the connection is established, the client starts listening to message events. As soon as a message has been received, it can be integrated into the current view. The serverside script sending events must respond with the MIME type text/event-stream. Each notification is sent as a text block and closed with line breaks. The event stream is a simple stream of UTF-8 encoded text data. Messages are separated by line breaks. Each message is a combination of the fields listed in Table I. [7] For our own events we use the event field to categorize the events and to execute specific functions. We transmit parameters via the corresponding data field. The ID is irrelevant in our context and as retry value we have chosen one second.

 TABLE I

 SERVER-SENT EVENTS MESSAGE FIELDS

FIELD NAME	DESCRIPTION
EVENT	Describes the type of event, in the form of
	a string. Specifying sends an event to the
	listener for the specified event name.
DATA	Contains the data of the message.
ID	Event ID for unique assignment.
RETRY	The time to use to send the event.

4) The holistic web-based architecture



Fig 3. The holistic web-based architecture.

Figure 3 shows the Laravel-based architecture we developed. A plugin manager has been added to the classic Laravel components. This allows the administration and integration of plugins and thus the expansion of the framework by additional functionalities. This allowed us to develop components for iBeacon communication and to add personalized content to the standard views. In addition, we have integrated a system on the client side that simplifies communication with the server. The Communication Handler manages the transfer of information and forwards it to the controller for processing. The controller generates a generic exchange format and makes it available to the UI or an API instance. This makes it easy to adequately display any information sent by the server, or to manipulate it on the client side. The architecture is based on the client-server principle. On the client side, a human user can interact with the client controller via UI and a machine via API. The client controller forwards requests to the interface to the server application, the communication handler. The handler handles sending and receiving messages via AJAX, direct connections, server-sent events and web sockets. The server consists of a main system and several plugins that can be added as required. The main system includes the complete management logic, including methods to manage users, protect application areas, generate views and integrate plug-ins. The use of plugins makes it easy to extend the functional core and avoids the risk of negatively affecting the core logic. The system router of the main system receives all requests sent by the client and passes them either to the system controller or to plug-in router instances. The system controller contains the application logic and uses models from the system model handler to generate data structures that are made available to the system view generator. The system model handler forms the interface between the application logic and the database. A model is generated for each table and made available as an object for the application logic. The System View Generator uses the data structures created in the controller and converts them into a standardized exchange format such as HTML, JSON or XML. The generated formats are then transferred to the Communication Handler of the HTTP client in response. Besides the variant to pass requests to the system controller, it is also possible to pass them on to a plugin. Each plugin integrated into the system consists of a router, model handler, view generator and controller like the main system. Each request is forwarded from the plugin router to the plugin controller. The data structures created in the Plugin Controller using the Plugin Model Handler are then converted into any exchange format using the Plugin View Generator and integrated into the view of the System View Generator.

B. iBeacon

iBeacon is a technology introduced by Apple in 2013 that can be used for indoor navigation. The iBeacon devices emit a low energy Bluetooth network for this purpose. Using the Proximity Beacon Advertising Package (PBA package) sent as a broadcast, mobile applications are able to detect their position on a microlocal level and provide the user with contextual content depending on the location. The structure of a PBA package is shown in the following table. [2] According to Table II, the following three parameters of the PBA package must be adapted for location detection and contextual display of content: Proximity UUID, Major and Minor.

1) Usage in the framework:

ALBIGen enables the visualization of personalized information using iBeacon technology. For this purpose, an iBeacon is placed near the interface, the values for "Proximity UUID", "Major" and "Minor" have to be adjusted as follows:

- UUID: The 16 byte flag is used as a service-specific identifier. With this ID a mobile application is able to detect the iBeacon.
- Major, Minor: Major and Minor are unsigned

integers between 0 and 65535 used to identify iBeacons with greater accuracy than UUID alone. ALBIGen uses the major value to identify the location (e.g. building number). The minor value is still used to distinguish the interface.

 TABLE II

 PROXIMITY BEACON ADVERTISING PACKET

BYTES	NAME	VALUE	NOTES
0	Flags[0]	0x02	Data length in first
			AD structure 2 bytes
1	Flags[1]	0x01	AD type
2	Flags[2]	0x06	Bit 0 (ON): LE
			Limited
			Discoverable Mode
			Bit 2 (ON): BR/EDR
			Not Supported
3	Length	0x1A	Data length 26 bytes
4	Туре	0xFF	Data type
			manufacturer
			specific data
5-6	Company ID	0x004C	Manufacturer data
			0x004C == Apple
7-8	Beacon Type	0x0215	iBeacon
			advertisement
			indicator
9-24	Proximity	0xnnnn	Set user UUID
	UUID		
25-26	Major	0xnnnn	Set major value
27-28	Minor	0xnnnn	Set minor value
29	Measured	0xnn	Signal power value
	Power		

Figure 4 visualize this example with use of Major and Minor values.



Fig. 4. URL composition for displaying personal data.

2) Extended final state machine

We have prototypically developed an iOS app to detect the iBeacon on a smartphone. The iOS enables the detection of iBeacons using two mechanisms, monitoring and ranging. Monitoring enables the detection of iBeacons even if the application is closed. However, only the detection of entering and leaving the iBeacon area is supported. If the application is not active in the foreground during a detection, iOS starts the app in the background for a few seconds to handle the event. Time consuming actions can therefore not be carried out. Ranging, on the other hand, can only be used if the application is active or has only recently been in the background. Ranging also make it possible to carry out time-consuming actions. The ranging mechanism can also be used to determine the distance between the smartphone and iBeacon. [8] We use both mechanisms to display personalized information, monitoring when the app is closed and ranging when it is in the background

or active. When using Ranging, the PBA packet byte "Measured Power" is also used to determine the distance to the interface. This enables us to only display personalized information if the user is not more than 3m away from the interface, for example.



Fig. 5 The extended finite state machine for iBeacon detection.

The programming for handling the iBeacon detection (entering, leaving) and the resulting actions was implemented in the iOS app as an extended finite state machine (EFSM). This means that the detection of multiple iBeacons does not lead to multiple HTTP requests, because each state can only be exited after defined input signals. Figure 5 shows the actions for displaying personal information based on the enter and exit events of iBeacon tracking. The URL is generated in the same way as in figure 4. To create the session token, the mobile device automatically performs authentication and then adds the session token received as a POST parameter to the URL request. Therefore, the choice of the authentication method is independent of ALBIGen.

C. Speech control

In addition to control using iBeacon technology, we also implement voice control. The used Web Speech API is a specification of the Speech API Community Group and enables the use of functions for speech recognition by JavaScript at the client. The procedure is specified. The speech recognition functions are available via the Speech Recognition class. Speech Recognition can be initiated by creating and configuring a Speech Recognition object using the start method. During configuration you can specify the language or a grammar. Whenever a result is generated, an event is triggered to provide the recognized text, alternatives and data on their confidence. It is not defined how the functions of the API are made available. We use Google Chrome and the Google Speech API. Spoken words are sent to external servers, analyzed by neural network models and the results are returned as text. The API recognizes over 110 languages and provides their

transcription. In the framework we have currently integrated voice commands to highlight certain elements, reorganize the view and display specific content, such as video streams, room plans and more.

D. Remote Control

The ability to send control commands to ALBIGen via HTTP makes it possible to use any REST client or systems with HTTP client support for remote control. The remote-control systems can therefore be used completely autonomously by ALBIGen. The only requirement in our configuration is the support of HTTP-POST-Request. Based on this architectural concept, any IoT management system can also perform automated control processes.

For the definition of the action to be performed, its identifier is transmitted as a post parameter. We use a post-request with the following parameters: APIKEY: Authentication for the ALBIGen API Action-Type: Defines the action to be performed (e.g. showLessons) Start Action Inferface-ID https://<IP>/api/registeraction/14-Stop Action https://<IP>/api/unregisteraction/14-To visualize whether the action type is activated or deactivated, its state can be checked .: APIKEY: Authentication for the ALBIGen API Action-Type: Defines the action to be checked (e.g. showLessons) get action state Inferface-ID https://<IP>/api/statusregisteraction/14-



For remote control of ALBIGen we have used for example the Smart Home System "ZWay". ZWay is a software solution available for various operating systems and can be used as a ZWave controller. Z-Wave is an international ITU-T standard for home automation. [6], [10] ZWay also offers user management and a convenient user interface that can be used to create rules for automation and virtual devices (e.g. switch) for sending HTTP requests. Especially the connection of sensors and actuators is an interesting application. For example, a smoke detector can automatically transmit an HTTP request to ALBIGen when it detects smoke, in order to display information in case of an emergency (e.g. escape plan). Two URLs are available for displaying or hiding information on the interfaces using HTTP requests. (see figure 6) For visual representation of the current setting of an action, its status can be checked via an additional URL. As a prototype we used a ZWay plugin "HTTP Device" for remote control of the interfaces. This plugin creates a switch in the user interface that allows different HTTP requests to be sent when switching on or off. In our configuration, this enables or disables the display of a video stream at the defined interface.

III. RESULT AND DISCUSSIONS

Based on the developed framework we were able to develop a prototype and integrate it on four clients. Two clients in the main building and two in the library of the University of Applied Sciences Mittweida. The main building displays the campus plan, a video stream, the restaurant plan, the news and the current room occupancy. Furthermore, current Livestreams can also be displayed. The library also displays information on book locations. The elements of all clients can be addressed and changed by voice commands. By using iBeacon technology, users can also view their own lesson plan as they are located within three meters of a client. Furthermore, it is possible to display any messages on the entire display, which allows alarm notices and evacuation plans to be displayed. The system has been running stable for 2520 hours without interruption. As shown in figure 7, 5187 dynamically personalized content was generated and displayed over 15 weeks. On average there were 346 requests to all clients generated by one of the integrated interaction options.



By using a ticket system set up for the system, we enable every user to report any bugs that may occur and to provide feedback. At present, we have not found any serious errors using this system.

A. Advantages of the framework

The developed prototype showed that the developed architecture of the framework offers several advantages over conventional proprietary solutions. The system is modular, expandable and allows control by other software using a standardized REST API. Furthermore, the development is open source and can therefore be integrated into heterogeneous infrastructures without further costs. The modular structure and abstraction of the logic by means of plugins also enables simple extensibility and platform independence. This is also supported by the MVC framework used, as this enables the logic to be divided into model, view and controller. Furthermore, Laravel offers a stable and well documented basis. By distributing the logic, it is possible to regenerate content dynamically and make adjustments at runtime. We were able to minimize the administration and maintenance effort by using generic routes and the instantiated, controllable interface instances and by a central and automated control. Bidirectional communication between mobile devices and the individual clients also enables the display of personalized, target group-oriented content. In this way, the development of an inexpensive system and a high target group accuracy can be realized in the context of the displayed information.

B. Deficits and room for improvement

Currently, we see deficits and potential for improvement in the detection time for detecting the iBeacon area and in the use of voice control in public areas.

1) iBeacon Discovery Time

To avoid unnecessarily reducing the battery life of the smartphone, the operating system automatically selects how short the polling times for finding iBeacons should be. This time depends on several factors, as for example the number of active apps, the number of apps running in the background or when the iBeacon app was last opened. Therefore, it cannot be guaranteed that the personalized information will be displayed immediately when entering the iBeacon area. So, the waiting time can last up to 10-15s. Fast detection of an iBeacon can only be guaranteed if the app has activated the ranging mechanism. However, the ranging mechanism is only activated after entering an iBeacon area. This means only after an iBeacon has been detected by monitoring. The placement of several additional iBeacons in the environment of the interface with a different UUID may therefore enable the app to be set in the ranging mechanism before entering the interface area. This would also allow the app to quickly find the iBeacon used to display the personalized information. However, this procedure only works if the app is active or running in the background. Currently, only one personalized information can be displayed per interface and time slot. If there are several people in front of an interface, there is competition between these people. To present personalized information, the principle "first come, first served" applies. This means that an interface is blocked for the duration of the display of personalized information, for further requests. In this case, the app receives a time range from when the interface is available again (see figure 4). If the person is still in front of the display within this time range, a further attempt is made to display the personalized information. This procedure continues until the server sends a confirmation message to the smartphone or the person has left the iBeacon area. If a person does not leave the iBeacon area after displaying the information, the request is resent after a definable interval, which we have currently set to 3 minutes.

2) Speech control

The control of the interface by voice commands is strongly dependent on the application environment. The microphones built into standard computers only allow voice operation in quiet environments. For use in public areas and the associated ambient noise, special microphone arrays are required to reduce the effects of reverberation and noise and to identify the source of the voice using beamforming.

IV. CONCLUSIONS

The developed framework enabled us to develop a digital signage system at Mittweida University. The system has already been stable in use for several months and was able to display specific information on several clients. The developed architecture proved to be extremely stable, low-maintenance and easily expandable. Currently, the system has displayed specific information for about 5200 users for 16 weeks as expected and we could not detect any interruptions. From this it can be concluded that the framework can be used to develop and individually adapt reliable digital signage systems for public institutions. The identified disadvantages are currently being processed successively in further iteration processes and the framework is being continuously improved. In the future, the framework could be extended with additional functionalities such as a manageable warning system, employee information systems and indoor navigation by using iBeacon location determination.

REFERENCES

- Anif, M., Dentha, A. & Sindung, H. (2017) Designing internship monitoring system web based with Laravel framework. IEEE International Conference on Communication, Networks and Satellite (Comnetsat).
- [2] Apple (2015) Proximity Beacon Specification Release R1. https://developer.apple.com/ibeacon. (6th April 2018)
- [3] Ardiyani, R., Arham, Z. & Rustamaji, E. (2016) The Development of a Web-based Spatial Information System Utilization of Forest Area (Case Study : Sulawesi Island). International Conference on Cyber and IT Service Management.
- [4] Chen, X., Ji, Z., Fan, Y. & Zhan, Y. (2017) Restful API Architecture Based on Laravel Framework. Journal of Physics: Conference Series, Volume 910, conference 1.
- [5] Chin, T., Chuang, Y, Fan, Y, Jiang, Y, Kang, Y, Kuo, W, To, T & Nishino, H. (2017) Prototyping digital signage systems with high-low tech interfaces. In SIGGRAPH Asia 2017 Posters (SA '17), ACM, New York, NY, USA, Article 18.
- [6] G.9959 (2015) Short range narrow-band digital radiocommunication transceivers PHY, MAC, SAR and LLC layer specifications. ITU-T.
- [7] Hickson, I. (2015) Server-Sent Events W3C Recommandation. https://www.w3.org/TR/eventsource. (06.04.2018)
- [8] Koehne, M. & Sieck, J. (2014). Location-based Services with iBeacon Technology. Second International Conference on Artificial Intelligence, Modelling and Simulation, S. 315-321
- [9] Park, Y., Yang, H., Dinh, T. & Kim, Y. (2017) Design and implementation of a container-based virtual client architecture for interactive digital signage systems. International Journal of Distributed Sensor Networks, Volume 13
- [10] SDS11847 (2017) Z-Wave Plus Device Type Specification. Sigma Designs.
- [11] Yu, H. R. (2014) Designing internship monitoring system web based with Laravel framework. International Conference on Computer Science and Electronic Technology (ICCSET 2014).

A Multi-agent System for Voltage Control of Power Distribution System Corresponding to Large-scale Renewable Energy Sources

Takeshi Nagata, Member, IEEE

Abstract—In this paper, we propose a multi-agent based distribution feeder voltage control method by combination of power factor controls and switching between feeders. The proposed system is composed of two-layer multi-agent and consists of three types of agent: feeder agent (FA), feeder-group agents (FGAs), and bus agents (BAs). At first, FA carries out a self-contained voltage control in the feeder, and then FGA tries to interconnect between the feeders. The voltage control strategies are implemented as the class definition of Java into the system. In order to verify the performance of the proposed method, it has been applied to a model of power distribution system. The simulation results show that the system is able to control the voltage to make effective use of the distributed generations. In addition, by changing the system configuration, it shows a possibility of reducing the amount of the power factor control.

Index Terms— Voltage control, distribution system, renewable energy source, multi-agent

I. INTRODUCTION

oday, the introduction of renewable energy source (RES) is A advanced worldwide to realize a low-carbon society. According to the International Energy Agency (IEA), the introduction volume of PV in Japan is 7.80 million kW in 2016 and the total generation capacity is 42.40 million kW [1]. Considering the situation where PV is introduced in large quantity, voltage fluctuation due to the output fluctuation of PV occurs in various places, and as a result, the influence on the voltage profile of the distribution line feeder is expected to be extremely complicated. That is, it is also expected that situations such as an upper limit voltage deviation occur in the feeder 1 and a lower limit voltage deviation occurs in the feeder 2. Under such circumstances, it is impossible to take countermeasures with LRT (load ratio transformer) which uniformly raise and lower the sending voltage. In addition, it is expected that the tap management of the SVR (step voltage regulator) introduced in the middle of the distribution line becomes complicated and the new installation of the SVR is also necessary.

Numerous studies on voltage control have been carried out so

far [2-11]. When these methods are classified by paying attention to control equipment, they can be roughly divided into three types. The first method is based on advanced use of LRT and SVR [2-7]. In the second method, equipment such as SVC (static reactive power compensator) is used directly to inject reactive power into the distribution system [8, 9]. The third method uses the PSC (power conditioner system) attached to photovoltaic power generation apparatus to control its power factor [10, 11]. Among these methods, the first method uses equipment developed in the past era without reverse power flow, so it is impossible to respond to the new environment of the power system. In the second method, there is a problem in the selection of installation position of the apparatus and installation cost. Especially, in the reference [10], it was demonstrated on a practical level that prototype PCS with reactive power generation function added could suppress the voltage fluctuation without increasing PCS capacity by demonstration test. However, there is still a study of cooperation scheme for PCS installed in multiple places in the distribution line feeder.

Therefore, in this paper, focusing on the third method, we propose a multi-agent based distribution feeder voltage control method by combination of power factor controls and switching between feeders. In constructing the system, we adopt the multi-agent system (MAS) [12] which is one of the systems which can easily construct the distributed system.

The proposed system is a two-level MAS. The MAS has "feeder agents (FAs)" corresponding to the distribution line feeders, "feeder group agent (FGA)" for managing a plurality of FAs, and "bus agents (BAs)" corresponding to customers or PV power plants. This method is a voltage control method that keeps the voltage profile of the bus proper by cooperative operation of these agents.

II. PROPOSED VOLTAGE CONTROL METHOD

A. Outline of the Proposed Method

This paper deals with two distribution line feeders connected to the distribution substation. The basic idea of the proposed method will be explained using the simple distribution line feeder shown in Fig. 1. In this figure, SS is a distribution substation, Feeder #1 and #2 are distribution line feeders, B1 to B14 are buses corresponding to residential areas, factories,

Takeshi Nagata is with the Hiroshima Institute of Technology, 2-1-1, Miyake, Saeki-ku, Hiroshima, 731-5193, JAPAN (corresponding author to provide e-mail: t,nagata.wp@ it-hiroshima.ac.jp).

offices, etc. CS1 to CS3 are feeder interconnection switches, the black arrow shows the load demand, and the white arrow shows the power output of the distributed power source.

In this method, by changing the power factor of the distributed power source indicated by the white arrow in Fig. 1, the voltage profile of the connected feeder is appropriately maintained. Fig. 2 shows the concept of voltage control by power factor control of a distributed power supply. In this method, when the voltage deviates from the lower limit, the power factor of the distributed power source is controlled, and lagged reactive power is supplied from the power supply to the power system side to attempt deviates from the upper limit, the power factor of the distributed power source is controlled, and lagged reactive power is supplied from the power supply to the power factor of the distributed power source is controlled, and lead reactive power is supplied from the power source to the power system side to attempt deviation elimination.

The features of this proposed method are as follows.

- This system is a two-layer voltage control MAS in which FA performs self-contained voltage control within the feeder and FGA manipulates CS (interconnection switch between feeders) as necessary.
- In FA and FGA, multiple algorithms for optimizing voltage profiles are implemented in the form of knowledge modules. In other words, they are "voltage improvement modules" for voltage control of customers and PV buses, and "power factor improvement modules" for improving power factor of distributed power supply.
- When the reactive power in the feeder is insufficient, FGA receives the connection request from FA and closes the designated CS, thereby enabling the power feed from the other feeder.
- Each knowledge module predicts the situation after control execution by solving the power flow problem on "virtual power distribution system". The appropriate control value is then determined by repeating this calculation. Solving this power flow problem is hereinafter referred to as "voltage estimation calculation" in this paper. Since the controlled variable (power factor) in this method is determined by trying this voltage estimation calculation a plurality of times, implementation of complicated control laws and tuning of control parameters are unnecessary.
- In the current interconnection regulation in Japan, it is an interconnection with a fixed power factor of 0.85 or more, however in the proposed method, when deviation from the control target voltage occurs, this power factor would be dynamically changed. As a result, this method can be advantageous to consumers as compared with the case of interconnection with fixed power factor.

Generally, the distribution system is roughly divided into a radial type and a loop type, and the radial type is widely adopted. Comparing the two types, the loop type configuration has less voltage drop and higher reliability; however problems such as complicated protection system are pointed out. In this method, as an initial study of the control method of multiple feeders, the loop operation is adopted for CS operation.

B. Implementation with Multi-agent

Here, a method of realizing the voltage control method by the MAS will be described. The proposed system consists of three types of agent: one feeder group agent (FGA), two feeder agents (F0, F1), and several bus agents (B1 to B14). BA corresponds to large-scale solar power plants, factories, offices, and residential districts such as housing complexes. Also, BA can have both load and distributed power supply.

The functions and operations of each agent will be described below.

(a) Feeder agent (FA): FA performs voltage improvement and power factor improvement of interconnected bus voltage. In this system, the lower limit of the power factor of the dispersed power source of the interconnecting BA is coordinated with 0.900. First, in voltage improvement, when the interconnected bus voltage deviates from the allowable range, the voltage control strategies are applied. For the evaluation of each bus voltage, the allowable voltage width (*V*-*Limit*) which corresponds to 101 ± 6 V, and the operational target voltage width (*V*-*Target*) are used. If the voltage is not within the *V*-*Target* even if the power factor is controlled to 0.900, it means that the reactive power in the feeder is insufficient.

Next, at the stage of improving the power factor, it is examined whether or not the power factor improvement is possible for the operation performed with the above voltage improvement. This is because there is a possibility that the power factor can be improved according to changing supply and demand situation.

(b) Feeder group agent (FGA): FGA closes one of CS 1 to CS 3 in Fig. 1 when receiving even one connection request from FA. The location of interconnection is determined by voltage estimation calculation. Specifically, paying attention to the voltage of BA (Target-Bus) in which the voltage violation cannot be resolved within the feeder, FGA selects the CS that can achieve the largest voltage improvement.





Fig. 2. Conceptual diagram of voltage control by power factor (P/sqrt(P'2+Q2)).

(c) Bus agent (BA): Based on the request from FA, BA returns the bus voltage of the current time, the demand and the generated power (active power and reactive power) at the next control time to FA. Here, the active power and reactive power of the generated power are calculated from the current power factor. Also, the output of the distributed power supply is changed based on the instruction value (power factor) from FA.

(d) Agent processing: The agent is implemented as a state transition machine. Fig. 3 shows the state transition diagram of the agent. The circle in the figure represents the state. An important thing in constricting a multi-agent system is to simplify the structure of the agent. By doing so, it becomes possible to reduce the size of the object and to implement a large number of agents in the system without using a high-performance computer. State 0 (S0) corresponds to initialization at startup, state 1 (S1) corresponds to decision making, and state 3 (S3) corresponds to action on the environment. In this method, all agents are composed of four states as shown in this figure.

C. Voltage Control Method

FA controls the power factor of the distributed power supply when the bus voltage of BA deviates from the allowable voltage range. Voltage control of this method is realized by using a knowledge module. The knowledge module is defined in the Java class and implements six modules (S2_V1 to S2_V3, S2_R1 to S2_R3) as shown in Fig. 4. This figure shows that FGA and FA are realized by state transition of such a knowledge module. In Fig. 4, S2_V1 to S2_V3 are voltage improvement modules, and S2_R1 to S2_R3 are power factor improving modules.



Fig. 3. State transition diagram of agent.



Fig. 4. State transition diagram of agent.

III. SIMULATION

In order to demonstrate the validity of the proposed system, a multi-agent simulation system was developed using the Java language, and simulations were performed.

A. Simulation Conditions

In this simulation, the distribution system model shown in Fig. 5 was used. The voltage of this model is 6,600 V, the line impedance between the buses is 0.14 + j 0.35 ohms/km, and the total length of the distributed line is 10.5 km. In this figure, SS is a substation, R1 to R4 are housing developments and other residential areas, F1 to F2 are factories, O1 is an office complex, and PV1 to PV7 are large-scale PV power plants.

The following values were used for the voltage tolerance.

- Permissible voltage range (V_Limit): 0.941-1.059 p.u.
- Operation target voltage range (V_target): 0.98-1.020 p.u.

It is assumed that the power factor of the distributed power source is set to 1.0 at the start of the simulation and the power factor can be controlled up to 0.900 per 0.001 based on the request from FA. The simulation control time interval was set to 60 seconds, and the simulation was performed for 24 hours (1,440 steps). Three cases will be described below.

B. Case Studies

Case-1 is the case without control; Case-2 is the case with power factor control and without CS control. Case-3 is a case where the power factor control and the CS control are performed at the same time.

(a) Case-1 (without control)

First, the trend of each bus voltage without control is shown in Fig. 6. In this figure, the horizontal axis represents time and the vertical axis represents voltage (p. u.). As shown in the figure, Feeder # 1 connected to B1 to B7 deviates from the lower limit (*V-Limit*) in many time zones. In addition, Feeder # 2 connected to B8 to B14 deviates from the upper limit (*V-Limit*) in many time zones because PV generated power is excessive.

(b) Case-2 (with power factor control; without CS control)

In this case, only the power factor control is performed under the same condition as in Case-1. The trend in each bus voltage is shown in Fig. 7. As shown in the figure, all bus voltages are within *V-Limit* at all times. However, in Feeder # 1, it deviates from the operation target lower limit (*V-Target*) in some time zones.



Fig. 5. A test system.

(c) Case-3 (with power factor control; with CS control)

In this case, both of the power factor control and the CS control are performed under the same condition as in Case-1. The trend in each bus voltage is shown in Fig. 8. As shown in the figure, all the bus voltage is within the operation target voltage (V-Target) at all times. In this case, FGA received a connection request from Feeder # 1 at 7:23 and closed CS2, the system configuration was changed. After that, we received a disconnection request from Feeder # 1 at 18: 33 and opened CS2.

IV. CONCLUSION

In this paper, a voltage control method of the distribution feeder combining the power factor control of distributed power supply and the interconnection switch (CS) is studied. This system is a two-layer multi-agent system that performs self-contained voltage control within the feeder and then performs interconnection control between feeders. From the simulation results, it shows the possibility of reducing the power factor control amount by changing the system configuration in addition to effective use of distributed power supply.

Future tasks include comparative examination with radial system configuration, evaluation on practical scale distribution system using real data, and examination of combined system with other system.







Fig.7. Simulation results (Case-2: with control (without CS)).



REFERENCES

- International Energy Agency (IEA), Available: http://www.iea-pvps.org. [1]
- Y. Liu and X. Qui, "Optimal Reactive Power and Voltage Control for [2] Radial Distribution System", in Proc. IEEE Power Engineering Society Summer Meeting, pp 85-90, 2000.
- [3] K. Kabemura, K. Yonekura, T. Tsukamoto, K. Hashimoto and M. Hara, "Application of Dispersed Autonomous Voltage Control System to a Real High Voltage Distribution Netwaork", IEEJ Trans. PE, Vol.122-B, No.12, pp. 1252-1261, 2002.
- J-young Park, S. Nam, and J-keun Park, "Control of a ULTC Considering [4] the Dispatch Schedule of Capacitors in a Distribution System", IEEE Trans. on Power Systems, Vol. 22, No. 2, pp. 755-761, 2007.
- [5] Y. Nakachi, A. Kato and H. Ukai, "Voltage and Reactive Power Control Taking into Account of Economy and Security by using Tabu Search", IEEJ Trans. PE, Vol.128, No.1, pp. 305-311, 2008.
- M. Oshiro, T. Senjyu, A. Yona, N. Urasaki and T. Funabashi, "Voltage [6] Control in Distribution Systems Considered Reactive Power Output Sharing", IEEJ Trans. PE, Vol.130, No.11, pp. 972-980, 2010 .
- [7] S. Yoshizawa, Y. Yamamoto, Y. Hayasi, S. Sasaki, T. Shigeto and H. Nomura, "Dynamic Updating Method of Optimal Control Parameters of Multiple Advanced SVRs in a Single Feeder", IEEJ Trans. PE, Vol.135, No.9, pp. 550-558, 2015.
- G. W. Kim and K. Y. Lee, "Coordination Control of ULTC Transformer [8] and STATCOM Based on an Artificial Neural Network", IEEE Trans. on Power Systems, Vol. 20, No. 2, pp. 550-558, 2015.
- [9] D. Iioka, K. Sakakibara, Y. Yokomizu, T. Matsumura and N. Izuhara: "Distribution Voltage Rise at Dense Photovoltaic Power Generation Area and its Suppression by SVC", IEEJ Trans. PE, Vol.126, No.2, pp. 153-158, 2006.
- [10] N. Uchiyama, H. Miyata, T. Ito and H. Konishi, "Reactive Power Control Method for Reducing Fluctuation in Large-scale Photovoltaic Systems", IEEJ Trans. PE, Vol.130, No.3, pp. 297-303, 2010.
- [11] S. Kawasaki, N. Kanemoto, H. Taoka, J. Matsuki and Y. Hayashi, "Cooperative Voltage Control Method by Power Factor Control of PV Systems and LRT", IEEJ Trans. PE, Vol.132, No.4, pp. 309-316, 2012.
- [12] C. Rehtanz, "Autonomous Systems and Intelligent Agents in Power System Control and Operation", Springer-Verlag Berlin Heidelberg, 2003.

An Investigation of Micro- and Macro-Interaction for 3D Manipulation using Dual-Hand Controller in Virtual Reality Environments

Difeng Yu, Juntao Zhu, Wenge Xu, Hai-Ning Liang*, Charles Fleming, and Yong Yue

Abstract—In this research, we conduct a two-fold user study to investigate micro- and macro-interaction for 3D manipulation using dual-hand controllers in virtual environments. The chosen device is the HTC VIVE controller due to its richness and variety of features and similarity to other dual-hand devices for virtual reality systems. The first study evaluates whether the HTC VIVE Controller supports basic micro- and macro-interactions and the second study aims to find out whether micro- or macro-interactions are more frequently used to perform common manipulation tasks.

Index Terms—Bimanual interaction, dual-hand controller, HTC VIVE, micro- and macro-interaction, user-elicitation study, virtual reality, 3D manipulation.

I. INTRODUCTION

Virtual Reality (VR) is an emergent technology. It provides users the sense of immersion within virtual environments and enables the user to interact with the 3D objects inside them. For example, by using HTC VIVE (one of the most common VR systems), participants are able to grab, throw, rotate, and translate the objects in a predefined 3D virtual environment (see Fig. 1a).

This research investigates two types of interaction for manipulating 3D objects using the HTC VIVE Controller in VR environments: micro- and macro-interactions. Micro-interaction requires relatively subtle, small movements such as finger and wrist movements [1, 3], while macro-interaction usually requires large physical movement like full-arm movement [4]. Given the rapid introduction of commercial VR systems like the HTC VIVE and Oculus RIFT systems, it is important to know whether the system supports basic micro- and macro-interactions and how people tend to perform manipulations when interacting with 3D objects in virtual environments.

Difeng Yu, Juntao Zhu, Wenge Xu, Hai-Ning Liang, Charles Fleming, and Yong Yue are with the Department of Computer Science and Software Engineering, Xi'an Jiaotong-Liverpool University.

Hai-Ning Liang is also with Research Institute of Big Data Analytics, Xi'an Jiaotong-Liverpool University.

*Corresponding author: Hai-Ning Liang; HaiNing.Liang@xjtlu.edu.cn.



Fig. 1. (a) The VR environment setup for the study; (b) HTC VIVE Headset; (c) HTC VIVE Controller; and (d) a user is performing the task with HTC VIVE VR System and Noitom Perception Neuron.

In this paper, we present two users studies built on each other and provide our discussions based on the results.

II. USER STUDY 1

Our main objective for this study is to evaluate whether users are able to perform some basic micro- and macro-interaction using HTC VIVE VR systems.

A. Participants and Apparatus

Thirty participants (11 females) between the ages of 18-22 were recruited from a local university campus to take part in this study. According to the results of a pre-experiment questionnaire, 2 participants were left-handed.

The experiment was conducted on an Intel Core i7 processor PC with an NVIDIA GTX 1070 graphics card. The program was developed in C#.NET and was run within the Unity3D platform. An HTC VIVE headset was used to immerse the user into the 3D virtual world and the HTC VIVE Controller was used as the input device for the user to interact with objects in the virtual environment (Figure 1a-c). We used the Noitom Perception Neuron, a body worn tracking device, to capture the body, head, and hand motions of the users (Figure 1d).

B. Task, Procedure, and Experiment Design

Participants were asked to perform a series of gestures via the HTC Vive Controller to manipulate 3D objects. Seven tasks (see Table 1) were predefined in this experiment [2]. For each

This work was supported in part by XJTLU Key Program Special Fund (KSF-P-02), and XJTLU Research Development Fund.

 TABLE I

 The 3D tasks given to participants in the user study 1

Manipulation	Animation Descriptions	Widget
Throwing	(1) Throw the Object forward	
Translation	(2) Move the Object along X-axis	У
	(3) Move the Object along Y-axis	V
	(4) Move the Object along Z-axis	
Rotation	(5) Rotate the Object along X-axis	
	(6) Rotate the Object along Y-axis	•
	(7) Rotate the Object along Z-axis	

task, participants would first select the object and then perform three types of gestures: (1) Wrist Gesture: moving wrist to perform the task while keeping elbow and shoulder fixed; (2) Elbow Gesture: moving elbow to perform the task while keeping wrist and shoulder fixed; and (3) Shoulder Gesture: moving shoulder to perform the task while keeping wrist and elbow fixed. Note that the first type of movements was considered as micro-interaction while the second and third types of movements were counted as macro-interaction.

Before the experiment started, participants were asked to fill in a pre-experiment questionnaire with their demographic information and were given time to familiarize themselves with the virtual environment and the HTC VIVE Controller. Then, they would proceed to carry out the manipulation tasks one by one (with the order of presentation balanced). After the experiment, participants were instructed to provide some comments on three types of gestures they had performed. To assess whether users are able to perform these tasks, we record the whole movements using the Noitom Axis Neuron Software and analyze the completion rate result after.

C. Results and Subjective Feedback

Based on our analysis, nearly all 3D manipulation tasks defined can be completed by the participants (see Table 2). However, one exception was performing the rotation task along the Z-axis using Wrist and Shoulder. According to video recordings, we found that participants could actually finish the task if properly guided but they just could not figure out how to do it by themselves. This might indicate that rotation along the Z-axis using Wrist and Shoulder are unnatural interactions and sounds unfamiliar to users.

Most participants suggested that some gestures made them feel uncomfortable when performing although they were able to complete the tasks. This motivated us to conduct the second user study—to explore natural, intuitive micro- and

TABLE 2 THE COMPLETION RATE FOR THE TASKS.						
Manipulation	Axis	Wrist	Elbow	Shoulder		
Throwing	-	30/30	30/30	30/30		
Translation	X-axis	30/30	30/30	30/30		
	Y-axis	30/30	30/30	30/30		
	Z-axis	30/30	30/30	30/30		
Rotation	X-axis	30/30	30/30	30/30		
	Y-axis	30/30	30/30	30/30		
	Z-axis	19/30*	30/30	17/30*		

macro-interactions when performing 3D manipulation tasks.

III. USER STUDY 2

This study aimed to find out whether micro- or macro-interactions were more frequently used by users when performing common manipulation tasks with the HTC VIVE Controller.

A. Participants and Apparatus

Fifteen participants (7 females) between the ages of 18-22 were recruited for this experiment. They all had some VR experience before. This experiment employed the same apparatus and materials as the previous experiment.

B. Task, Procedure, and Experiment Design

We used similar task design and experimental procedure as the previous experiment. However, instead of instructing the participants to finish the task by Wrist, Elbow, and Shoulder successively, we let them use the gestures which they felt most comfortable with to complete the tasks.

C. Task, Procedure, and Experiment Design

Fig. 2 shows the frequency distribution of participants performing the 3D manipulation tasks using Wrist, Elbow, or Shoulder Gestures. For both throwing tasks and rotation tasks, participants favored using Elbow Gestures the most, followed by Wrist Gestures. Few used Shoulder Gestures in these two types of tasks. Participants tended to avoid large macro-motion in these cases. On the other hand, for translating tasks, most participants favor macro-interactions (Elbow Gestures and Shoulder Gestures) than micro-interactions (Wrist Gestures). All participants used Shoulder Gestures when translating the object along Z-axis. This was probably intuitive to understand since the Wrist Gestures were not able to move the object for a long distance and may rarely be used to translate an object in normal life. We noticed that micro-gestures would always cooperate with macro-gestures to complete a translation task.



Fig. 2. The frequency distribution of participants performing Task (1)-(7) –using Wrist, Elbow, or Shoulder.

IV. DISCUSSION AND FUTURE WORK

According to our user studies, we summarize the following lessons we learned:

- The HTC VIVE Controller supports most basic micro- and macro-interactions. However, some interaction gestures seem to be unnatural for users and should be carefully considered when designing VR applications.
- VR environments should minimize requiring users to carry out large full-arm movements with the HTC VIVE Controller, and similar devices, to finish throwing and rotation tasks.
- Users tend to favor macro-interaction when translating an object using dual hand controllers like the HTC VIVE Controller.

In the future, we can either conduct elicitation studies on more complex 3D manipulation tasks which include multiple basic gestures or explore how we can transfer large macro movements (for example, Shoulder movement) to fine-level middle-range interactions (like Elbow movement).

V. CONCLUSION

This research focuses on investigating micro- and macro-interactions for manipulating objects using dual hand controllers like the HTC VIVE Controller in 3D virtual environments. We have conducted two studies that are built on each other to explore if the HTC VIVE Controller supports basic micro- and macro-interactions and whether micro- or macro-interactions are more frequently used by users to perform 3D manipulation tasks.

ACKNOWLEDGMENT

The authors wish to thank the participants for their time and the reviewers for their comments and feedback that have helped us improve our paper.

REFERENCES

- Edwin Chan, Teddy Seyed, Wolfgang Stuerzlinger, Xing-Dong Yang, and Frank Maurer. "User elicitation on single-hand microgestures." In *Proc. CHI* pp. 3403-3414. ACM, 2016.
- [2] Joseph J. LaViola Jr, Ernst Kruijff, Ryan P. McMahan, Doug Bowman, and Ivan P. Poupyrev. 3D user interfaces: theory and practice. Addison-Wesley Professional, 2017.
- [3] Katrin Wolf. "Microinteractions for supporting grasp tasks through usage of spare attentional and motor resources." In *Proc. ECCE*, pp. 221-224. ACM, 2011.
- [4] Richard Tang, Xing-Dong Yang, Scott Bateman, Joaquim Jorge, and Anthony Tang. "Physio@ Home: Exploring visual guidance and feedback techniques for physiotherapy exercises." In *Proc. CHI*, pp. 4123-4132. ACM, 2015.

First A. Author (M'76–SM'81–F'87) and the other authors may include biographies at the end of regular papers. Biographies are often not included in conference-related papers. This author became a Member (M) of IEEE in 1976, a Senior Member (SM) in 1981, and a Fellow (F) in 1987. The first paragraph may contain a place and/or date of birth (list place, then date). Next, the author's educational background is listed. The degrees should be listed with type of degree in what field, which institution, city, state, and country, and year degree was earned. The author's major field of study should be lower-cased.

The second paragraph uses the pronoun of the person (he or she) and not the author's last name. It lists military and work experience, including summer and fellowship jobs. Job titles are capitalized. The current job must have a location; previous positions may be listed without one. Information concerning previous publications may be included. Try not to list more than three books or published articles. The format for listing publishers of a book within the biography is: title of book (city, state: publisher name, year) similar to a reference. Current and previous research interests end the paragraph.

The third paragraph begins with the author's title and last name (e.g., Dr. Smith, Prof. Jones, Mr. Kajor, Ms. Hunter). List any memberships in professional societies other than the IEEE. Finally, list any awards and work for IEEE committees and publications. If a photograph is provided, the biography will be indented around it. The photograph is placed at the top left of the biography. Personal hobbies will be deleted from the biography.

Abnormal Pedestrian Behavior Detection System Using Multiple Deep Convolutional Neural Networks

Lim Kuoy Suong, Yun Borin, Lee Sunwoo and Kwon Jangwoo

Abstract—Two different deep convolutional neural networks (CNNs) models along with Motion Histogram Image (MHI) are proposed to improve the surveillance camera system in public area such as parking lot, train station, retail stores, malls, and etc. The studies are conducted using Single Shot Detector (SSD), You Only Look Once version 2 (YOLOv2) Detector (YOLOv2), and MHI technique combined with YOLOv2 to detect human face, people, and abnormal human behavior respectively in camera. The experiment evaluation showed that the proposed approach achieved over 90% of accuracy for big public datasets.

Index Terms—Abnormal behavior, CNNs, MHI, Object detection.

I. INTRODUCTION

A RTIFICIAL Intelligence is one of the active fields in the technology industry. Deep Convolutional Neural Networks (CNNs) have made many computer vision tasks possible. Additionally, there are many state-of-the-art research studies which have promising results for real-life application. In recent years, deep learning has become best known for its ability to learn from experience, and is used in complex problems. Noticeably, deep convolutional neural networks (CNNs) have made tremendous progress in large-scale object recognition [6], [10], [14], and in detection problems [13], [11], [12].

In the attempt to achieve a fully autonomous vehicle, many researchers have applied a deep CNN to extract information about the road and to understand the environment surrounding the vehicle, ranging from detecting pedestrians [1], cars [18], and bicycles, to detecting road signs [9] and obstacles [5].

Moreover, a huge number of surveillance camera systems still requires human supervision. A dramatic efficiency gains could be seen by utilizing the recent advancement in computer vision and artificial intelligence to embed in those camera systems so as to ensure public security such as criminal

Lim Kuoy Suong¹ is with Department of Computer Science and Engineering, Inha University, South Korea (email: limkuoysuong@gmail.com)

Yun Borin¹ is with Department of Computer Engineering, Inha University, South Korea (email: borinyun@gmail.com)

Lee Sunwoo is with Department of Computer Engineering, Inha University, South Korea (email: x21999@inha.ac.kr)

Prof. Kwon Jangwoo is an Inha University Professor at the Department of Computer Science and Engineering (email: jwkwon@inha.ac.kr)

¹ equal contribution authors

activities prevention and investigation, accident monitoring, people protection, public properties guarding, and etc. [15] In 2015, there is an estimate 600,000 surveillance cameras in Tianjin, China, and one camera produces around 50 petabytes of data every day. High quality of camera's resolution and large volume of video data generated by the long time span has put high pressure on data storage. For this reason, innovation in surveillance camera system is needed if we were to find or detect objects in 50 petabytes of video data in a single day.

By using CNNs, we could rely solely on the imagery input from the camera which provides a few benefits regarding designing a reliable system for surveillance camera. Firstly, costly sensors can be cut down from the device which would make it more affordable to own the device. Secondly, reducing sensors is also reducing the complexity of different data modality integration and how they share resources. As a result, this would reduce the computation power, battery life, and battery weight.

In this research study, we would like to use CNNs model with the help of MHI technique to detect people (people's face, people as a whole, and specific abnormal human behavior) presented in the public areas so as to improve the surveillance system.

II. METHODS

A. Face Detection

In the case of face detection, SSD500 (Single Shot Detector) [11] has superior performance among real-time detection networks in term of detecting relatively small object. SSD is an algorithm with a balance between object detection speed and accuracy. It executes the convolutional network on the input image only once and calculates the feature map. In order to predict the bounding box and object classification probability, a small 3×3 kernel is used to run on several feature maps. This method can detect objects of various scales. There is a trade-off between speed and accuracy in detecting objects. Hence, it is necessary to apply appropriate algorithms according to the application purpose.

Most of the other networks, on the other hand, preprocess the image to see if the objects they wanted to detect exist in the picture. The images were transformed into Image Pyramid, Sliding Window, or Regional Proposal Network (RPN) [13]. However, the fact that the objects obtained through this process have to be detected one by one in the network has caused the

processing speed of the image to be slowed down significantly. to sol Using VGG-16 network as the backbone structure, SSD is able

to solves this problem through a single-shot learning.



Fig. 2. Single Short Detection Model. (a) Image and ground truth are used as input during training phase. (b) and (c) default boxes with different aspect ratios are evaluated at different scale (e.g. 8 x 8 and 4 x 4) then predict both the shape offsets and their relative confidences (c1, c2, ..., cp). Image taken from [11]

B. People Detection

1) YOLO Object Detection

YOLO takes an approach different from other networks that use a region proposal or a sliding window; instead, it reframes object detection as a single regression problem. YOLO looks at the input image just once, and divides it into a grid of S x S cells. Each grid cell predicts B bounding boxes, a confidence score representing the intersection over union (IOU) with the ground truth bounding box, and the probability that the predicted bounding box contains objects:

$$Confidence = Pr(object) * IOU_{pred}^{truth}$$
(1)

 IOU_{pred}^{truth} denotes intersection over union between the predicted box and ground truth. Each cell also predicts C conditional class probabilities, Pr(object). Both confidence score and class prediction will output one final score telling us the probability that this bounding box contains a specific type of object.

2) YOLOv2 Architecture

In conducting our training to detect people in the images, we used YOLOv2 [12]. It contains 31 layers in which 23 are convolutional layers with a batch normalization layer before leaky rectified linear unit (ReLu) activation, and a maxpool layer at the 1st, 3rd, 7th, 11th, and 17th layers. The architecture is based on the Darknet architecture of YOLOv2 [see Table I]. In order to train our own dataset, we needed to reinitialize the final convolutional layer so it outputs a tensor with a 13 x 13 x 30 shape, where 30 = 5 (bounding boxes) x 6 (4 coordinates + 1 confidence value + 1 class probability).

IABLE I F1 Architecture. Table Adapted From [5]					
Layer	Туре	Filters	Size/Pad/Stride	Output	
0	Convolutional	32	3 x 3 / 1 / 1	416 x 416	
1	Maxpool		2 x 2 / 0 / 2	208 x 208	
2	Convolutional	64	3 x 3 / 1 / 1	208 x 208	
3	Maxpool		2 x 2 / 0 / 2	104 x 104	
4	Convolutional	128	3 x 3 / 1 / 1	104 x 104	
5	Convolutional	64	1 x 1 / 0 / 1	104 x 104	
6	Convolutional	128	3 x 3 / 1 / 1	104 x 104	
7	Maxpool		2 x 2 / 0 / 2	52 x 52	
8	Convolutional	256	3 x 3 / 1 / 1	52 x 52	
9	Convolutional	128	1 x 1 / 0 / 1	52 x 52	
10	Convolutional	256	3 x 3 / 1 / 1	52 x 52	
11	Maxpool		2 x 2 / 0 / 2	26 x 26	
12	Convolutional	512	3 x 3 / 1 / 1	26 x 26	
13	Convolutional	256	1 x 1 / 0 / 1	26 x 26	
14	Convolutional	512	3 x 3 / 1 / 1	26 x 26	
15	Convolutional	256	1 x 1 / 0 / 1	26 x 26	
16	Convolutional	512	3 x 3 / 1 / 1	26 x 26	
17	Maxpool		2 x 2 / 0 / 2	13 x 13	
18	Convolutional	1024	3 x 3 / 1 / 1	13 x 13	
19	Convolutional	512	1 x 1 / 0 / 1	13 x 13	
20	Convolutional	1024	3 x 3 / 1 / 1	13 x 13	
21	Convolutional	512	1 x 1 / 0 / 1	13 x 13	
22	Convolutional	1024	3 x 3 / 1 / 1	13 x 13	
23	Convolutional	1024	3 x 3 / 1 / 1	13 x 13	
24	Convolutional	1024	3 x 3 / 1 / 1	13 x 13	
25	Route[16]	512		26 x 26	
26	Convolutional	64	1 x 1 / 0 / 1	26 x 26	
27	Reorganize	256	2 x 2 / 0 / 2	13 x 13	
28	Route[27][24]	1280		13 x 13	
29	Convolutional	1024	3 x 3 / 1 / 1	13 x 13	
30	Convolutional	30	1 x 1 / 0 / 1	13 x 13	

C. Human Behavior Detection

Analyzing human motion is a difficult research issue due to a large change in human motion and shape, camera viewpoint, and environment setting. Motion History Image [2] is a simple yet robust way of representing movements, recognizing action, and analyze motion.



Fig. 3. Motion History Image Example. Image taken from [2]

III. EXPERIMENT

A. Dataset

1) Face Detection Dataset

To perform Face Detection, 2 different datasets has been used.

Training set:

• WIDER FACE: A Face Detection Benchmark: The dataset consists of 32,203 images and 393,703 labels. It contains rich annotations, including occlusions, poses, event categories, and face bounding boxes [16].

Testing set:

• Labeled Faces in the Wild: a database of face photographs designed for studying the problem of unconstrained face recognition [7]. The data set contains more than 13,000 images of faces collected from the web. Each face has been labeled with the name of the person pictured. However, for our study, we used only 3000 images as the testing set.

2) People Detection Dataset

We used two different datasets in conducting people detection:

Training set:

• INRIA Person Dataset: contains images from personal digital image collections taken over a long time period, and few of images are taken from the web using google images [8]. There are 3,548 images in total (1132 images are 70 x 134 pixels dimension, and 2416 images are 96 x 160 pixels dimension)

Testing set:

• Stanford 40 Actions: contains images of humans performing 40 actions. However, for the purpose of this research, we used only 3 action images (applauding, blowing bubbles, and brushing teeth) as the testing set which accounts for around 1000 images [17].

3) Human Behavior Detection Dataset

We took a video of people walking back and forth while some people stop and did some abnormal behavior of waving hands for help. After collection the data, we did some preprocessing of the video dataset. First, we converted the color video to MHI video through a process that we will describe in the following section [see section a]. Secondly, we converted the MHI video into 3698 frames. We used 3158 images for the training set and reserved 540 images for the testing set. Finally, in each frame, we labelled and annotated the bounding box representing the size and the location of the abnormal behavior. We did this by ourselves using LabelImg [20], a graphical image annotation tool freely available on GitHub.

a) Convert color video to MHI video

We propose a method to track motion recognition by extracting motion information from color image obtained from image data. The background image and the foreground image are generated using the background modeling technique to extract the Motion.



Fig. 4. System flow diagram of MHI method for motion recognition. Diagram taken from [2]



Fig. 5. MHI video segmentation by 3fps.



Fig. 6. Annotation the abnormal behavior. We drew bounding box only the area with the hands waving—from the chest of the person upward.

B. Training

Our implementation on people and human behavior detection is based on an open source YOLO framework called Darkflow [19]. We trained the two detection problems

using pre-trained weights on PASCAL VOC and/or Microsoft Common Objects in Context (MS COCO).

Training was performed on a GeForce GTX 1080 with 10 GB RAM. In all our training sessions, we used the Adam optimizer because of its tendency for towards expedient convergence.

Training for people was started with a learning rate of 1e-5 to quickly reduce loss for 100 epochs. At this point, we validated the training with test images that the model had never seen; the performance was not good, with a number of false positive and false negative bounding boxes. For this reason, we changed the learning rate to 1e- 6 for another 100 epochs so as to ensure finer granularity and proper convergence of our models.

Training for human behavior was a little bit different. The same learning rate of 1e - 5 was used for the first 100 epochs while we increased the number of epochs to 120 for 1e - 6 learning rate.

Like people detection, we trained the SSD detection at the same learning rates of 1e - 5 and 1e - 6 for 100 epochs and 100 epochs, respectively. As far as the validation process, we continued to train for another 30 epochs at the 1e - 6 learning rate. However, the model's performance got slightly worse from overfitting, so we reverted to the previous checkpoint.

IV. RESULT

We evaluated the performance of the networks using precision and recall scores with the following formulas:

Average Precision =
$$\frac{\sum_{i=1}^{n} TP/(TP+FP)}{n}$$
 (2)
Recall = $\frac{\sum_{i=1}^{n} TP/(TP+FN)}{n}$ (3)

TP, FP, and FN denote true positive, false positive, and false negative, respectively, with n representing the total number of testing images.

TABLE III Result of All Detection Model's performance						
Abnor Score Face (SSD) People Huma (YOLOv2) Behavior + YOLO						
Precision	0.9426	0.9802	0.9132			
Recall	0.9907	0.9824	0.9166			
F1-Measure	0.9661	0.9813	0.9149			

Based on the result table, we can see that SSD and YOLOv2 detector with MHI technique produced a very high precision, recall, and F1 measure score on the problem of

Face detection, People detection, and abnormal human behavior detection respectively.



Fig. 7. Example of applying Face Detection model on images of people in parking lot.



Fig. 8. Example of applying Human Detection model on images of people walking in parking lot.





Fig. 9. Example of applying Human Behavior model on images of people waving hands in office area, and their related MHI images.

V. CONCLUSION

We present SSD, MHI technique, and YOLOv2, which are ones of the current state-of-the-art object detection, to improve surveillance camera systems and to be used in real-life application. We showed that a deep-learning architecture designed can indeed be applied to detecting problem in finding face, people, and unusual behavior with promising results in accuracy. For future work, we would work more to reduce the model size and computational complexity of the network. Secondly, for this study, we set out to detect only the abnormal human behavior of waving hands. It would be a better idea if we used the system that we designed to apply on different datasets to detect other different types of unusual human behavior.

ACKNOWLEDGMENT

This research was supported by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education (2010-0020163) and by the MSIT (Ministry of Science and ICT), Korea, under the ITRC (Information Technology Research Center) support program(IITP2014-1-00729) supervised by the IITP.

REFERENCES

- [1] [Angelova et al. 2015] Angelova, A., Krizhevsky, A., Vanhoucke, V., Ogale, A. S., & Ferguson, D. (2015, September). Real-Time Pedestrian Detection with Deep Network Cascades. In *BMVC* (Vol. 2, p. 4).
- [2] Ahad, M. A. R., Tan, J. K., Kim, H., & Ishikawa, S. (2012). Motion history image: its variants and applications. *Machine Vision and Applications*, 23(2), 255-281.
- [3] Bansal, A., Nanduri, A., Castillo, C. D., Ranjan, R., & Chellappa, R. (2017, October). Umdfaces: An annotated face dataset for training deep networks. In *Biometrics (IJCB), 2017 IEEE International Joint Conference on* (pp. 464-473). IEEE.
- [4] Dalal, N., & Triggs, B. (2005, June). Histograms of oriented gradients for human detection. In *Computer Vision and Pattern Recognition*, 2005. CVPR 2005. IEEE Computer Society Conference on (Vol. 1, pp. 886-893). IEEE.
- [5] Hadsell, R., Sermanet, P., Ben, J., Erkan, A., Scoffier, M., Kavukcuoglu, K., ... & LeCun, Y. (2009). Learning long-range vision for autonomous off-road driving. *Journal of Field Robotics*, 26(2), 120-144.

- [6] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770-778).
- [7] Huang, G. B., Mattar, M., Berg, T., & Learned-Miller, E. (2008, October). Labeled faces in the wild: A database forstudying face recognition in unconstrained environments. In Workshop on faces in'Real-Life'Images: detection, alignment, and recognition.
- [8] INRIA Person Dataset. (n.d.). Retrieved from http://pascal.inrialpes.fr/data/human/
- 9] John, V., Yoneda, K., Qi, B., Liu, Z., & Mita, S. (2014, October). Traffic light recognition in varying illumination using deep learning and saliency map. In *Intelligent Transportation Systems (ITSC), 2014 IEEE 17th International Conference on*(pp. 2286-2291). IEEE.
- [10] Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In Advances in neural information processing systems (pp. 1097-1105).
- [11] Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C. Y., & Berg, A. C. (2016, October). Ssd: Single shot multibox detector. In *European conference on computer vision* (pp. 21-37). Springer, Cham.
- [12] Redmon, J., & Farhadi, A. (2017). YOLO9000: better, faster, stronger. *arXiv preprint*.
- [13] Ren, S., He, K., Girshick, R., & Sun, J. (2015). Faster r-cnn: Towards real-time object detection with region proposal networks. In Advances in neural information processing systems (pp. 91-99).
- [14] Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., ... & Rabinovich, A. (2015). Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1-9).
- [15] Xiao, J., Liao, L., Hu, J., Chen, Y. and Hu, R., 2015. Exploiting global redundancy in big surveillance video data for efficient coding. *Cluster Computing*, 18(2), pp.531–540.
- [16] Yang, S., Luo, P., Loy, C. C., & Tang, X. (2016). Wider face: A face detection benchmark. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 5525-5533).
- [17] Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., ... & Rabinovich, A. (2015). Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1-9).
- [18] Xiao, J., Liao, L., Hu, J., Chen, Y. and Hu, R., 2015. Exploiting global redundancy in big surveillance video data for efficient coding. *Cluster Computing*, 18(2), pp.531–540.
- [19] Yang, S., Luo, P., Loy, C. C., & Tang, X. (2016). Wider face: A face detection benchmark. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 5525-5533).

Person Tracking and Statistical Representation of Person Movements in Surveillance Areas

B. Platte¹, R. Thomanek¹, R. Rolletschke¹, C. Roschke¹ and M. Ritter¹ ¹Hochschule Mittweida – University of Applied Sciences, 09648 Mittweida, Germany

Abstract—Video surveillance of security-critical areas is being used more frequently. Evaluations and the timely recognition of problematic movements are limited by personnel resources. Automated systems offer an approach to the solution. Neural networks for the recognition of object classes achieve ever higher detection rates. Tracking of detected objects over time causes problems. In the first part, an algorithm is developed to track people across video frames. In the second part, possibilities of evaluation and presentation are presented using the information obtained. Furthermore, the movements of all and individual persons are displayed statistically. As a result, movements during a time window in the video can be captured at a glance without having to look at the video itself.

Index Terms—Index Terms—Video surveillance, Monitoring, Machine learning, Statistical learning.

I. INTRODUCTION

OVER the last few years, the number of surveillance cameras used has increased worldwide. Through "Closed Circuit Television" (CCTV) more and more pictures and data are recorded. Mostly these data are sighted only after an event and searched for clues. In areas of traffic safety (monitoring of intersections) or other sensitive areas, the videos should be evaluated sensibly and at the same time in a resource-saving manner. Monitoring movements for hours is a monotonous task requiring high concentration. People tire quickly. Another effect is "inattentional blindness": at high focus on one object or one activity unexpected objects are not recognized even in the middle of a field of vision ("Gorillas in Our Midst", [14]). A manual evaluation is therefore limited by mainly 2 points:

- human resources
- human error rate due to monotonous work

Earlier work trains feature vectors based on reference data to

B. Platte is with the Hochschule Mittweida, Mittweida, Germany phone: +49 3727 58 1147; e-mail: platte@hs-mittweida.de.

T. Rolletschke is with the Hochschule Mittweida, Mittweida, Germany



perform tracking. In the big data environment with thousands of hours of video material from a wide variety of starting situations, learning how to handle special situations is difficult.

The contribution of this work is the development of an algorithm for tracking detected persons without prior training using reference material. The focus here should be on universal applicability. This enables person tracking in various previously unknown scenarios. The algorithm brings detected but unconnected individual objects into a temporal dependency. This is the basis for evaluations at a higher level of abstraction. Without watching the video, typical movements and crowds of people within time slots or highly frequented points can be quickly detected at a glance. In the present test arrangement, the results have been validated manually. The positional data have been processed with the algorithm presented here for tracking persons. Individual files for all the frames (in this case between 1500 and 14000 files per video) have been summarized into one file per video with person *rediscovering across frames* data.

R. Thomanek is with the Hochschule Mittweida, Mittweida, Germany, phone: +49 3727 58 1407; e-mail: rthomane@hs-mittweida.de.

C. Roschke is with the Hochschule Mittweida, Mittweida, Germany

M. Ritter is with the Hochschule Mittweida, Mittweida, Germany.

This makes people traceable over time. Using this data, we have derived relevant information on movements within the individual scenes using aggregation methods.

II. RELATED WORK

For motion tracking of multiple objects ("multi target tracking") a wide range of approaches are available. **Particle filter** Particle filters were introduced to find targets in a state space. The particles are "greedy" for certain properties. During the iteration process, they accumulate in objects that have desired properties [5]. Dearden, Demiris, and Grau track soccer players by using a collection of SIR particle filters on color collections [6] with manual tracking start. In some cases, tracking is performed without manual start points [3].

social interaction models. Movement prediction for pedestrians about the formation of typical social behaviour patterns of pedestrians are described by [15].Using manually annotated reference data, average motion distances and directions were included as predicted values in motion tracking. The required parameters were optimized based on the reference data and generate good predictions for similar scenarios.

Tracking of properties. Neural networks detect human body structures in images [10, 7, 2].

In order to be able to clearly identify persons in the video, individual characteristics such as "color histogram" or "Local Binary Pattern" are also used for classification [13]. This makes it possible to track people across occlusions. The feature vectors must first be trained on reference material to ensure global optimization by position, size and similarity [9, 1]. Methods based on learning of patterns of reference data are very useful for similar scenes. In the area of big data, it is often necessary to work without reference data. Here it is important to effectively track down evaluable data in large quantities of unknown material.

III. DATA

Basis. There is no internationally recognised benchmark for multi-person tracking. Most of the related publications work with their own video sequences. Our investigation was based on our own video sequences. Figure 1 shows a selection of the sequences used, the respective start frame can be seen. The videos in different lengths show public space, both indoors and outdoors. People are moving. The goal was to automatically track and statistically represent the movements of the persons.

IV. PERSON TRACKING

The tracking algorithm creates a relationship between individual and unrelated item data.

A. Problem

By splitting the videos into frames quasistatic images are created, schematically represented in fig. 3. The framework Openpose considers each frame individually without any reference to the previous one. This means that any information about the keypoints of the persons is purely static. Person numbers represent only an item in the list of this single frame and are not transferable to another frame.

B. Requirements and basic idea

The coordinates of the keypoints change with their movement. Since the frames of a video run very fast (usually 25 to 30 frames per second), the pure coordinate shift from frame to frame is relatively slow.



Figure 2: Keypoints as output of the framework Openpose and tracking relevant points.

People do not usually move abruptly, but tend to glide. The concept is to buffer all coordinates of each person and compare them with the corresponding coordinates of other frames. The parameters required have been determined experimentally.



Figure 3: Schematic representation of a video, here divided into 4 frames. Several problems can be seen here: disappearance and return (person 5 in frame 3 to 4), a lot of movement without a big change of position (persons 1 and 2), occlusion (person 3 in frame 3).

C. Selecting relevant points

The further away the points are from the center of the body, the faster are their movement and consequently the change in coordinates. Figure 3 shows e.g. playing children in the foreground. Their position data of the wrists or elbows carry out large coordinate changes without greatly changing the position of the person himself. Due to its mass, the torso has the greatest inertia and moves more slowly in space than the wrists or ankles.

The key points of the human torso as well as the head fulfill the requirement of the least movement. Therefore, a set of trunk keypoints is used, consisting of shoulders, hips and the head keypoints of nose, neck and ears. In total, these 10 keypoints are used for further calculation, also shown in Figure 2:

Shoulder le	Hip le	Ear le	Eye le	Nose
Shoulder ri	Hip le	Ear ri	Eye ri	Neck

This provides a reduced data set for each person in each frame, as shown in Table I. Please note that not all points are visible in every frame.

Table I: Example data set for 1 person in 1 frame. Such a data record is created for each person in the frame.

Keypoint	x	y	Р			
Nose	173	259	0.67			
Neck						
LShoulder	167	167	0.86			
:	:	:	:			
LEar		•	•			
x, y Coordinate in video image						
P Coo	rdinate C	onfidence	-			

D. tracking algorithm

The body point detector returns a set of predictions $\mathbb{V}_f = \{v_i\}$ for each frame and as a consequence $\mathbb{V} = \{v_{f,j}\}$ for all frames. Each individual point prediction v consists of

$$v_{f,j,\kappa} = (f, x_{j,\kappa}, y_{j,\kappa}, P_{j,\kappa}) \tag{1}$$

f = frame index

(x, y) =coordinates for keypoint κ

A person is at this point only represented by his or her keypoints. Which person in the current frame matches which person from the previous frame is unknown at this point. \mathbb{Q} is the set of predictions $v_{f,j}$ for the set of all keypoints belonging to a person κ_i . Table I represents \mathbb{Q} in an example. After preselection of relevant points according to sec. IV-C the subset $\mathbb{Q}' \subset \mathbb{Q}$ results. Every single set of \mathbb{Q}' is filtered to reduce false detection:

$$\mathbb{Q}^{\prime\prime} = \{ v \in \mathbb{Q}^{\prime} | P_{\kappa} > \tau \}$$
⁽²⁾

The further calculation includes only those elements from \mathbb{Q}'' whose prediction probability is higher than the plausibility threshold τ . In the present study $\tau = 0.5$ was used. All predictions in a frame are given by the set of all \mathbb{Q}_j'' with *m* as number of persons in the frame:

$$\mathbb{R} = \bigcup_{j=1..m} \mathbb{Q}_j^{\,\prime\prime} \quad (3)$$

All predictions in frame 1 are represented by \mathbb{R}_1 , and all predictions in frame 2 by \mathbb{R}_2 , and so on. The x-coordinate of a keypoint κ_i for the current person in the current frame is

$$x_{f,i} = x(\{v \in \mathbb{Q}_{j,f}^{\prime\prime} | \kappa_i\})$$
(4)

and $x_{f-1,i}$ is the corresponding x-coordinate of a keypoint of another person in the previous frame. For a person represented by the prediction set \mathbb{Q}'' , the most suitable person in the previous frame is searched. The Euclidean distances in the 2D space of the video and the confidence of each keypoint κ_i to the related κ_i of all predictions of a person in the previous frame (f-1) are calculated:

Table II: distance array for 1 person in current frame to all persons of previous frame (example)

d distance

- *m* person count in previous frame
- \mathbb{D} all keypoint distances of one person
- \mathbb{E} set of all \mathbb{D} of one person in current frame

Ø Value does not exist ("NaN"-Value)

crossed out by filtering

Keypoint	distance $d(\mathbb{D}_1)$	d to all perpendicular $P(\mathbb{D}_1)$	ersons $d(\mathbb{D}_2)$	$P(\mathbb{D}_2)$		$P(\mathbb{D}_m)$
Nose Neck RShoulder	30 22 24 27	0.65 0.76 0.45	4 5 2 2	0.71 0.63 0.34		Ø Ø Ø
Eshoulder E LEar	27 : 0.23	0.71 : 0.23	2 : 3	0.86 : 0.91	 ` 	© : Ø
mean value $\overline{\mathbb{E}}$ Minimum $M =$	25 $\min d(\overline{\mathbb{E}})$	0.71	4	0.64 0.68		
$d < \epsilon$?			\checkmark			

$$d_i = \sqrt{(x_{f,i} - x_{f-1,i})^2 + (y_{f,i} - y_{f-1,i})^2}$$
(5)

$$P_i = \frac{(P_{f,i} - P_{f-1,i})}{2} \tag{6}$$

All person's distances with their confidences to another person form the distance set $\mathbb{D} = \{\{d_1, b_1\} \dots \{d_m, b_m\}\}$. All Sets \mathbb{D} for all prediction sets of the previous frame \mathbb{R}_{f-1} are calculated and define the set

$$\mathbb{E} = \bigcup_{j=1..m} \mathbb{D}_j \tag{7}$$

where: m = set of persons within frame

The resulting set \mathbb{E} now contains all point distances *of one person* to all persons in the previous frame, exemplified in table II. In the further course, within the quantity \mathbb{E} the minimum is determined from the set of the mean values of each subset $\mathbb{D}_1 \dots \mathbb{D}_m$:

$$\overline{\mathbb{E}} = \sum_{i=1}^{m} \{ d(\mathbb{E}_i), P(\mathbb{E}_j) \}$$
(8)
$$M = \min d(\overline{\mathbb{E}})$$
(9)

In the example in table II the minimum value will be $\frac{4}{4}$

Pixel and occurs at person 2 in column $d(\mathbb{D}_2)$. These minimum of all eligible persons is compared to a fixed threshold \in *(enquoteexclusion circle)*. Only if that minimum is below this threshold, the person is considered to be found again (\checkmark in table II). This results in the \in parameter forming a robust *feedback exclusion circle* around the coordinates. If an object disappears and another object appears in another position in the following frame, the exclusion circle prevents the new object from being assigned to the object that has just disappeared.



Figure 4: searching for suitable sets \mathbb{E} in previous frames h



The described procedure is iteratively performed *for each person* in each frame. Furthermore, the distance sets \mathbb{E} are calculated for frames reaching further into the past. The "history depth" is specified by a parameter *h* ("history depth") as shown in Figure 4.

V. RESULTS/DISCUSSION

In the videos, the persons in the individual frames were tagged with a frame-spanning identifier, a unique ID. To achieve a clearer display, the frames were aggregated over whole seconds.

A. Persons during the run time

Figure 5 shows the number of persons detected and tracked every full second of the running video. In video 1 around 10 people are visible at the beginning, after 17 seconds the number of people decreases. After 30 seconds the number of persons rises again. In video 2, people were detected only during the first 23 seconds. At the beginning 2 persons were detected, later a bit more and again decreasing. Apparently, a group was walking through the picture.

B. Retrieval after occlusion

Figure 8 schematically displays the retrieval after occlusion. Figure 6 shows the total number of detection gaps of the specific person and thus the "re-locking" of the tracking algorithm, controlled by the parameter history depth (fig. 4). The "Noise" in the range from 0 to 5 is caused by detection dropouts in the detector data used. During the manual check it



Figure 6: Number of frme gaps for each person, counted in frames



Figure 7: Number of missing frames for each person in 2 scenes

was visible that especially in the background, in blurred areas or in side views the detection of persons for individual frames is interrupted. These dropouts become visible here. The higher values were caused by occlusion, i.e. one person was repeatedly occluded by another person during the duration of the video. The detector temporarily couldn't find a hidden person, the tracking algorithm found them again after the detection. Prerequisites for retrieval ("person relock") are

- The occlusion duration remains within the historical search range h: the person is detected again before the history depth parameter limit is reached (fig. 4, fig. 8).
- The covert person does not move out of the exclusion circle during the occlusion, defined by the parameter "exclusion circle" ∈.

Figure 7 shows another view onto the locking of the tracking algorithm. The accumulated number of missing frames is shown here. In scene 1 not all people have missing frames. For persons with missing frames, the average value is less than 25 frames, i.e. less than 1 second. Less than half of the persons can be seen continuously without interruption during their visibility period. There are longer gaps in some of them. This can be explained by the passing of other persons or temporary occlusion.

C. Movement aggregation by heatmaps

Figure 9 shows summarized movement data of the heads. For this evaluation, the video image was divided into 10x10 pixel sized fields. As a result, a matrix is created that reproduces the movements in the image plane in reduced resolution, both spatially and temporally.



Figure 8: Tracking based on previous frames. Blurred persons are unknown. Sharply depicted persons were found in previous frames.

In the summary, the keypoints "nose" and "neck" were considered further. The elements of these one-second heatmaps were averaged over the runtime of the video and visualized in fig. 9. The result represents the aggregated illustration of typical movements during the video runtime. The heatmaps in fig. 9 visualizes locations of movement. The persons walk e.g. in scene 1 under avoidance of the lantern from the foreground to the house or take the way back. They walk around the flowerbed on the right. In scene 2, the heatmap indicates that the walkway is typically used in an arc instead of a straight line. With this one image per video the typical movement locations of people can be identified without seeing the video itself.

D. Motion aggregation through tracklet generation

The vector representation in fig. 9 allows conclusions on the direction of movement. For this purpose, the distances of each person are added separately over selectable time slots. Figure 9 show these movements in row 2 as overlay of the first frame. Each person will be painted in a new colour. Scene 1 shows that much more people move through the scene during the video's runtime than can be seen in the first frame. The difference to

the Heatmap becomes clear: a typical direction becomes recognizable. In this way it can be clearly distinguished that many people move towards the house, and some move away from the house. Scene 3 shows a person waiting at the elevator (blue vectors). Further on, according to the green vectors, a person comes out of the elevator and walks into the foreground. Another person, represented by red vectors, comes through the door on the left, turns and walks to the elevator. Typical movement patterns become visible at a glance.

VI. ASSUMPTIONS, LIMITATIONS, FUTURE WORK

The algorithm doesn't detect objects in videos itself. These are based on the output of detectors such as Detectron [8], YOLOv3 [12, 11] or Openpose [4]. Improvements in these detectors will improve tracking. Tracking is performed on the 2D projection plane of the video. However, surveillance videos are generally filmed from a raised position and show a spatial perspective. Objects further away from the camera seem to move more slowly at the same speed than objects closer to the camera. The algorithm is currently not taking this into account. When moving an object, the algorithm searches for the object in the constant exclusion circle surrounding the keyppoints of the last position. At a constant speed of an object, the perspective-related circle of retrievability in the foreground increases. The parameter \in additionally depends on the resolution of the underlying video data. For higher-resolution videos, motion includes larger pixel distances. The parameter must also be adjusted accordingly.

VII. FUTURE WORK

The tracking of objects, in this case persons, is based on the search for the minimum distances of specific points. A prediction algorithm might improve tracking accuracy: A track fragment could be extrapolated from the previous frame and used for calculating a probable target. Extending the feature vector by additional characteristic properties is expected to yield further improvements.

VIII. CONCLUSION

In this paper, an algorithm for efficient tracking and temporal attributing of persons in videos was presented. The attributing of persons has been based on a feature vector of 10 points of the human torso plus neck and head. The feature vectors are calculated for each person with parameterizable historical depth. This deep temporal grouping of all individual elements of the feature vectors results in a robust recognition of the respective person. Detection gaps or short-term occlusions are successfully noticed and rectified. The algorithm works without specific learning methods. This makes it applicable to unknown video material.

REFERENCES

 Mykhaylo Andriluka et al. "PoseTrack: A Benchmark for Human Pose Estimation and Tracking". In: *arXiv:1710.10000 [cs]* (Oct. 2017). arXiv: 1710.10000 [cs].



Figure 9: Aggregated visualization of all movements during the runtime of the respective video. Line 1 shows the grid for the heatmap following in line 2, line 3 the vectorized representation.

- [2] Hou Beiping and Zhu Wen. "Fast Human Detection Using Motion Detection and Histogram of Oriented Gradients". en. In: *Journal of Computers* 6.8 (Aug. 2011).
- [3] M. D. Breitenstein et al. "Online Multiperson Trackingby-Detection from a Single, Uncalibrated Camera". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 33.9 (Sept. 2011), pp. 1820–1833.
- [4] Zhe Cao. OpenPose: Real-Time Multi-Person Keypoint Detection Library for Body, Face, and Hands Estimation. May 2018.
- [5] J. W. Choi and J. H. Yoo. "Real-Time Multi-Person Tracking in Fixed Surveillance Camera Environment". In: 2013 IEEE International Conference on Consumer Electronics (ICCE). Jan. 2013, pp. 125–126.
- [6] A. Dearden, Y. Demiris, and O. Grau. "Tracking Football Player Movement from a Single Moving Camera Using Particle Filters". In: *European Conference on Visual Media Production (CVMP)*. IET, Nov. 2006.
- [7] Á Garcia-Martin and J. M. Martinez. "People Detection in Surveillance: Classification and Evaluation". In: *IET Computer Vision* 9.5 (2015), pp. 779–788.
- [8] Ross Girshick et al. Detectron. 2018.
- [9] Y. Li, C. Huang, and R. Nevatia. "Learning to Associate: HybridBoosted Multi-Target Tracker for Crowded Scene". In: 2009 IEEE Conference on Computer Vision and Pattern Recognition. June 2009, pp. 2953–2960.
- [10] Alejandro Newell, Kaiyu Yang, and Jia Deng. "Stacked Hourglass Networks for Human Pose Estimation". In: arXiv:1603.06937 [cs] (Mar. 2016). arXiv: 1603.06937 [cs].
- [11] Joseph Redmon and Ali Farhadi. "YOLOv3: An Incremental Improvement". In: *arXiv* (2018).
- [12] Joseph Redmon et al. "You Only Look Once: Unified, Real-Time Object Detection". In: arXiv:1506.02640 [cs] (June 2015). arXiv: 1506.02640 [cs].
- [13] G. Shu et al. "Part-Based Multiple-Person Tracking with Partial Occlusion Handling". In: 2012 IEEE Conference on Computer Vision and Pattern Recognition. June 2012, pp. 1815–1821.
- [14] Daniel J Simons and Christopher F Chabris. "Gorillas in Our Midst: Sustained Inattentional Blindness for Dynamic Events". en. In: *Perception* 28.9 (Sept. 1999), pp. 1059–1074.

[15] Kota Yamaguchi et al. "Who are you with and where are you going?" English (US). In: 2011 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2011. 2011, pp. 1345–1352.

Wireless Sensor Network Traffic Modeling and Anomaly Simulation based on OPNET

Weihao Liu, Haoyang Wang, Limin Yu, Mark Leach and Fei Ma

Abstract—The interaction between nodes in a multihop wireless sensor network is a complex issue. Markov random field (MRF) theory has been used as an important mathematical tool to model local statistical interactions in a sensor network. This paper focuses on the OPNET simulation of a typical broadcasting network of distributed control under shortest path routing. Due to imposed power constrains, each sensor node only communicates with its neighboring nodes forming an MRF. The network traffic intensity in both normal and anomaly situations is simulated in OPNET. The simulation results demonstrate the validity of the MRF based inference structure and point out a way to improve the model parameter optimization and anomaly detection.

Index Terms—Wireless Sensor Network, Traffic Modeling, OPNET, Markov Random Field, Network Anomaly Analysis.

I. INTRODUCTION

THE wireless sensor networks (WSNs) have been applied widely in different areas including remote monitoring of the environment, target tracking, satellite communication and distributed control [1]. The sensor nodes of WSNs could be made smaller, cheaper while even more intelligent compared with the traditional wired sensor networks [2]. The information collected by the sensor nodes would, in most cases, be sent via the intermediate nodes to the sink node for further processing. The interaction between nodes is a complex issue and most studies of multihop networks at different layers of the protocol stack face the issue of understanding this complex process [3, 4].

As an important mathematical tool to study statistical dependencies of a set of random variables, Markov random

This work was supported by the National Science and Foundation of China (NSFC) under Grant 61501380 and in part by the Research Institute of Big Data Analytics, Xi'an Jiaotong-Liverpool University.

Weihao Liu and Haoyang Wang are with the Department of Electrical and Electronic Engineering, Xi'an Jiaotong-Liverpool University, Suzhou 215123, China.

Limin Yu is with the Department of Electrical and Electronic Engineering, Xi'an Jiaotong-Liverpool University, Suzhou 215123, China. (corresponding author, phone: 86-512-88161414; e-mail: limin.yu@xjtlu.edu.cn).

Mark Leach is also with the Department of Electrical and Electronic Engineering, Xi'an Jiaotong-Liverpool University, Suzhou 215123, China. (e-mail: Mark.Leach@xjtlu.edu.cn).

Fei Ma is with the Department of Mathematical Sciences, and Research Institute of Big Data Analytics, Xi'an Jiaotong-Liverpool University, Suzhou 215123, China. (e-mail: fei.ma@xjtlu.edu.cn). field (MRF) theory has been utilized to model local statistical interactions to solve large scale inference problems in sensor networks [5, 6]. In [6], MRF theory was applied to model traffic intensity of three types of sensor networks with lattice structure under the shortest path routing. The three network structures included uniform communication, central gathering and border gathering. The dynamic variation of the traffic load in practical situation was simulated by adding the Gaussian noise to the traffic distribution. Network traffic prediction was implemented by a MRF smoothing algorithm.

This paper is an extension of the work in [6]. It aims to use the network simulation tool, the OPNET, to simulate the network traffic of one type of the sensor networks described in [6], where each sensor node possesses the information from all the other nodes under broadcasting and shortest path routing as in a distributed control network. The simulation results demonstrate the validity of the MRF based traffic inference algorithm. It also proves that the original model should be adjusted to fit more realistic network situations. The location of anomaly nodes and its influence on the traffic distribution are simulated and analyzed.

II. METHODOLOGY

A. The structure of the sensor nodes network

In this paper, the uniform communication [6] structure is adopted to set up the OPNET simulation. As shown in Fig. 1, all sensor nodes are the sink node and each generated source data packet will be sent to all the other nodes as the destination. With power constraints, each sensor node can only communicate with its immediate neighbors for data transfer or packet relay, e.g., there will be 4 nearest neighbors while the sensor node is not on the edge of the network.

In OPNET simulation, we specify a square area of 300 by 300 meters. As demonstrated in Fig. 2, 81 sensor nodes are uniformly deployed in this given area in a 9 by 9 lattice.

The distances between two neighboring nodes is set to approximately 25 meters. The distance between two nearest sensor nodes at the diagonal line is about 35 meters. Therefore, it is important to set in the OPNET that the transmission distance for single node is limited to 30 meters to avoid the communication between two sensor nodes that are not adjacent in the network.



Fig. 1. Uniform communication model.



Fig. 2. WSN model in the OPNET with 81 nodes.

The generated data packet has the size of 1024 bits. The data flow lasts for 5 minutes in the simulation. The traffic data at each node and the total traffic are measured and updated in the OPNET every 3 seconds. The data could be visualized conveniently in the system. Fig. 3 illustrates the data flow of a single node (node 40) to the 4 neighboring nodes (node 31, 39, 41 and 49) which makes one clique in a MRF structure. To be noted that, the 4-point cliques will be adjusted to the 2-point clique and 3-point clique for nodes at the edge of the network.



Fig. 3. Data transmission in one clique where data packages are sent from a sink node to four adjacent sink nodes with the equal probabilities

B. Setting of a single network node

A typical structure of a sink node is shown in Fig. 4. The node has 6 modes including source, sink, network layer, 802_15_4 mac mode and wireless transmitter/receiver. It enables specific network settings in this layered structure. The blue arrows represent up and downlink of the data flow across the layers.



Fig. 4. Inner structure of a sink node (6 modes and 2 types of arrow in a single sink node)

The source mode produces the new data package based on the Poisson distribution. The time intervals for arriving new data packages conform to the exponential distribution. The network layer is set as the broadcasting mode (from each sink node to all the other nodes). 802_15_4_Mac mode regulates the link layer implementing ZigBee MAC protocol. Current simulation settings fit in the general ZigBee network scenarios with transmission distance between 10 to 100 meters for applications in embedded sensing, industrial control or building monitoring and automation.

III. RESULTS AND ANALYSIS

A. Data traffic intensity at a single node and the traffic distribution of the network

In the OPNET simulation, we pick 5 nodes with the ID number 0, 10, 20, 30 and 40, located at the diagonal line from network edge to the center for detailed analysis as shown in Fig. 5.

The traffic intensity of the selected nodes resulted from a 5 minutes simulation is shown in Fig. 6. The traffic data is fluctuating which reflects the real network traffic statistics with incidents of delay, congestion and temporary failure. The average traffic intensity is shown in Table I. A plot of the traffic



Node Numbers	0	10	20	30	40
Average Traffic	35023	35850	36895	37620	37993
Intensity	(bit/s)	(bit/s)	(bit/s)	(bit/s)	(bit/s)

Table I. Table of average traffic amount of 5 nodes



Fig. 7. Simulated Data traffic distribution in uniform communication network

The simulation illustrates that the amount of traffic increases from edge to the center of the network which is consistent with the theoretical results in [6]. However, in comparison with the ideal traffic intensity profile as shown in Fig. 8, more realistic traffic intensity simulated by OPNET reveals that the traffic variation in percentage at different locations of the network is much smaller than the theoretical model. The node 40 at the central position, which has the maximum traffic is only 7.8% larger than the edge (node 0) as shown in Fig. 7.

This result suggests that the power consumption for sensor nodes at different locations varies but the power management strategy may not differ significantly. Applying different buffer



Fig. 8. Ideal model of the traffic distribution in uniform communication network

size according to the distance of sensor nodes to the center (larger in the center and smaller on the edge) provides a simple solution to improve the general performance in this type of network.

B. Anomaly detection and the anomaly influence

To better understand the influence of network anomaly on the network traffic, OPNET simulation is conducted by adding anomaly nodes at different areas of the network. 9 nodes of a 3 by 3 lattice in 4 different locations are selected as the anomaly nodes for testing as shown in Fig. 9.

The selected anomaly nodes are not able to generate and relay data packets. The anomaly is set to occur at 100 seconds


Fig. 9. Anomaly nodes in 4 locations: (a) Anomaly at the top left corner, (b) at inner positions (c) at inner positions (closer to the center) (d) at the central positions.

after the start of simulation and the traffic is dropped to zero at these nodes. The failure of the selected nodes affects the surrounding nodes, and the total traffic intensity in the 4 scenarios are illustrated in Fig. 10.

Table II provides a numerical comparison of the total traffic in normal and 4 anomaly situations. The own traffic generated by the 9 anomaly nodes (in original normal situation) is about 10%. It is observed that the relay traffic is affected significantly by the location of the anomaly nodes. More relay traffic is affected when the anomaly nodes are approaching the center of the network. The total traffic reductions increases from 16.1% to 27.4% as shown in Table II.

To be noted that, current anomaly simulation assumed the complete failure of selected nodes. Protocol design may provide an alternative for anomaly detection. However, if abnormal nodes experience a connection problem in certain time period (i.e., abnormal nodes are not completely failed), the MRF inference based on the change of traffic profile would be a more effective way for anomaly detection.

IV. CONCLUSION

The paper uses OPNET to simulate a typical wireless network scenario where, with power constraints, each sensor node is only able to communicate with the neighboring nodes forming a MRF of 4-point cliques. Under the normal networking settings with broadcasting and shortest path routing, the simulation demonstrated the validity of the traffic distribution in theory and the MRF model for traffic prediction. The simulation and the



Fig. 10. Plot of traffic intensity of 4 anomaly scenarios

theoretical results differ quantitatively which infers that the cost function parameters could be adjusted in the MRF model for better prediction

The simulation also suggested that applying different buffer size according to the distance of sensor nodes to the center provides a simple solution to improve the network performance. The anomaly simulation with OPNET clearly indicated that the failure of central node leads to more significant traffic loss. The numerical connection between the amount of traffic loss and locations of the anomaly nodes provides a way to solve the inference problem of network anomaly detection which is worth further investigation.

REFERENCES

- J. Sun and E. Modiano, "Routing strategies for maximizing throughput in leo satellite network," *IEEE journal on selected areas in communications*, 22(2), Feb. 2004.
- [2] A. Bharathidasan, V. A. S. Ponduru, "Sensor Networks: An Overview," Department of Computer Science University of California. Available: http://www.csun.edu/~andrzej/COMP529-S05/papers/sensorNetworksS urvey.pdf
- [3] B. Zeng, Y. Dong, J. He, D. Lu, "An energy-efficient TDMA scheduling for data collection in wireless sensor networks," in *Proc Communications* in *China (ICCC)*, 2013.
- [4] C. Intanagonwiwat, R. Govindan, D. Estrin, J. Heidemann, and F. Silva, "Directed diffusion for wireless sensor networking," *IEEE/ACM Transactions on Networking*, 11(1):2–16, February.
- [5] D. Ganesan, A. Cerpa, Y. Yu, W. Ye, J. Zhao, and D. Estrin, "Networking issues in sensor networks in the journal of parallel and distributed computing (jpdc)," *Invited paper in Special issue on Frontiers in Distributed Sensor Networks*, Elsevier Publishers, 2003.
- [6] Y. Cai and L. Yu, "Sensor network traffic load prediction with Markov random field theory," in *Proc 4th International Conference on Computer Science and Network Technology (ICCSNT)*, Harbin, 2015, pp. 967-971.

Anomaly Situation	Normal	1	2	3	4
Network Total Traffic	5597173(bit/s)	4698631 (bit/s)	4501410 (bit/s)	4276461 (bit/s)	4066656 (bit/s)
Data Loss	0%	16.1%	19.6%	23.5%	27.4%

Table II. Comparison of network traffic in normal and 4 anomaly scenarios

A Low Multiplicative S-Box for a Stochastic Random Number Generator

M. L. Dennis Wong and Jia Jun Tay

Abstract— Stochastic computing is an emerging paradigm for fault tolerant digital computation. Unlike conventional computing, stochastic computing has better tolerance towards high noise floor. This is an important figure of merit as we move towards lower power platforms, area optimised design is desired to achieve low power consumption. In this paper we present an area optimised design for a S-box based Pseudo Random Number Generator for stochastic computing.

Index Terms—Low multiplicative complexity, random number generator, stochastic computing, S-box.

I. INTRODUCTION

WITH the advent of Internet of Things, the demand for low power IC design continues to grow than ever before. This leads to the exploration of alternative computing paradigms to enable IoT devices to operate at extremely low power. One of the alternative paradigms being increasingly looked at is Stochastic Computing (SC) introduced by [1] in 1967. This is primarily owing to SC's low power requirement and its fault tolerance ability.

This technique, lent itself from the probability theory, has been proven to be better in handling computational uncertainties [2]. Furthermore, SC is particularly attractive in IC design as it requires low complexity computation blocks. In general, SC is seen as a promising alternative in comparison to its conventional binary computing counterpart, which usually has a higher computational cost.

Though SC has been known for decades, very few physical realizations have been proposed. Initially, SC applications were limited to the field of neural networks [3] and machine controls [4]. Until recent years, it was discovered that SC efficiently mathematical simplifies some functions which are computational expensive in binary computation. These functions can be efficiently approximated using stochastic logic with minimal hardware requirements and without significant accuracy degradation. Ever since, SC implementation has been extended to image processing [5]-[7], error control coding applications [8] and digital filter design [9]-[12].

However, the drawbacks of SC include computation

inaccuracies, additional computation times compared to conventional computing, etc. Besides, to produce stochastic numbers (SN) from the operands, one would require stochastic number generators (SNG). SNs are strings of binary bits in which the probability of "1"s is determined by the value of the operand it represents. Typically, each copy of input SN applied to a stochastic circuit requires its own independent SNG. Furthermore, additional SNGs are also required to reduce the correlation between the SNs during computation. In [13], it was stated that the hardware needed for the SNGs can take up more than 80% of a stochastic circuit's area.

This work studied a recently proposed SNG and proposed an optimised implementation for its S-Box. Using a low multiplicative approach, it is found that one can further reduce its circuit area and in turn its dynamic power.

II. S-BOX BASED STOCHASTIC NUMBER GENERATOR

Recently, a new S-box based pseudo-random number generator (SBoNG) for stochastic circuits is proposed by [14]. The new design combined an LFSR with a non-linear S-box function. As it does not interfere with decorrelation and thus it can be shared efficiently by multiple stochastic number generator. It was stated by the authors that SBoNGs can scale effectively in stochastic circuits operating on a large number of input variables.

The S-box used by the authors was taken from a Mini-Advanced Encryption Standard (AES) Algorithm proposed in [15]:

which could be represented as the followings, denoting input, $X = \{x_0, \dots, x_3\}$ and output, $Y = \{y_0, \dots, y_3\}$:

TABLE I S-Box Table								
Х	0	1	2	3	4	5	6	7
Y	6	В	5	4	2	Е	7	А
Х	8	9	А	В	С	D	Е	F
Y	9	D	F	С	3	1	0	8

 $y_3 = \bar{x}_2 x_3 + x_0 \bar{x}_1 \bar{x}_3 + x_0 x_1 x_2$

 $y_2 = x_0\overline{x}_1x_2\overline{x}_3 + x_1\overline{x}_2 + \overline{x}_0\overline{x}_2\overline{x}_3 + x_0\overline{x}_2x_3 + \overline{x}_0x_1\overline{x}_3$

 $y_1 = \bar{x}_0 x_1 \bar{x}_2 x_3 + \bar{x}_1 \bar{x}_3 + x_2 \bar{x}_3 + x_0 \bar{x}_1 x_2$

 $y_0 = x_0 \bar{x}_1 \bar{x}_2 + \bar{x}_0 x_1 \bar{x}_3 + \bar{x}_1 x_3 + \bar{x}_0 \bar{x}_2 x_3$

The proposed S-box was implemented diredtly in [14] as depicted in Figure 1. The S-boxes were then combined with

Manuscript received July 26, 2018.

M. L. D. Wong is with Heriot-Watt University Malaysia, Precinct 5, Putrajaya, 62200, Malaysia. (e-mail: D.Wong@hw.ac.uk).

J. J. Tay is with Swinburne University of Technology Sarawak Campus, Kuching, Sarawak, 93350, Malaysia. (e-mail: jtay@swinburne.edu.my).



Fig. 1. A direct implementation of the S-box using multiple-input gates. implementation as shown in Figure 2.



Fig. 2. SBoNG implementation as proposed [14].

III. LOW MULTIPLICATIVE DESIGN AND BOYAR-PERALTA ALGORITHM

By definition, multiplicative complexity $c_{\wedge}(f)$ is a measure of the minimum number of multiplications (AND gates) necessary to compute a function *f* over the logic basis (AND, XOR, NOT). It is similar to the notions of circuit complexity introduced in [16]. Generally, determining the exact multiplicative complexity of a given function is a difficult problem. Among the important studies on this topic [17], [18], the propositions from [19] and [20] allow us to reduce the range down to $d-1 \le c_{\wedge}(f) \le n-1$ for functions with $n \le 5$ input variables.

Remark (Degree). The degree *d* of a function *f* is the highest order of the individual terms in its XOR sum-of-products (XSOP) expression. For example, a function $f = x_1x_2x_3 \oplus x_1x_4$ has a degree of d = 3.

We discuss the notion of multiplicative complexity as it has significant implications on deriving a low gate count circuit for a desired function. Boyar *et al.* [21] demonstrated how the concept of optimal multiplicative complexity can be leveraged to construct low gate count circuit for the AES $GF(2^4)$ inversion circuit in [22]. The approach is based on the heuristic which states that a circuit constructed with the minimum number of AND gates results in close to optimal gate count. The actual optimisation process is handled by a two-step algorithm [23] which includes: (a) AND-minimisation step and (b) XOR-minimisation step.

The AND-minimisation step is the foundation of the Boyar-Peralta algorithm and is responsible for minimising the number of AND gates in non-linear (contains both XOR and AND gates) portions of the target circuit. Fundamentally, the AND-minimisation step from [23] is a *randomised* selection algorithm. The algorithm initialises a sample space with all input variables of the target function. This sample space is then expanded through additions and multiplications (XOR and AND operations) on randomly selected pairs of elements from the sample space. The expansion process is repeated until a sum or product that computes the target function is discovered. To ensure optimal multiplicative complexity, any randomly generated sum or product that has an AND gate requirement exceeding the multiplicative complexity is discarded.

Once an optimal multiplicative complexity implementation is discovered, the XOR-minimisation step is applied to reduce gate count in linear (contains XOR gates only) portions of the circuit. First, the algorithm identifies the target signals to be computed and calculates the distance vector (number of XOR required for each signal). In each iteration, the algorithm identifies a pair of available signals whose sum allows a maximum reduction in the magnitude of the distance vector. This process is repeated until the distance vector is reduced to zero. This approach is equivalent to solving a shortest linear path (SLP) problem as implied in [21].

IV. ALGEBRAIC DECOMPOSITION ALGORITHM

The Boyar-Peralta algorithm has shown great results in optimising digital circuits. The approach is fairly flexible as shown in [24] where it is tweaked to reject sums/products that exhibit high circuit depth for a more speed-focused optimisation. However, it also comes with consistency issues due to the randomised selection process occurred in the original AND-minimisation step. Most complex functions have multiple implementations that are optimal in terms of multiplicative complexity. To identify a *good* solution from many alternatives, a large number of solutions need to be gathered using the algorithm. Even then, the potential for a better solution cannot be denied.

Enhancements to the randomised search algorithm can be done to improve the odds of finding *good* solutions and for a more consistent computational time for the algorithm [25]. However, we propose a deterministic approach to the low multiplicative complexity heuristic based on algebraic decomposition. This new approach keeps the two-step procedure of the heuristic but replaces the randomised selection algorithm in the AND-minimisation step with a deterministic tree search algorithm [26]. Sometimes, the two approaches present the same results as in the case of this paper.

Given an algebraic expression, decomposition or factorisation is a method to reduce the amount of multiplications in said expression. The same concept can be applied to functions described over the basis (AND, XOR, NOT) in which the AND and XOR operations are *GF*(2) multiplication and addition respectively. Consider a function $f = x_1x_2 \oplus x_1x_3$ which requires two AND gates and one XOR gate to implement at first glance. After decomposition, the function can be implemented with only one AND gate and one XOR gate as shown below:

$f=x_1(x_2\oplus x_3)$

Following the example above, we designed an algorithm to factorise each literal in the expression in order to identify all optimal solutions achievable through decomposition. However, decomposition on single function is not sufficient in solving most practical problems that produce multiple outputs. A multiple-output function can be described with an algebraic expression for each individual output. Optimising each expression individually often does not yield optimal result for the function as a whole. To solve for low multiplicative complexity implementation for multiple-out functions, *product sharing* between expressions are essential. We first establish the four functions of the S-box in Positive Polarity Reed-Muller (PPRM) form as follows (they compute the same function described in Table I):

$$y_0 = x_0 x_1 x_2 \oplus x_0 x_1 x_3 \oplus x_0 x_2 x_3 \oplus x_1 x_2 x_3 \oplus x_0 x_2 \oplus x_0 x_3 \oplus x_1 x_3$$
$$\oplus x_0 \oplus x_1 \oplus x_3$$
$$y_1 = x_0 x_1 x_3 \oplus x_0 x_2 x_3 \oplus x_1 x_2 x_3 \oplus x_1 x_2 \oplus x_2 x_3 \oplus x_1 \oplus x_3 \oplus 1$$

 $y_2 = x_0 x_1 x_2 \oplus x_0 x_2 x_3 \oplus x_0 x_1 \oplus x_1 x_2 \oplus x_1 x_3 \oplus x_2 x_3 \oplus x_0 \oplus x_2 \oplus x_3 \oplus 1$

 $y_3 = x_0x_1x_2 \oplus x_0x_1x_3 \oplus x_0x_1 \oplus x_0x_3 \oplus x_2x_3 \oplus x_0 \oplus x_3$

As an illustrative example, we can construct y_0 and y_2 in optimal multiplicative complexity as follow, each requiring two AND gates:

 $y_0 = x_0 \oplus x_1 \oplus (x_0 \oplus x_3) \times (x_3 \oplus (x_0 \oplus x_1) \times (x_1 \oplus x_2))$

 $y_2 = (x_0 \oplus x_1 \oplus x_3) \times (x_0 \times x_2 \oplus x_0 \oplus x_2 \oplus x_3) \oplus x_2 \oplus 1$

However, if the products in y_0 is shared with y_2 , we can construct y_2 with only one additional AND gate as shown below: $y_2 = (x_0 \oplus x_2 \oplus x_3 \oplus p_1) \times (x_0 \oplus x_1 \oplus p_1 \oplus p_2) \oplus x_2 \oplus 1$

Where p_1 and p_2 designate the products shared with y_0 : $p_1 = (x_0 \oplus x_1) \times (x_1 \oplus x_2)$

 $p_2 = (x_0 \oplus x_3) \times (x_3 \oplus (x_0 \oplus x_1)(x_1 \oplus x_2))$

As such, the tree search algorithm is designed to form a collective set of product(s) from solved expressions. When solving for subsequent expressions that are part of the same multiple-output function, these *free* products can be XOR-ed into the expression during different stages of decomposition to enable solutions with less AND gates in total.

Not all solutions that are optimal in multiplicative complexity exhibit the same circuit size. As such, a final selection step is required in the algorithm to return the solution with the lowest gate count. The overall flow of the proposed algorithm can be summarized in Figure 3.

V. RESULTS

In this Section, we present the results obtained using the



Fig. 3. Flow of the proposed deterministic algorithm for low multiplicative complexity circuit.

above two-step optimisation algorithm. The optimised non-linear S-box for stochastic RNG is as presented in Figure 4.

The original S-box is illustrated as a circuit of 15 AND gates and 4 OR gates. However, a number of the gates required are multiple-input AND or OR gates. To enable a fairer comparison, we converted the circuit into one with only two-input gates. This resulted in a circuit with 22 AND and 12 OR gates. Table II gives a comparison of the number of

OK gates.	Table II gives a comp	barison of the number of
$t_1=x_0\oplus x_1$	$t_2 = x_1 \oplus x_2$	$t_3 = t_1 \times t_2$
$t_4=x_3\oplus t_3$	$t_5 = t_4 \oplus t_2$	$t_6 = x_0 \oplus x_3$
$t_7 = t_6 \times t_4$	$t_8 = x_0 \times t_4$	$t_9 = t_6 \times t_5$
$t_{10} \!=\! t_3 \oplus t_9$	$y_0 = t_7 \oplus t_1$	$y_1 = (t_8 \oplus t_{10})'$
$t_{11} = t_3 + y_0$	$t_{12} = t_1 \oplus t_5$	$t_{13} = x_0 \oplus t_3$
$t_{14} = t_{12} \times t_{11}$	$t_{15} = t_{12} \times t_{10}$	$y_2 = (x_2 \oplus t_{14})'$
$y_3 = t_{15} \oplus t_{12}$	3	

Fig. 4. A 19-gate implementation of the non-linear S-box. The 4-bit inputs are $X = \{x_0, x_1, x_2, x_3\}$ and the 4-bit outputs are $Y = \{y_0, y_1, y_2, y_3\}$.

two-input gates required by both implementations. Results in NAND2 gate equivalent are also provided for each circuit for better comparison.

To observe the power and speed performance of both implementations, two 8-bit SBoNG number generator circuits are implemented using the original S-box and the optimised S-box respectively. The designs are synthesized for Intel FPGA Cyclone IV EP4CE6E22C8 using Quartus Prime Version 17.1.0. Table III tabulates the results for both designs.

VI. CONCLUSION

From the experimental results, we showed that using the proposed algorithm we were able to produce an equivalent design for the same RNG with 17.5% smaller in GE counts. The corresponding dynamic power for the final SBoNG was also reduced by 24%. Given that the areas of the individual SBoNG can potentially occupy up to 80% of the final SC application, the reduction of areas and power were deemed as substantial.

REFERENCES

 B. R. Gaines, "Stochastic computing," in *Proceedings of the April 18-20,* 1967, Spring Joint Computer Conference, New York, NY, USA, 1967, TABLE II

COMPARISON OF CIRCUIT SIZE BETWEEN THE ORIGINAL AND PROPOSED S-BOX

			CIRCUIT	S		
	Gate count					Gate
Work	AND	OR	XOR	XNOR	Total	equivalent (GE)
Original	22	12	0	0	34	80
This work	6	0	11	2	19	66

TABLE III FPGA IMPLEMENTATION RESULTS FOR BOTH 8-BIT SBONG NUMBER GENERATORS

GENERATORS						
Work	Logic	Fmax	Power (mW)			
	element (LE)	(MHz)	Static	Dynamic	Total	
Original	32	358.94	43.19	0.17	43.33	
This work	30	383.00	42.82	0.13	42.95	

AFIPS '67 (Spring), pp. 149-156, ACM.

- [2] A. Alaghi and J. P. Hayes, "Survey of stochastic computing," ACM Tans. Embed. Comput. Syst., vol. 12, no. 2s, pp. 92:1-92:19, May 2013.
- [3] B. D. Brown and H. C. Card, "Stochastic neural computation. i. computational elements," *Computers, IEEE Transactions on*, vol. 50, no. 9, pp. 891-905, Sep 2001.
- [4] A. Dinu, M. N. Cirstea, and M. McCormick, "Stochastic implementation of motor controllers," in *Industrial Electronics*, 2002. *ISIE 2002. Proceedings of the 2002 IEEE International Symposium on*, 2002, vol. 2, pp.639-644.
- [5] W. Qian, X. Li, M. D. Riedel, K. Bazargan, and D. J. Lilja, "An architecture for fault-tolerant computation with stochastic logic," *Computers, IEEE Transactions on*, vol. 60, no. 1, pp. 93-105, Jan 2011.
- [6] A. Alaghi, C. Li, and J. P. Hayes, "Stochastic circuits for real-time image-processing applications," in *Design Automation Conference* (DAC), 2013 50th ACM/EDAC/IEEE, May 2013, pp. 1-6.
- [7] P. Li and D. J. Lilja, "Using stochastic computing to implement digital image processing algorithms," in *Computer Design (ICCD), 2011 IEEE* 29th International Conference on, Oct 2011, pp. 154-161.
- [8] A. Naderi, S. Mannor, M. Sawan, and W. J. Gross, "Delayed stochastic decoding of ldpc codes," *Signal Processing, IEEE Transactions on*, vol. 59, no. 11, pp. 5617-5626, Nov 2011.
- [9] Y. Chang and K. K. Parhi, "Architectures for digital filters using stochastic computing," in *Acoustics, Speech and Signal Processing* (*ICASSP*), 2013 IEEE International Conference on, May 2013, pp. 2697-2701.

- [10] K. K. Parhi and Y. Liu, "Architectures for iir digital filters using stochastic computing," in *Circuits and Systems (ISCAS)*, 2014 IEEE International Symposium on, June 2014, pp. 373-376.
- [11] Y. Liu and K. K. Parhi, "Lattice fir digital filter architectures using stochastic computing," in Acoustics, Speech and Signal Processing (ICASSP), 2015 IEEE International Conference on, April 2015, pp. 1027-1031.
- [12] N. Saraf, K. Bazargan, D. J. Lilja, and M. D. Riedel, "Iir filters using stochastic arithmetic," in *Design, Automation and Test in Europe Conference and Exhibition (DATE), 2014*, March 2014, pp. 1-6.
- [13] W. Qian, X. Li, M. D. Riedel, K. Bazargan, and D. J. Lilja, "An architecture for fault-tolerant computation with stochastic logic," *IEEE Transactions on Computers*, vol. 60, no. 1, pp. 93-105, Jan 2011.
- [14] F. Neugebauer, I. Polian, and J. P. Hayes, "S-box-based random number generation for stochastic computing," *Microprocessors and Microsystems*, vol. 61, pp. 316-326, 2018.
- [15] M. Gay, J. Burchard, J. Horácek, A. M. Ekossono, T. Schubert, B. Becker, M. Kreuzer, and I. Polian, "Small scale aes toolbox: algebraic and propositional formulas, circuit-implementations and fault equatuions," 2016.
- [16] M. Sipser, Introduction to the Theory of Computation, International Thomson Publishing, 1st edition, 1996.
- [17] J. Boyar, R. Peralta, and D. Pochuev, "On the multiplicative complexity of Boolean functions over the basis (∧, ⊕, 1)," *Theoretical Computer Science*, vol. 235, no. 1, pp. 43-57, 2000.
- [18] J. Boyar and R. Peralta, ^{ter}Tight bounds for the multiplicative complexity of symmetric functions," *Theoretical Computer Science*, vol. 396, no.1, pp. 223-246, 2008.
- [19] C. P. Schnorr, "The multiplicative complexity of Boolean functions," in *Applied Algebra, Algebraic Algorithms and Error-Correcting Codes*, Teo Mora, Ed., Berlin, Heidelberg, 1989, pp. 45-58, Springer Berlin Heidelberg.
- [20] M. T. Sönmez and R. Peralta, "The multiplicative complexity of Boolean functions on four and five variables," in *Lightweight Cryptography for Security and Privacy*, Thomas Eisenbarth and E. Öztürk, Eds., Cham, 2015, pp. 21-33, Springer International Publishing.
- [21] J. Boyar, P. Matthews, and R. Peralta, "Logic minimization techniques with applications to cryptology," *Journal of Cryptology*, vol. 26, no. 2, pp. 280-312, 2013.
- [22] D. Canright, A Very Compact S-Box for AES, pp. 441-455, Springer Berlin Heidelberg, Berlin, Heidelberg, 2005.
- [23] R. C. Peralta and J. Boyar, "Method of optimizing combinational circuits," Nov. 20 2012, US Patent 8,316,338.
- [24] J. Boyar and R. Peralta, "A small depth-16 circuit for the aes s-box," in Information Security and Privacy Research, Dimitris Gritzalis, Steven Furnell, and Marianthi Theoharidou, Eds., Berlin, Heidelberg, 2012, pp. 287-298, Springer Berlin Heidelberg.
- [25] J. J. Tay, M. L. D. Wong, M. M. Wong, C. Zhang, and I. Hijazin, "Low multiplicative complexity logic minimisation over the basis (and, xor, not)," *Electronics Letters*, vol. 52, no. 17, pp. 1438-1440, 2016.
- [26] J. J. Tay, M. L. D. Wong, M. M. Wong, C. Zhang, and I. Hijazin, "A tree search algorithm for low multiplicative complexity logic design," *Future Generation Computer Systems*, vol. 83, pp. 132-143, 2018.

Banana (*Musa Acuminata Triploid AAA*, *Cavendish*) Sweetness Measurement by Digital Image Processing Technique

P. Chomtip, V. Adhisaya, L. Kanhokthorn, L. Phoptorn

Abstract— The objective of this project is to develop a computer system which can measure banana sweetness by using the digital image processing technique. The system is called "Banana Sweetness Measurement by Digital Image Processing or BSMDIP. The system consists of 5 subsystems, which are 1) image acquisition, 2) image preprocessing, 3) image evaluation, 4) displaying result, and 5) verification. The BSMDIP was conducted on 5 Thai banana species, which are 1) Musa acuminata Triploid AAA, Cavendish banana or Kluai Hom Thong in the Thai language, 2) Musa x paradisiaca Triploid ABB or Kluai Namwa, 3) Musa acuminata Diploid AA or Kluai Khai, 4) Musa acuminata Diploid AA or Kluai Lep Mu Nang, and 5) Musa ABB paradisiaca or Kluai Hak Muk. The system was evaluated on more than 1,600 banana peel images by using both color and texture features. The accuracy rate of the system is around 80.18 per cent, with the average access time of 1.51 seconds per image.

Index Terms— banana; banana sweetness; digital image processing; fruit sweetness; sweetness measurement

I. INTRODUCTION

B ananas are probably one of the world's oldest cultivated fruits. Not only could parts of banana trees be used in many traditional ceremonies, but the banana fruit also has high nutritive value, which can be eaten for a quick fix for flagging energy levels or as a diet meal in a day.

Bananas play an important role in Thai culture. Thai people use bananas, ranging from the tree, flower, to the leaf. Bananas are available in Thailand all year round. They come in many different shapes, sizes, and types. However, the most popular banana that is most commonly found in the fruit market is the 'Kluai Hom Thong' or fragrant banana.

Since Bananas are an important product for domestic and international export, Thai farmers will have to supply bananas that are not yet ripe for export purpose so that they reach the destination in a ripe state. However, farmers do not have an accurate method to determine the sweetness (ripeness) of each banana, except that they can only estimate them with their eyes.

Manuscript received July 23, 2018.

Chomtip Pornpanomchai is with the Faculty of ICT, Mahidol University Thailand (corresponding author to provide phone: 662-4410909; fax: 662-4410808; e-mail: chomtip.por@mahidol.ac.th).

Adhisaya Vajarobola is with the Faculty of ICT, Mahidol University Thailand (e-mail: adhisaya.vaj@student.mahidol.ac.th).

Kanhokthorn Leadkeattiwong is with the Faculty of ICT, Mahidol University Thailand (e-mail: kanhokthorn.lea@student.mahidol.ac.th).

Phoptorn Limpitigranon is with the Faculty of ICT, Mahidol University Thailand (e-mailphoptorn.lim@student.mahidol.ac.th).

There are many projects finding the banana ripeness from color and the outer appearance; however, there are few projects which talk about judging the banana sweetness from the color of its peel. That is why this banana sweetness evaluation system aims to create a system that predicts the sweetness based on the data that is stored in the database and the color of the banana peel. Initially, this project focuses only on Kluai Hom Thong. This project does not include other kinds of bananas in the database, since each banana species is different in size, color, sweetness and taste. The project objective is to develop a computer application to help farmers determine the sweetness of bananas by using a non-destructive method. The system evaluates banana sweetness by the color and texture of banana peel.

II. LITERATURE REVIEW

There are 2 main techniques of measuring the fruit sweetness, which are the destructive and non-destructive measurement. The brief details of each technique are as follows.

A. Destructive fruit sweetness measurement

Normally, researchers use a refractometer to measure fruit sweetness in the °Brix unit. A refractometer is a defective fruit and vegetable sweetness instrument. [1] It consists of a monocular eyepiece, which has a rugged exterior of plastic with metal to protect the optical lenses, as shown in Figure 1 (a). The refractometer woks by measuring the refractive index of light passing through liquid bends. Researchers test the fruit sweetness by dropping fruit juice on a main prism of the refractometer, as shown in Figure 1 (b). Then look into the eyepiece in light source direction, as shown in Figure 1 (c). Finally, read a number where the blue and white colors meet for the °Brix result, as shown in Figure 1(d). [2]



Fig. 1. Illustrate how to use a refractometer

Bumgarner and Kleinhenz demonstrated how to use a refractometer to measure cucumbers, leafy greens, sweet corn,

tomatoes and watermelons. [3] Iswari developed a mobile application to measure the sweetness of an apple, banana and melon. [4] Harker, et al, used a refractometer to evaluate sucrose in an apple (7-14 °Brix) and used gas chromatography to analyse malic acid in an apple (0.08-0.2 % wt./v). [5]

B. Non-destructive fruit sweetness measurement

Some researchers demonstrated two techniques of common non-destructive fruit sweetness measurement, which are an electronic sensing technique and an image processing technique, to measure fruit sweetness. Each technique has the following details.

1) The electronic sensing technique

Soltani, et al, evaluated banana quality during a ripening stage by using a capacitance sensing system. The system was equipped by a pair of parallel capacitor sensors. The system showed 7.8 ° Brix for a green banana and 18.6 ° Brix for a full-ripe banana. [6] Sankhe used a capacitive sensing system to inspect a banana for its ripeness level. Based on the experiment, it was concluded that the ripeness of bananas is directly proportional to the capacitance value. [7] Jamaludin, et al, employed an impedance analyzer board AD5933 with electrocardiogram (ECG) probe to measure the impedance value in a banana. The experiment shows that a ripe banana has an impedance value between 20.1 - 30.1 kHz, with 1.258 - 0.769 °Brix, while an unripe banana has a higher impedance value than a ripe and overripe banana. [8] Jam and Chia used five near infrared bands to classify the pineapple quality. [9]

2) The image processing technique

Many researches measure fruit sweetness by using various image features, namely: color, texture, histogram, shape, and size. Jackman and Sun used image textures (contrast, correlation, entropy) for food quantity assessment. [10] Ittatirut, et al, developed a computer system to predict apple sweetness by using RGB and HSV color features. The developed system has the precision rate of 79.03 per cent. [11] Bayarri, et al, employed the L*a*b*C*h* color features to measure peach, orange, kiwifruit and berry sweetness, with the °Brix sweetness being around 12.5, 7.6, 7.6 and 7.8, respectively. [12] Salvador, et al, investigated the relationship between image color, texture features, and storage time that affected the banana (Musa Cavendish AAA group and Musa Parasdisiaca L. AAB group) quality. The color and texture of bananas are variedly changed depending on the storage period. [13] Bejo and Kamaruddin determined Chokanan mango (Mangifera indica) sweetenss by using color (HSB, RGB) features. The experimental results show that the Chokanan mango has 4.0-17.0 °Brix sweetness.

Based on the previous related researches, this research adopts the image processing with color and texture features to predict the banana sweetness. The system analysis, design and implementation are presented in the next section.

III. METHODOLOGY

This part describes the process of analysis and design, which is presented by the BSMDIP conceptual diagram and system structure chart. The details of each element are described below.

A. System conceptual diagram

The BSMDIP starts with a user taking a banana image by using a mobile phone camera. Then the banana image is submitted to a computer system for evaluating the sweetness. After that, the system compares the banana image with all banana images in the system database. Then the BSMDIP displays the evaluation results. Finally, the user can verify banana sweetness by using a refractometer and update a new banana image with the sweetness in °Brix value into the BSMDIP system database, as shown in Figure 2.



Fig. 2. The BSMDIP system conceptual diagram

B. System Structure Chart

For better understanding, the system structure chart, which maps to the system conceptual diagram is shown in Figure 3. The structure chart consists of 5 subsystems, namely: 1) image acquisition, 2) image preprocessing, 3) image evaluation, 4) displaying result, and 5) verification. Each subsystem has the following details.



Fig. 3. The BSMDIP system structure chart



Fig. 4. Sample of Musa acuminata Triploid AAA, Cavendish banana image

1) *Image acquisition* - This module takes a banana image in the box with all black color side for controlling the system environment, and then submits the taken image to the BSMDIP system server computer. The sample of banana image taken in this module is shown in Figure 4.

2) *Image Preprocessing* - This module processes the banana image in three sub-modules, which are 1) image segmentation, 2) color extraction, and 3) texture extraction. Each sub-module has the following details.

a) *Image segmentation* – This sub-module segments the banana image into foreground and background. The banana is

the foreground and the black box color is the background. The BSMDIP system uses only the foreground image to process the sweetness evaluation

b) *Image color extraction* – the BSMDIP system extracts RGB and HSI of the banana image by calculating the mean value of red-green-blue color, including the mean value of hue-saturation-intensity.

c) *Image texture extraction* – The image texture is calculated based on the Gray-Level Co-occurrence Matrix (GLCM). The contrast, correlation, energy and homogeneity are four banana image textures, which are extracted in this sub-module. Each texture value has the following details [15].

Giving P i,j is the entry in a normalized gray-tone spatial-dependence matrix and N is the number of distinct gray levels in the quantized image.

The contrast texture feature measures the local variations in the GLCM. The contrast texture can be calculated using Equation 1:

$$\sum_{i,j=0}^{N^{-1}} P_{i,j}(i-j)^{2}$$
(1)

The correlation texture feature measures the linear dependency of the gray levels on those of neighboring pixels. The correlation texture can be calculated using Equation 2:

$$\sum_{i,j=0}^{N-1} P_{i,j} \left(\frac{(i-\mu_i)(j-\mu_j)}{\sigma_i \sigma_j} \right)$$
(2)

Where μi , μj , $\sigma i \& \sigma j$ = the mean and standard deviations of $p_{j, j}$. The energy texture feature measures the uniformity of the texture that is pixel pair represent. High energy values occur when the GLCM distribution has a constant or periodic form. The energy texture can be calculated using Equation 3:

$$\sum_{i,j=0}^{n-1} (P_{i,j})^2 \tag{3}$$

The homogeneity texture feature measures the closeness of the distribution of elements in the GLCM to the GLCM diagonal. The homogeneity texture can be calculated using Equation 4:

$$\sum_{i,j=0}^{N-1} \frac{Pi, j}{1 + (i-j)^2}$$
(4)

3) *Image Evaluation* - This function is used to compare the detail of an input banana image with the information in the database by using the Euclideance distance method, using Equation 5:

$$\sqrt{\sum_{i=1}^{n} (q_i - p_i)^2}$$
(5)

Where q_i is the features of a banana image and p_i is the training features in the BSMDIP database and N is the total number of training datasets.

4) *Displaying Result* - This sub-module shows the banana sweetness evaluation result. The graphic user interface (GUI) is shown in Figure 5, which has of the following details.

a) Display image box – There are 2 display image boxes, which are the input banana image box (label number 1) and display results box (label number 2).

b) *Command button* – There are 3 command buttons, which are 1) open test image (label number 3), 2) analyze sweetness button (label number 4), and 3) save verification of banana sweetness button (label number 5).

c) *Display value box* – There are 5 display value boxes, which are 1) display mean of RGB values box (label number 6), 2)

display mean of HSI values box (label numbar 7), 3) display mean of texture values box (label numbar 8), 4) banana sweetness prediction value box (label number 9), and 5) banana sweetness verification value, which is measured by a refractometer (label number 10).



Fig. 5. The GUI of the BSMDIP system

5) Verification

Farmers or users can verify banana sweetness by using a refractometer with the system results. Moreover, they can update new banana images with the °Brix sweetness values into the BSMDIP database. The system database stored in the Excel file consisting of 13 fields, which are: 1) banana image filename, 2) contrast texture, 3) correlation texture, 4) Energy texture, 5) homogeneity texture, 6) mean hue, 7) mean saturation, 8) mean intensity, 9) mean red, 10) mean green, 11) mean blue, 12) °Brix sweetness, and 13) access time.

IV. EXPERIMETAL RESULTS

The system was developed on a computer with an Intel Core i7-7500U central processing unit 2.70 GHz (Intel's headquarters in <u>Santa Clara, CA.</u> USA). An operating system used the Windows 10 Pro (Microsoft Corp.; Redmond, WA, USA). MATLAB R2017b. (License Number 40598465, MathWorks, Natick, MA, USA) A smartphone digital camera used in this research was the Samsung Galaxy Note 5 (Sumsung Corp.; South Korea).

The BSMDIP system was conducted on 1,682 banana images, which consisted of 1,226 banana images for a training dataset, 328 banana images for testing dataset and 128 banana images for an un-training dataset. The system was tested with Kluai Hom Thong (shown in Figure 6 (a)) for the training and un-training dataset. The precision rates of the BSMDIP system are 99.27 (1,217/1,226*100), 80.18 (263/328*100) and 31.25 (40/128*100) per cent for training, un-training and unknown dataset, respectively, as shown in Table 1.

TABLE I. THE PRECISION RATES OF BSMDIP SYSTEM

Banana Image Type	No. Testing	Prediction Right	Prediction Wrong
Training Dataset	1,226	1,217	9
Testing Dataset	328	263	65
Unknown Dataset	128	40	88

For an unknown dataset, the system was conducted on four different species of bananas, which are 1) Kluai Namwa (*Musa x paradisiaca Triploid ABB*), 2) Kluai Khai (*Musa acuminata Diploid AA*), 3) Kluai Lep Mu Nang (*Musa acuminata Diploid AA*), and 4) Kluai Hak Muk (*Musa ABB paradisiaca*), as shown in Figure 6 (b), (c), (d) and (e), respectively. The precision rates

of the un-training dataset are 23.53 (8/34*100), 78.13 (25/32*100), 0.0 (0/26*100) and 19.44 (7/36*100) per cent for Kluai Namwa, Kluai Khai, Kluai Lep Mu Nang, and Kluai Hak Muk, respectily (as shown in Table 2).



Fig. 6. Image of (a) Kluai Hom Thong (b) Kluai Namwa, (c) Kluai Khai, (d) Kluai Lep Mu Nang, and (e) Kluai Hak Muk,

TABLE II. THE BANANA SPECIES OF UNKNOWN DATASET

Banana S	pecies	No. Testing	Prediction Right	Prediction Wrong
Kluai Nai (Musa x p	mwa paradisiaca Triploid ABB)	34	8	26
Kluai Kha (Musa ac	ai uminata Diploid AA)	32	25	7
Kluai Leb (Musa act	o Mu Namg uminata Diploid AA)	26	0	26
Kluai Hul	k Muk 8 naradisiaca)	36	7	29

Based on the experimental results, the BSMDIP classifies banana sweetness in 7 levels. Each level has the relationship between sweetness, taste, texture and color, as the following: (as shown in Table 3)

TABLE III. THE RELATIONSHIP BETWEEN BANANA °BRIX, TASTE, TEXTURE

° Brix	Taste	Texture	Color
0-5	bitter	solid	all green
6-10	dull	hard	light green
11-15	bland	tough	half green half yellow
16-20	slightly sweet	firm	yellow more than green
21-25	sweet	tender	yellow with green tip & neck
26-30	super sweet	soft	light yellow
> 30	extremely sweet,	squishy	all yellow with brown flecks

V. CONCLUSION

The research paper fulfills the research objective, which is to develop a computer system plus a mobile camera, which is able to measure banana sweetness. The system is called "Banana Sweetness Measurement by Digital Image Processing or BSMDIP". The developed system has the precision rate of 80.18 per cent for un-training banana images with the average access time of 1.51 seconds/image.

Based on the experiment, the system classifies the relationship between banana sweetness, taste, texture and color features, as shown in Table 3. The raw banana with low sweetness (0 - 5 °Brix) has the all green color and the mature banana with all yellow color has sweet taste of around 21-25 °Brix.

Thailand has more than 30 kinds of fruits, which are harvest all year round. A list of Thai fruits includes 1) mango, 2) mangosteen, 3) rambutan, 4) durian, 5) pineapple, 6) papaya, 7) dragon fruit, 8) guava, 9) pomelo, 10) rose-apple, 11) jack-fruit, 12) custard apple, 13) langsat, 14) longan, 15) lychee, 16) sapodilla, 17) coconut, 18) banana, 19) snake fruit, 20) watermelon, 21) star fruit, 22) plum mango, 23) pomegranate, 24) kraton, 25) makok, 26) salacca, 27) longgong, 28) strawberry, 29) tangerine and 30) sweet orange, etc. Therefore, Thai people have different kinds of fruits every month. Moreover, this research will encourage Thai framers to evaluate fruit sweetness with non-destructive measurement. A mobile phone and a set of computer are simple equipment for, which is easy for farmers to seek. The farmers need to create the fruit database before evaluating the sweetness of each type of fruit.

VI. ACKNOWLEDGMENT

The authors thank the Faculty of Information Communication Technology, Mahidol University, Nakhon Pathom, Thailand for supporting the research.

VII. REFERENCES

- R.Harrill, "Using a refractometer to test the quality of fruits & Vegetables", Perfect blend organics, 1998:1-24.
- [2] Y.S.Jasmine, "Comparison of sugar content in bottled 100% fruit juice versus extracted juice of fresh fruit", Food and Nutrition Science, 2012(3):1509-1513.
- [3] N.R. Bumgarner, and M.D. Kleinhenz, "Instructions for measuring "Brix in cucumber, leafy greens, sweet corn, tomato, and watermelon", Fact sheet, Agriculture and Natural Resources, The Ohio State University, 2012:1-12.
- [4] N.M.S. Iswari, and W. Ranny, "Fruitylicious: mobile application for fruit ripeness determination based on fruit image", International conference on human system interaction, 17-19 July 2017, Ulsan, Korea, pp. 183-187.
- [5] F.R.Harker, K.B.Marsh, H. Young, S.H.Murry, F.A.Gunson and S.B.Walker, "Sensory interpretation of instrumental measurements 2:sweet and acid taste of apple fruit", Postharvest Biology and Technology, 2002 (24):241-250.
- [6] M. Soltani, R. Alimardani and M. Omid, "Prediction of banana quality during ripening stage using capacitance sensing system", Australian Journal of Crop Science, 2010 (6): 443-447.
- [7] D. Sankhe, "Ripeness inspection system for banana", International Journal of Computer Applications, 2015 (6): 0975-8887.
- [8] D. Jamaludin, S. A. Aziz and N.U.A. Ibrahim, "Dielectric based sensing system for banana ripeness assessment", International Journal of Environmental Science and Development, 2014 (5):286-289.
- [9] M.N.H.Jam, and K.S.Chia, "A five band near-infrared portable sensor in nondestructively predicting the internal quality of pineapples", International conference on signal processing & its applications, 10-12 March 2017, Penang, Malaysia, pp.135-138.
- [10] P.Jackman and D.W.Sun, "Recent advances in image processing using image texture features for food quality assessment" Trends in Food Science & Technology 2013(29):35-43.
- [11] T.Ittatirut, A. Lekhalawan, W.Tangjitwattanakorn and C.Pornpanomchai," Apple sweetness measurement by image processing technique", ICT International student project conference, Nakhon Pathom, Thailand, 27-28 May 2016, pp. 1-4.
- [12] S. Bayarri, C.Calvo, E.Costell and L.Duran, "Influence of color on perception of sweetness and fruit flavor of fruit drinks", Food Science Technoloy International, 2001;7(5):399-404.
- [13] A.Salvador, T.Sanz, and S.M.Fiszman, "Changes in color and texture and their relationship with eating quality during storage of two different dessert bananas", Postharvest Biology and Technology, 2007 (43):319-325.
- [14] S.K.Bejo, and S.Kamaruddin, "Determination of Chokanan mango sweetness (*Mangifera indica*) using non-destructive image processing technique", Australian Journal of Crop Science, 2014 8(4):475-480.
- [15] H. Zareiforoush, S. Minaei, M.R. Alizadeh, and A. Banakar, "Potential applications of computer vision in quality inspection of rice: a review", Food Ethics Rev. 2015(7), 321e345.

An Effective Feature Selection and Classification Model for High Dimensional Big Data Sets

Dingkun Li, Keun Ho Ryu^{*}, *Member, IEEE*, Erdenebileg Batbaatar, Hyun Woo Park, Seon Phil Jeone^{*}, Zhou Ye^{*}

Abstract—Never before in history is the data growing at such a high volume, variety and velocity. It not only provides multi-sources of information but also generates high dimensional data sets. As the dimensionality of the input data space (i.e., the number of predictors) increases, it becomes exponentially more difficult to find global optima for the parameter space, i.e., to fit models. This is one phenomenon of the "curse of dimensionality". Intuitively, data reduction provides a promising way, feature selection is a one method of data reduction. This paper proposes an effective model for feature selection and classification. It takes advantage of the divide-and-conquer idea combined with MapReduce paradigm to handle high dimensional big data sets. Several well-known classification methods have been used and compared in our experiment and the result illustrates that the proposed model outperforms all selected methods. And this model shows great power and flexibility for handling non-big data as well.

Index Terms—Feature selection, classification, high dimensional data, MapReduce

I. INTRODUCTION

T HE society is changing dramatically. Never before in history is the data growing at such a high volume, variety and velocity. It not only provides multi-sources of information

Manuscript was received on June 30, 2018. This work was supported by the Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Science, ICT & Future Planning (No.2017R1A2B4010826) and the MSIP(Ministry of Science, ICT and Future Planning), Korea, under the ITRC(Information Technology Research Center) support program (IITP-2018-2013-1-00881) supervised by the IITP(Institute for Information & communication Technology Promotion). And Karamay Central hospital (2017HZ006A).

D. Li was with Chungbuk National University, Chungbuk, Korea. He is now with the post-doctor station of Shanghai Tongji University cooperated with Karamay Central Hospital, Xinjiang China. (e-mail: 33644251@qq.com).

Corresponding author K.H. Ryu is with the Chungbuk National University, Chungbuk, Korea. (phone: +82-10-4930-1500; fax: +82-43-275-2254; e-mail: khryu@dblab.chungbuk.ac.kr).

Corresponding author S. P. Jeone is with the Beijing Normal University-Hong Kong Baptist University United International College (UIC), Guangdong, China. (phone: +86-135-3658-5425; e-mail: spjeong@uic.edu.hk).

Corresponding author Z. Ye is with the Karamay Central Hospital, Xinjiang China. (phone: +86-139-9930-7668; e-mail: yezhou126@126.com).

E. Batbaatar, H. W. Park are with the Chungbuk National University, Chungbuk, Korea. (e-mail: {eegii, hwpark}@dblab.chungbuk.ac.kr).

but also generates high dimensional data sets. It is becoming more common in many practical applications such as medical, microarray, business and etc. data analysis. Especially, microarray data, which has tens of thousands of features with small sample size, obtained from public API, is typical high dimensional data set.

As the dimensionality of the input data space (i.e., the number of predictors) increases, it becomes exponentially more difficult to find global optima for the parameter space, i.e., to fit models. This is one phenomenon of the "curse of dimensionality". The curse of dimensionality is a term introduced by Bellman to describe the problem caused by the exponential increase in volume associated with adding extra dimensions to Euclidean space [1]. This problem affects many aspects of data analytics, such as data preprocessing, classification, clustering and etc. The work [2] declares that the curse of dimensionality leads to invalid conclusions by some commonly used clustering methods. The work [3] declares that the prediction result of the accuracy will be unstable over the variations in the training set, especially in high dimensional data. A striking result has been found that the simple and popular Fisher linear discriminant analysis can be as poor as random guessing as the number of features increasing [4], [5].

Intuitively, feature selection provides a promising way for dealing with "curse of dimensionality". There were a lot of works [6], [7], [8], [9], [10], [11], [12]that have been done during the last decades of years. Nevertheless, most of them are lack of ability to handle big high dimensional data set, and the others describ some methods can only be used for the specified data sets.

The purpose of our work is to develop an effective feature selection and classification model taking advantage of statistics, data mining and MapReduce techniques on high dimensional big data set.

The rest of the paper is organized as follows: We describe the related work in Section 2. An overview of our system will be introduced in Section 3 and its design detail will be described in Section 4, which is followed by experiment result in Section 5. Finally, Section 6 concludes our work and predicts future work.

II. RELATED WORK

This section briefly describes the related platform, algorithms and some key techniques we used in our work.

A. Hadoop and Spark

Hadoop consists of HDFS (Hadoop Distributed File System), HBase, and Hadoop MapReduce which can analyze big data [13]. It is an open source framework that writes and implements an application program for processing big data.

B. Data Mining (DM)

Data mining is the process of discovering interesting patterns and knowledge from large amounts of data. The data sources include databases, data warehouses, the Web, other information repositories, or data that are streamed into the system dynamically [14].

Also, send a sheet of paper or PDF with complete contact information for all authors. Include full mailing addresses, telephone numbers, fax numbers, and e-mail addresses. This information will be used to send each author a complimentary copy of the journal in which the paper appears. In addition, designate one author as the "corresponding author." This is the author to whom proofs of the paper will be sent. Proofs are sent to the corresponding author only.

C. Feature Selection (FS) Methods

Feature selection aims at finding the most relevant features of a problem domain. It is one essential step for data mining which is defined as a multidisciplinary task to find out hidden nugget of information from data [15]. Primarily, there are three kinds of feature selection methods, filters, wrappers and embedded methods. The filter methods work fast but its result is not always satisfactory. While the wrapper methods guarantee good results but very slow when applied to wide feature sets which contain thousands or even hundreds of thousands number of features. The third one is embedded methods which reduce the computation time taken up for reclassifying different subsets which is done in wrapper methods.

Following figures compare the difference between filter, wrapper, and embedded methods [15]. The procedure of filter method is shown in Fig. 1.





Filter methods are generally used as a preprocessing step. The selection of features is independent of any machine learning algorithms. Instead, features are selected on the basis of their scores in various statistical tests for their correlation with the outcome variable. The correlation is a subjective term here.

The procedure of wrapper method is shown in Fig. 2.

In wrapper methods, it uses a subset of features and train a model using them. Based on the inferences that it draws from



Fig. 2. Wrapper method analysis procedure

the previous model, it decides to add or remove features from subset. The problem is essentially reduced to a search problem. These methods are usually computationally very expensive.

The procedure of embedded method is shown in Fig. 3.



Fig. 3. Embedded method analysis procedure

Embedded method combines the qualities of filter and wrapper methods. It's implemented by algorithms that have their own built-in feature selection methods.

The benefits for feature selection are reducing data analysis complexity and improve data analysis performance. Besides, there are more benefits such as accuracy improvement, expenditure reduction, etc.

III. FRAMEWORK OF PROPOSED MODEL

The brief framework of these two steps is given in Fig. 4.



Fig. 4. Proposed model framework

Firstly, after basic preprocessing step, the method integrates data from multi-source in term of year and region. Data sets of all data source should consist similar feature space so that they can be integrated without conflict. Or these data sets should be processed to achieve similar feature space for further integration.

Secondly, clinical data is natural to have a big amount of features and big amount of missing value. There are many ways to handle missing value such as using mean, median or user specified value to take place of the missing value. However, to our concern, the feature consists more than 80% of missing value that should be eliminated because too much missing value will reduce the importance of the feature. If a feature is blank or only has one value, it holds no meaning.

Thirdly, feature families (sub feature sets) are generated based on the domain or correlation between each other, or random combination of them.

Fourthly, feature selection methods are used to select features according to a certain kind of criteria such as voting or weight. In this paper, voting strategy has been used to generate the feature candidates. Next, only important features are selected and combined for further model training.

Finally, the prediction model will be generated and the result will be compared with existing widely used classification methods.

More detail will be described in the following sections.

IV. MODEL IMPLEMENTATION

Implementation detail is described in this section.

A. Feature Family Generation

Feature family is also called feature subset or column family. It is defined as:

$$\mathbf{X}_{i} = \{\mathbf{f}_{i1}, \mathbf{f}_{i2}, \dots, \mathbf{f}_{ii}, \dots, \mathbf{f}_{im}\}, i \in \{1, 2, \dots, n\}, j \in \{1, 2, \dots, m\}$$
(1)

Where n is the number of the families, m is the number of the features in each family. Feature family is a group of features generated from the whole feature space. But they are mutual exclusive, which means none of the feature in one family also belongs to another one.

There are three ways to generate the feature families, they are based on the feature domain, correlation or random combination. Features of the same domain are categorized into one family, such as name, sex, height, WC etc. which belong to observer's basic information family. For other dataset, which has no explicit domain information to categorize the features, some specified criteria such as correlation based methods or random feature combination methods are used to generate the feature families.

B. Feature Selection on Each Family

Two filter methods: Information Gain [16], ReliefF [17] and One wrapper method: Linear-SVM [18] are chosen and modified to be used on our data sets since the performance of these three methods are better than other feature selection algorithms according to our experiment. They are used for redundant feature elimination and candidate feature generation.

After obtained the result from these algorithms, several important features can be selected from each feature family. A voting strategy has been used to generate the candidates in each feature family. The idea is that the feature which has more than 2 votes among 3 methods is selected as the candidate for next step.

C. Feature Integration

All features selected from different feature families are integrated based on the primary keys to generate features candidates for classification method. Full outer join (one data integration operation among the four widely used operations which are left join, right join, inner join and full outer join) has been used to merge all selected features together, the joint result will be used for next step.

D. Classification and Result Evaluation

Rule-based decision tree classification method CHAID has been applied on feature candidates. The reason using this method is that, it achieves highest accuracy among other classification methods at this step. Besides, rule based method provides an easy and effective way for disease interpretation and prediction.

The format of disease rules is same as the IF-THEN rule; for example, IF (edu = elementary, B1 < = 0.86 mg/day, married), THEN (hypertension = yes). The purpose here is the mining of all of the disease-related rules from the training dataset for a further data analysis including disease prediction.

Nevertheless, there is no method that fits for all kinds of the data sets that can achieve highest accuracy and efficiency. One advantage for proposed method, at this step, the classification method can be replaced by other classification methods if some other data sets are used.

In order to compare efficiency of the proposed method, several existing methods such as LR, K-NN, AdaBoost, CHAID are used to compare with the proposed method. The result will be given in the experimental section.

E. Apply Proposed Model in MapReduce Environment

The proposed model is suitable to be applied in a parallel environment due to its distributed design. It be applied on both high dimensional big data sets and normal data sets.

Fig. 5 describes the procedure of the MapReduce framework.



Fig. 5. Framework of RFD-DP model running on MapReduce

It depicted from the Fig. 5 that MapReduce consists of mapping phase and reducing phase. During the mapping phase, HBase server t (t=1, 2,...,n) splits the data into different blocks. Map worker is assigned by master to process the data in this block. Its work includes feature family generation and feature

selection. During the reducing phase, the output of worker node is shuffled and sorted regarding to year and stored in intermediate data block D_m . Reduce worker is assigned to generate feature candidates of each feature family, integrates these candidates regarding to *id* to generate the training data set, and executes CHAID classification algorithm to train the classifier. The pseudo code is given in Algorithm 1.

	Algorithm 1	1. RFD-D	P model	l algorithm	running	on MapReduc
--	-------------	----------	---------	-------------	---------	-------------

Algori	ithm 1. RFD-DP model algorithm running on MapReduce	ru
Inpu	t: data set S,	•
Outp	ut: Classifier C	
Proc	edure:	
1	Mapper()	
2	For a specified disease: generate N feature families from S $\{\{id,X_1\}, \{id,X_2\},, \{id,X_k\},, \{id,X_N\}\}, k \in N, T_1, T_2, T_3 \leftarrow \emptyset // T_1, T_2, T_3$ are temp feature sets	
3	for each X_k , do	ar
4	initiate feature set $F_I \leftarrow \emptyset$	T
5	$F_1 \leftarrow \text{IG-FS}(X_k), T_1 \leftarrow T_1 \bigcup F_1 // \text{Information gain based}$ feature selection method	_
6	initiate feature set $F_2 \leftarrow \emptyset$	V
7	$F_2 \leftarrow \text{ReliefF-FS}(X_k), T_2 \leftarrow T_2 \cup F_2 // \text{ReliefF based feature selection method}$	or
8	initiate feature set $F_3 \leftarrow \emptyset$	Pr
9	$F_3 \leftarrow \text{LSVM-FS}(X_k), T_3 \leftarrow T_3 \bigcup F_3 // \text{ linear-SVM based}$ feature selection method	K
10	end for	(R
11	return { <i>id</i> , <i>T</i> ₁ , <i>T</i> ₂ , <i>T</i> ₃ }	
12	Reducer()	
13	//generate feature candidate Fc using voting strategy	fr
14	$F_C \leftarrow \text{vote} (id, T_1, T_2, T_3)$	da
15	$F_L \leftarrow$ Integrate F_C based on the <i>id</i> // F_L is the integrated feature set	or cc
16	$C_m \leftarrow \text{CHAID}(F_C) // \text{CHAID}()$ could be replaced by other classification algorithms	pr
17	$r_m \leftarrow \text{evaluate}(C_m) // C_m \text{ is rule-based classifier, } r_m \text{ is the evaluation result of } C_m$	
18	return C_m, r_m	
19	vote()	
20	initiate feature set $F_s \leftarrow \emptyset$	A
21	for each $f_i \in F_j$, where $j \in \{1,2,3\}$ do	be
22	if $\exists f_i \in F_k$, count $(f_i) ++, k \neq j, k \in \{1,2,3\}$	te
23	if count $(f_i) > 2$, return add f_i to F_s	
24	end for	cł
25	return <i>F</i> _s	fo
		is

As it described from the algorithm that for a specified disease such as HTN, before the MapReduce procedure, the data set S has been divided into several subsets in terms of the feature families, such as (*id*, X_K). Each Map work node is a filter and wrapper based feature selection work node. Algorithms IG-FS(), RFD-DP-ReliefF() and LSVM-FS() run in this node, the output is the features selected from each feature family.

After that, MapReduce master will start Reducing procedure

to integrate selected features into one intermediate data block regarding to year. Then on Reduce work node, the method vote() is used to generate feature candidates from the output of previous step. Then rule-based classification algorithms will run in each Reduce worker node.

The result of each classification algorithm will be compared and the best one is chosen as the classification algorithm on current data sets. In our work, CHAID has been chosen as the le-based classification algorithm of the proposed model.

V. EXPERIMENT AND RESULT

Experimental detail is described in this section.

A. Data Preparation

Data used in this paper comes from three data sources, they re KNHANES data [19], orlraws10P [20] and Prostate [21]. he basic information has been summarized in Table 1.

TABLE 1	. BASIC INFORMATION	ON SELECTED DATA SET	S
---------	---------------------	----------------------	---

Dataset	Instances	Features	Classes	Size
KNHANES(6)	15,587	727	3	184 MB
orlraws10P	100	10,304	10	3.7 MB
Prostate	102	12,600	2	5.6 MB
KNHANES (R3-R6)	733,530	1,194	3	6.5 GB

One data set called KNHANES(6) comes from [19] dates om 2013 to 2015. Another data set called KNHANES(R3-R6) ates from 2009 to 2015 and randomly sampled 10 times in rder to generate a distinguishable relative bigger data set ompared with other 3 data sets.

The data preprocessing step has been described in our revious work [12].

B. Compare Proposed Model with Existing Classification Methods

Existing methods such as Logistic Regression (LR), daBoost, Naive Bayesian (NB) and Decision Tree (DT) have een used to compare with the proposed model in order to stify the efficiency of our model.

The reason why these classification methods have been nosen is that they are not only the most widely used methods or classification but also: 1. For LR, if the signal to noise ratio low, LR is likely to perform best. And it is the best method fit for both big data and small data set among other regression models in most cases. 2. For K-NN, it is an instance-based lazy learning method, can produce arbitrarily shaped decision boundaries fit for many kinds of data sets. 3. For AdaBoost, it achieves the best result among other boost algorithms. 4. For decision tree algorithms, CHAID is the best one among DT algorithms regarding our experiments.

We have run all these methods 3 times, the average

performance is given in Table 2.

From the Table 2 we can see that the proposed method achieves the highest score regarding AUC and F-score among all the other methods. However, K-NN achieves better result than the proposed method regarding to accuracy on KNHANES(6) data set. Nevertheless, the problem for accuracy is that it is not suitable to be used as evaluation standards of classifier since the classification is uneven [22]. Accuracy works best if false positives and false negatives have a similar cost. Predictive accuracy is a misleading performance measure for highly imbalanced data [23]. AUC holds more value based on this data set.

Dataset	Methods	Sensitivity	Specificity	Precision	Accuracy	F-score	AUC
KNHANES(6)	M_1	0.866	0.539	0.807	0.764	0.835	0.826
	M_2	0.847	0.522	0.798	0.745	0.821	0.807
	M ₃	0.708	0.838	0.617	0.803	0.660	0.742
	M_4	0.633	0.806	0.538	0.760	0.581	0.822
	PM	0.804	0.676	0.890	0.774	0.845	0.849
orlr	M ₁	1.0	0.0	0.909	0.909	0.952	0.990
aws10P	M_2	NA	NA	NA	NA	NA	NA
	M ₃	NA	NA	NA	NA	NA	NA
	M_4	1.0	1.0	1.0	1.0	1.0	1.0
	PM	1.0	1.0	1.0	1.0	1.0	1.0
Pro	M ₁	0.905	0.918	0.923	0.911	0.914	0.987
state	M_2	0.960	0.961	0.960	0.960	0.960	0.981
	M ₃	NA	NA	NA	NA	NA	NA
	M_4	0.819	0.951	0.961	0.872	0.885	0.943
	PM	0.974	0.960	0.960	0.970	0.970	0.990
KN	PM M ₁	0.974 NA	0.960 NA	0.960 NA	0.970 NA	0.970 NA	0.990 NA
KNHAN	PM M1 M2	0.974 NA NA	0.960 NA NA	0.960 NA NA	0.970 NA NA	0.970 NA NA	0.990 NA NA
KNHANES(PM M1 M2 M3	0.974 NA NA NA	0.960 NA NA NA	0.960 NA NA NA	0.970 NA NA NA	0.970 NA NA NA	0.990 NA NA NA
KNHANES(R3-1	PM M1 M2 M3 M4	0.974 NA NA NA NA	0.960 NA NA NA NA	0.960 NA NA NA NA	0.970 NA NA NA NA	0.970 NA NA NA NA	0.990 NA NA NA NA

•"NA" Not available, algorithm run failed on related data set

 \bullet "M1~M4" stand for LR, AdaBoost, K-NN, and CHAID classification methods separately

•"PM" stands for the proposed method

VI. CONCLUSION AND FUTURE WORK

In this paper, we have developed an effective feature selection and classification model for high dimensional data set. According to the experimental result, the proposed model achieves the best performance among other well-known classification methods on all data sets. Moreover, the key idea of our model is using divide-and-conquer strategy to handle high dimensional dataset, it perfectly matching the mechanism of MapReduce. Thus it can be used for high dimensional big data set.

In future, more advanced models are planning to be compared with ours model. For the sake of applying this model to the real practice, continues improvement of this model is needed as well regarding to it algorithms.

ACKNOWLEDGMENT

Dingkun Li thanks corresponding author Keun Ho Ryu, Seon Phil Jeone, Zhou Ye for their academical and financial support as well as Erdenebileg Batbaatar, Hyun Woo Park for their great help while writing this paper.

REFERENCES

- R. E. Bellman, and E. D. Stuart, "Applied dynamic programming", *Princeton university press*, 2015.
- [2] D. Y. Orlova, A. H. Leonore, and W. Guenther, "Science not art: statistically sound methods for identifying subsets in multi-dimensional flow and mass cytometry data sets", *Nature Reviews Immunology* 18.1 2018, p. 77.
- [3] H.J. Kim, S. C. Byong, and Y. H. Moon, "Booster in high dimensional data classification." *IEEE transactions on knowledge and data engineering*, 28.1, 2016, p. 29-40.
- [4] P. J. Bickel, and L. Elizaveta, "Some theory for Fisher's linear discriminant function, naive Bayes', and some alternatives when there are many more variables than observations", *Bernoulli* 10, no. 6, 2004, p. 989-1010.
- [5] J. Fan, and Y. Fan, "High dimensional classification using features annealed independence rules." *Annals of statistics*, 2008, p. 2605.
- [6] T. Abeel, et al, "Robust biomarker identification for cancer diagnosis with ensemble feature selection methods." *Bioinformatics*, 26.3, 2009, p. 392-398.
- [7] A. J. Ferreira, and A. F. MáRio, "Efficient feature selection filters for high-dimensional data." *Pattern Recognition Letters*, 33.13, 2012, p. 1794-1804.
- [8] I. Guyon, and E. André, "An introduction to variable and feature selection." *Journal of machine learning research*, 3, 2003, p.1157-1182.
- M. Hilario, and K. Alexandros, "Approaches to dimensionality reduction in proteomic biomarker studies", *Briefings in bioinformatics*, 9.2, 2008, p. 102-118.
- [10] J. Hua, D. T. Waibhav, and R. D. Edward, "Performance of feature-selection methods in the classification of high-dimension data", *Pattern Recognition*, 42.3, 2009, p. 409-424
- [11] Y. Piao, and K. H. Ryu, "A Hybrid Feature Selection Method Based on Symmetrical Uncertainty and Support Vector Machine for High-Dimensional Data Classification." Asian Conference on Intelligent Information and Database Systems 2017, Springer, Cham.
- [12] D. Li, et al, "Application of a Mobile Chronic Disease Health-Care System for Hypertension Based on Big Data Platforms." *Journal of Sensors*, 2018.
- [13] Welcome to ApacheTM Hadoop, Available: http://hadoop.apache.org/
- [14] J. Han, J. Pei, and M. Kamber, *Data mining: concepts and techniques*, Elsevier, 2011.
- [15] S. Kaushik, (2016), Introduction to Feature Selection methods with an example, Available: https://www.analyticsvidhya.com/blog/2016/12/introduction-to-featureselection-methods-with-an-example-or-how-to-select-the-right-variables
- [16] M. C. Thomas, A. T. Joy, *Elements of information theory*, John Wiley &Sons, 2012.
- [17] I. Kononenko, "Estimating attributes: analysis and extensions of RELIEF." *European conference on machine learning*, Springer, Berlin, Heidelberg, 1994.
- [18] S. Paul, M. Malik, and D. Petros, "Feature selection for linear SVM with provable guarantees", *Artificial Intelligence and Statistics*, 2015, p. 735-743.
- [19] KNHANES, (Nov, 2017), Available: http://knhanes.cdc.go.kr.
- [20] F. S. Samaria, and C. H. Andy, "Parameterisation of a stochastic model for human face identification." *Applications of Computer Vision*, 1994, p. 138-142.

- [21] D. Singh, et al, "Gene expression correlates of clinical prostate cancer behavior." *Cancer cell*, 1.2, 2002, p. 203-209.
- [22] C. X. Ling, H. Jin, and H. Zhang, "AUC: a better measure than accuracy in comparing learning algorithms", *Conference of the canadian society for computational studies of intelligence*, Springer, Berlin, Heidelberg, 2003, p. 329-341.
- [23] J. Akosa, "Predictive Accuracy: A Misleading Performance Measure for Highly Imbalanced Data." *Proceedings of the SAS Global Forum*, 2017.



Dingkun Li received the Ph.D degree in computer science from Chungbuk National University, Korea. He received the B.S. degree in software engineering from Harbin Institute of Technology (HIT), Heilongjiang, China, in 2009, and another M.S. degree in Computer Science from Chungbuk National University

(CBNU), South Korea, in 2014. Mr. Li's recent awards include the Best Poster Award at FITAT'14, Thailand; the Best Paper Award at ACIT'15, Japan; and the Best Paper Award at FITAT'16, China. His research interests are focused in the fields related to software engineering, machine learning, data mining, and big-data analysis.



Keun Ho Ryu received the Ph.D. degree in Computer Science and Engineering from Yonsei University, Korea, in 1988. He is also an honorary doctorate of the National University of Mongolia. He is currently a Professor at Chungbuk National University (CBNU), South Korea, and has been a leader in the Database and Bioinformatics

Laboratory, South Korea, since 1986. Also, he has served the Korean Army as an ROTC. He has worked at the University of Arizona, U.S.A., as a postdoctoral Research Scientist, and also at the Electronics and Telecommunications Research Institute, South Korea, as a Senior Researcher. He was a Vice-President of the Personalized Tumor Engineering Research Center.

Professor Ryu has served on numerous program committees including as a demonstration co-chair of the VLDB, a panel and tutorial co-chair of the APWeb, and a general co-chair of the FITAT. He has published or presented more than 1,000 referred technical articles in various journals and international conferences and is the author of several books. His research interests include temporal databases, spatiotemporal databases, temporal GIS, stream data processing, knowledge-based information retrieval, data mining, biomedicine, and bioinformatics. Professor Ryu has been a member of the IEEE and the ACM since 1983.



Erdenebileg Batbaatar received the B.S. degree in Software Engineering from the National University of Mongolia in 2013. Currently, he is pursuing the M.S. degree in computer science at Chungbuk National University (CBNU), South Korea. His research interests include biomedical and

chemical text mining and deep learning.



Hyun-Woo Park received the B.S. degree in Industrial Systems Engineering from Kongju University, Korea, in 2011. Currently, he is pursuing the Ph.D. degree in computer science at Chungbuk National University (CBNU), South Korea. His research interests include big data, biomedical, and data mining.



Seon-phil Jeong Currently, he is an associate professor at BNU-HKBU United International College in Zhuhai, China. He received his Ph.D. degree in MIS from Chungbuk National University, Korea. He has many experience in developing and implementing e-business systems for SMEs.

His research interests are Project Management, Recommend Systems, Business Intelligence and e-Business Applications.



Zhou Ye is President of the Karamay Central Hospital in Xinjiang, Chief Physician, Professor of Xinjiang Medical University, Graduatse Tutor of Shihezi University Medical School, the First Member of Standing Committee of China Medical Information Technology Specialized Committee, Expert Leader of

Karamay City Information Subject. Mr. Ye's recent awards include Golden Tripod Award (Prominent Personalities) powered by China Medical Health Information Technology in 2013, "Person of innovation of the year powered by China Health Service" in 2014, the Award of Outstanding Contribution to Science and Technology in Karamay city in 2016, the Award of Outstanding Management Talents powered by Xinjiang Health and Family Planning Commission in 2016.

RNNs for Lithuanian Multiword Expressions Identification

I. Bumbulienė, J. Mandravickaitė, A. Bielinskienė, L. Boizou, J. Kovalevskaitė, E. Rimkutė, L. Vilkaitė-Lozdienė, K. L. Man, T. Krilavičius

Abstract—We discuss an experiment on automatic identification of multiword expressions (MWEs) in Lithuanian corpus. Our training dataset was annotated morphologically (POS tagger). It was manually annotated with MWEs by 4 linguists as well. We also used word embeddings in our feature set. Deep learning methods are widely used in many NLP tasks and applications including MWEs identification. Thus, our experimental setup included deep learning methods (Recurrent Neural Networks; RNNs) and was used for automatic identification of contiguous and non-contiguous MWEs of different length. Best results (44.9% F1-Score) were achieved with RNNs and Stochastic Gradient Descent as optimizer together with Categorical Cross Entropy as loss function.

Index Terms—contiguous multiword expressions, identification, Lithuanian, non-contiguous multiword expressions, recurrent neural networks.

I. INTRODUCTION

A MULTIword Expression (MWE) is a sequence of at least two words frequently used together [1]. MWE, among other features, *acts as a single unit at some level of linguistic analysis* [2]. It also has *idiosyncratic interpretations that cross word boundaries (or spaces)* [3]. Accurate identification and processing of MWEs is one of the main problems for the

This research was funded by the Research Council of Lithuania (No. LIP-027/2016)), see <u>www.mwe.lt</u> for more details.

I. Bumbulienė is with Baltic Institute of Advanced Technology, Vilnius, Lithuania (corresponding author e-mail: ieva.bumbuliene@bpti.lt).

J. Mandravickaitė is with Baltic Institute of Advanced Technology (email: justina@bpti.lt).

A. Bielinskiene is with Vytautas Magnus University, Kaunas, Lithuania (email: agne.bielinskiene@vdu.lt).

L. Boizou is with Vytautas Magnus University, Kaunas, Lithuania (email: loic.boizou@vdu.lt).

J. Kovalevskaitė is with Vytautas Magnus University, Kaunas, Lithuania (email: jolanta.kovalevskaite@vdu.lt).

E. Rimkutė is with Vytautas Magnus University, Kaunas, Lithuania (email: erika.rimkute@vdu.lt).

L. Vilkaitė-Lozdienė is with Baltic Institute of Advanced Technology, Vilnius, Lithuania (email: laura.vilkaite@bpti.lt).

K. L. Man is Xi'an Jiaotong-Liverpool University, China and Swinburne University of Technology Sarawak, Malaysia (e-mail: ka.man@xjtlu.edu.cn).

T. Krilavičius is with Baltic institute of Advanced Technology, Vilnius, Lithuania and Vytautas Magnus University, Kaunas, Lithuania (email: t.krilavicius@bpti.lt) development of large-scale, effective and precise NLP technologies with applications such as foreign language acquisition, machine translation, text analytics and retrieval.

The difficulty in terms of identification and processing MWEs comes due to their highly heterogeneous behavior at different levels - lexical, syntactic and semantic. MWEs cover variety of linguistic constructions [3]: idioms (*kick the bucket, kill some time*), fixed phrases (*by and large, status quo*), noun compounds (*traffic light, fish tank*), compound verbs (*draw a conclusion, make a decision*), etc. Although MWEs are used easily by native speakers [4], they are a significant challenge for computational systems because of their flexibility and heterogeneity [5]. So, MWEs are one of the key issues for natural language parsing and generation, as well as real-life applications such as machine translation or text retrieval.

Nowadays, deep learning (DL) methods are widely used in many natural language processing (NLP) tasks and applications including identification of MWEs. Traditional methods that use lexical association measures and/or classification techniques (e.g. [15-20]) use n-grams and are limited to n size, e.g. limited to bigram MWEs [15, 16, 18], thus deep learning methods show to be more versatile. Usage of Recurrent Neural Networks (RNN) enables to detect MWEs in various formats. RNNs are very popular type of DL methods used for NLP as it enables to use and parse sequential information and its semantics, for example, to identify a difference of meaning between dog and hot dog [6]. Experiments reveal that bidirectional RNNs perform better than singular direction for NLP tasks [7]. As RNNs seem to be relevant for MWEs identification tasks, it was decided to use a bi-LSTM type of RNNs for Lithuanian MWEs identification.

II. CORPORA AND EMBEDDINGS

Main 70 million words corpus was collected from Lithuanian news portal **delfi.lt**¹. Articles were collected from 11 different categories: people, projects, science, auto, sport, life, news, citizen, business, fit, other category (did not fit in the mentioned categories or category was not identified). Small corpus (72 thousand words) was constructed out of the main corpus by randomly selecting 1% of the articles for manual MWEs tagging. In Fig. 1 a comparison of full corpus and annotated part is shown. Visualization (Fig. 1) shows that a distribution of



Fig. 1. Small corpora and all corpora comparison.

sentences and words in small and full corpus is very similar. Also, the distribution of articles in different categories is similar the whole and annotated corpora. The only exception is the category *other*. The higher priority was given to clearly defined categories and this resulted in smaller number of articles in *other* category in the annotated corpus.

Manual MWEs tagging was performed by 4 annotators -2 pairs of annotators tagged the same articles independently. Annotation of MWEs was performed using Brat Rapid Annotation Tool². As MWEs are highly heterogeneous, 2 different categories of MWEs as separate classes were chosen to be marked – idioms and collocations. Idioms are defined by their integrity of meaning or semantic indivisibility while collocations – by their consistency of use and semantic indivisibility [11]. Altogether 4335 MWEs there were tagged: 283 idioms and 4052 collocations.

Total agreement of both pairs of annotators was \sim 37% for both classes of MWEs. The main reasons of disagreements were the following:

- 1. annotator did not tag MWE,
- 2. annotators tagged MWEs as different types,
- 3. annotators tagged different boundaries of MWEs,
- criteria of identifying MWEs were understood differently by the annotators.

For training deep artificial neural networks (ANNs), it was decided to merge 2 classes of MWEs into one as each class – idioms and collocations – consisted of small number of annotated MWEs. Overall, 4695 sentences were annotated with MWEs and formatted into BIO scheme [14] where insertion of O between B-MWE and I-MWE is allowed due to the MWEs with interchangeable middle word(s). B-MWE points to the

first word of MWE, I-MWE – to the other words that belong to that same MWE and O marks word(s) that is not a part of MWE. For example, in snippet of data shown in Fig. 2 can be found 3 different MWEs:

- 1. consisting of 2 words (en: *serial killers*): *serijiniai* ADJ B-MWE *žudikai* NOUN I-MWE
- 2. consisting of 3 words (en: *without any doubt*): *be* ADP B-MWE *jokios* PRON I-MWE

abejonės NOUN I-MWE

3. consisting of 2 words with a gap of 1 word (en: *obsessed* with [insert: the same] mania):

apsėsti VERB B-MWE vienos PRON O manijos NOUN I-MWE

Overall, the training data contained MWEs of different length (from 2 to 6 words) and with inserts of various length (from 0 to 5 words).

Our main corpus was morphologically annotated to include part-of-speech (POS) tags for training RNNs. Some part of speech tags also can be seen in snippet of training data in Fig. 2.

The morphological annotator Semantika.lt we used is based on Hunspell open source platform (it has a lexicon and the affixes) which was supplemented with the statistical method (Markov model + Viterbi algorithm) for the morphological Garsiausi ADJ O pasaulio NOUN O serijiniai ADJ B-MWE žudikai NOUN I-MWE be ADP B-MWE jokios PRON I-MWE abejonės NOUN I-MWE buvo VERB O apsėsti VERB B-MWE vienos PRON O manijos NOUN I-MWE - PUNCT O Fig. 2. Snippet of training data.

disambiguation. This new analyzer has 429 groups of rules; 1,518 explicit tags for flexing/non-flexing characteristics; 5,832 rules for suffix and affix alternation in 16,734 alternation cases [22].

For the reasons of simplicity, for our experiments original output of Semantika.lt morphological analyzer was converted into the format with Universal POS tags only (see more at the website of Universal Dependencies Project³).

POS tags for training was used in One Hot Encoding format, where data is encoded in binary manner indicating presence or absence of each possible value from the original data [12]. For example, if 1st, 2nd and 3rd columns represent NOUN, VERB and PUNCT and the analyzed word is a noun, encoding is 1 0 0. However, for encoding words it was decided to use word vectors instead of One Hot Encoding. Thus, word embeddings were used for quantitative representation of text for deep ANNs. It was decided to build GloVe word embeddings [8]. Words vectors were generated using the following parameters:

- minimum frequency of word: 5 (vocabulary contains 331 000 words; all punctuation marks are not deleted);
- 2. window size: 5;
- dimension of 200 and 100 iterations. GloVe word vectors were built using model implementation scripts developed by Stanford NLP group [9].

III. EXPERIMENTAL EVALUATION

Sequence labeling model proposed in [10] was applied. It uses bidirectional LSTM cells, pretrained GloVe word embeddings, character embeddings and linear Conditional Random Field method in the output layer. Data (in BIO format) was separated into three sets:

- 1. training (2917 sentences),
- 2. validation (879)
- 3. testing (899).

Parameters of used model:

- 1. *batch size* 20;
- 2. Adam optimization;
- 3. *learning rate* 0.001.

Hyper parameters:

- 1. hidden LSTM size 100 on character embeddings;
- 2. hidden LSTM size 300 on word embeddings.

The results on test corpus reached 26% F1-Score.

Secondly, experiments were performed using bi-LSTM with part-of-speech tags in One Hot Encoding format for training, excluding character embeddings and with different optimization and loss functions. The same data as with the first model was used. Setup consisted of:

- (1) *batch size* 20;
- (2) *dropout* 0.1;

(3) hidden LSTM size 100 on word embeddings;

(4) *optimizers*: RMSprop, Stochastic Gradient Descent, Adagrad, Adadelta, Adam, Adamax, Nadam; and

(5) *loss functions*: Categorical Cross Entropy, Categorical Hinge.

F1-Score values varied from 37.8% to 44.9% depending on the optimizers and loss functions (see distribution of those values in Fig. 2). The best result (44.9% F1-Score) was achieved with Stochastic Gradient Descent as optimizer and Categorical Cross Entropy as loss function.



Fig. 3. F1-Score values distribution.

IV. MANUAL EVALUATION

After experimental evaluation by F1-Score, a manual evaluation was done by the same annotators who manually tagged our small corpus with MWEs. It was decided to analyze training data once again in order to improve its quality, hoping that this may lead to even better results. Thus, all data (that was manually annotated at the first step) was annotated by the trained model which has shown the best result. Then False Positives were collected for analysis, except the ones that contained only 1 word (the latter ones were discarded as being obviously incorrect). The number of False Positives was 1894. This list of False Positives was analyzed by the annotators and 445 of candidates were marked as incorrect whereas the other 1449 were recognized as MWEs.

Most of the discarded False Positives fell into 3 categories:

- incomplete MWEs that were actually a part of longer MWEs (e.g. *atitikti laikmečio* [missing: *dvasiq*] → en: comply with [missing: the spirit] of the times);
- co-occurrences (automatically recognized statistically significant word pairs or longer sequences [21]; e.g. *investuoja i* → en: investing in);

³ More about the Universal Dependencies Project: <u>http://universaldependencies.org/</u>

• free compounds (not MWEs but regular word sequences; e.g. *ivestas euras* → en: introduction of the euro).

Therefore, 76% of False Positives seemed to be relevant MWEs. These findings are very promising. It signifies that training data can be extended.

V. CONCLUSIONS

Values of F1-Score is decent (e.g. in shared task for automatic identification of verbal MWEs F1-Score for Lithuanian dataset was 28% [13]) and used methods seem to be very promising. Usage of RNNs enables to detect MWEs of various lengths. Traditional methods using lexical association measures and classification techniques are dependent on n-grams and are limited to n size, thus deep learning methods show to be more versatile. In the future we plan to test different types of deep neural networks and their architectures for identifying Lithuanian multiword expressions: change LSTM to Gated Recurrent Unit (GRU), apply Convolutional Neural Network (CNN).

Also, we plan to combine ML models to a system where different models may identify different types of MWEs and consequently may lead to even better results. Moreover, additional linguistic rule-based filtering may be added to ensure that found MWEs are grammatically correct.

REFERENCES

- R. Marcinkevičienė, "Tradicinė frazeologija ir kiti stabilūs žodžių junginiai", *Lituanistica*, 4(48), 2001, pp. 81--98.
- [2] N. Calzolari, Ch. J. Fillmore, R. Grishman, N. Ide, A. Lenci, C. MacLeod, and A. Zampolli, "Towards best practice for multiword expressions in computational lexicons", In LREC, 2002.
- [3] I. A. Sag, T. Baldwin, F. Bond, A. Copestake, D. Flickinger, "Multiword expressions: A pain in the neck for NLP", *Computational Linguistics and Intelligent Text Processing*, Springer Berlin Heidelberg, 2002, pp. 1--15.
- [4] F. Boers, J. Eyckmans, J. Kappel, H. Stengers, and M. Demecheleer, "Formulaic sequences and perceived oral proficiency: Putting a lexical approach to the test", *Language teaching research*, 10(3), 2006, 245-261.
- [5] O. C. Acosta, A. Villavicencio, and V. P. Moreira, "Identification and treatment of multiword expressions applied to information retrieval", In *Proc. of the Workshop on Multiword Expressions: from Parsing and Generation to the Real World*, ACL, June 2011 June, pp. 101-109.
- [6] T. Young, D. Hazarika, S. Poria, and E. Cambria, "Recent Trends in Deep Learning Based Natural Language Processing", *CoRR*, abs/1708.02709, 2017.
- [7] W. Yin, K. Kann, M. Yu, and H. Schütze, "Comparative Study of CNN and RNN for Natural Language Processing", *CoRR*, abs/1702.01923, 2017.
- [8] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global Vectors for Word Representation", *EMNLP*, 2014.
- [9] "Github GloVe", https://github.com/stanfordnlp/GloVe, accessed: 2018-05-11.
- [10] G. Genthial, "Sequence Tagging with Tensorflow" https://guillaumegenthial.github.io/sequence-tagging-with-tensorflow.ht ml (2017 April), accessed: 2018-05-11.
- [11] E. Rimkutė, A. Bielinskienė, J. Kovalevskaitė, and L. Vilkaitė, "Towards The Criteria For Identification Of Idioms And Collocations", *Studies About Languages*, (31), 2017, 83-101.
- [12] J. E. Beck and B. P. Woolf, "High-level student modeling with machine learning", In *International Conference on Intelligent Tutoring Systems*, Springer, Berlin, Heidelberg, 2000, pp. 584-593.
- [13] A. Savary, C. Ramisch, S. Cordeiro, F. Sangati, V. Vincze, B. QasemiZadeh, M. Candito, F. Cap, V. Giouli, and I. Stoyanova, "The PARSEME shared task on automatic identification of verbal multiword expressions", *Proceedings of the 13th Workshop on Multiword*

Expressions (MWE 2017), Association for Computational Linguistics, Valencia (Spain), 2017, 31–47.

- [14] L. A. Ramshaw and M. P. Marcus, "Text chunking using transformation-based learning", In *Natural language processing using* very large corpora, Springer, Dordrecht, 1999, pp. 157-176.
- [15] J. Mandravickaitė, E. Rimkutė, T. Krilavičius, "Hybrid Approach for Automatic Identification of Multi-Word Expressions in Lithuanian", *Proceedings of the Seventh International Conference Baltic HLT 2016*, Amsterdam, Berlin, Tokyo, Washington, DC: IOS Press, 2016, 153–159.
- [16] J. Mandravickaitė, T. Krilavičius, "Identification of Multiword Expressions for Latvian and Lithuanian: Hybrid Approach", *Proceedings* of the 13th Workshop on Multiword Expressions (MWE 2017), Association for Computational Linguistics, Valencia (Spain), 2017, 97– 101.
- [17] L. Zilio, L. Svoboda, L. H. L. Rossi, and R. M. Feitosa, "Automatic extraction and evaluation of MWE", In *Proceedings of the 8th Brazilian Symposium in Information and Human Language Technology*, 2011, 214-218.
- [18] P. Pecina, "Lexical association measures and collocation extraction", Lang. Res. & Eval. Special Issue on Multiword expression: hard going or plain sailing, 44(1-2), 2010, 137–158.
- [19] C. Ramisch, A. Villavicencio, and C. Boitet, "Multiword expressions in the wild?: the mwetoolkit comes in handy", In *Proceedings of the 23rd International Conference on Computational Linguistics:* Demonstrations, Association for Computational Linguistics, August 2010, pp. 57-60.
- [20] C. Hashimoto and D. Kawahara, "Construction of an idiom corpus and its application to idiom identification based on WSD incorporating idiom-specific features", In *Proc. of EMNLP*, Honolulu, Hawaii, USA, 2008, pp. 992–1001.
- [21] S. Granger, M. Paquot, "Disentangling the phraseological web", *Phraseology. An Interdisciplinary Perspective* (red. S. Granger, F. Moinier), Amsterdam: John Benjamins, 2008, 27–49.
- [22] J. Kapočiūtė-Dzikienė, E. Rimkutė, and L. Boizou, "A Comparison of Lithuanian Morphological Analyzers", In *International Conference on Text, Speech, and Dialogue*, Springer, Cham, August 2017, pp. 47-56.

Long Short-Term Memory Encoder-Decoder for Traffic Flow Prediction

Qi Chen, Wei Wang, and Xin Huang

Abstract—Traffic flow prediction is a fundamental problem in transportation modeling and management. Neural networks have been widely applied to traffic flow prediction in the past few years. However, existing studies only focus on predicting the traffic flow at next time step, while travelers may need a sequence of predictions to make better travel decisions. To address the above limitation, this paper introduces a long short-term memory encoder-decoder architecture for traffic flow prediction. Experiments demonstrate that the proposed method for traffic flow prediction has superior performance and can be further used for traffic anomaly detection.

Index Terms—Deep learning, recurrent neural network encoder-decoder, long short-term memory, traffic flow prediction.

I. INTRODUCTION

Traffic flow information is needed for individual travelers, business sectors, and government agencies to make better travel decision, alleviate traffic congestion, and improve traffic operation efficiency [1]. With the rapid development and deployment of intelligence transportation systems (ITSs), traffic flow prediction has gained increasing attention in recent years. However, the traffic flow prediction is challenging as it is a non-linear and stochastic problem that mixes with many uncertain factors such as traffic incident, traffic event, weather condition, etc.

Deep learning has drawn a lot of academic and industrial interests in recent years. It has been continuously improving the state-of-the-art in speech recognition [2], computer vision [3], natural language processing [4] and many other domains. It allows computational models that are composed of multiple layers to learn features of data by using the backpropagation algorithm [5]. As the traffic prediction problem is complicated in nature, deep learning algorithms have been applied to capture traffic patterns and provide traffic predictions even without prior knowledge about the external factors.

Current studies [6-9] have shown the promising results of applying deep learning in traffic flow prediction. However, these studies only focus on predicting traffic flow at next time

Qi Chen, Wei Wang and Xin Huang are with the Department of Computer Science and Software Engineering, Xi'an Jiaotong Liverpool University, China. Email: qi.chen,wei.wang03,xin.huang@xjtlu.edu.cn step (e.g. after 15 mins). As travelers may need a sequence of traffic flow predictions (e.g. traffic flow in the next few hours) to make better travel decision, we propose a Long Short-Term Memory Encoder-Decoder (LSTM-ED) architecture. To further improve the accuracy of prediction results, time information such as time of the day, day of the week and national holiday is added to the model input. The proposed model shows the superior performance on traffic flow prediction and can be further used for real-time traffic anomaly detection.

The rest of the paper is organised as follows. Section II reviews the existing studies on traffic flow prediction. Section III presents the LSTM-ED architecture for traffic flow prediction and provides details of the design. Section IV shows the experimental and evaluation results. The conclusion and future work are described in Section V.

II. RELATED WORK

Over the past few decades, a number of traffic flow prediction models have been developed. As early as 1970s, the autoregressive integrated moving average (ARIMA) model was applied to short-term traffic flow prediction [10]. Due to the stochastic and nonlinear nature of traffic flow, nonparametric approaches, e.g. k-NN [11], SVR [12] and ANN [13], attracted more attention from researchers for traffic prediction problems. As the recent advances in deep learning, researchers have shown the great potential of deep neural networks in traffic prediction. Huang et al. [6] firstly used a deep belief network to learn features of traffic flow and proposed a multi-task learning architecture for traffic flow prediction. Similarly, Lv et al. [7] proposed a stacked autoencoder model for traffic flow prediction. Ma et al. [8] applied Long Short-Term Memory to capture long-term memory of traffic speed for short-term traffic speed prediction. Tian and Pan [9] applied Long Short-Term Memory to memorize long historical traffic flow data for short-term traffic flow prediction. Since traffic flow prediction and machine translation share some similar characteristics, it is natural to consider the state-of-the-art models in machine translation, e.g. recurrent neural network encoder-decoder [14], [15] for traffic flow prediction. In this paper, we explore the use of LSTM-ED model for traffic flow prediction.

III. LSTM Encoder-Decoder Architecture

- A. Long Short-Term Memory (LSTM)
- A recurrent neural network (RNN) is a type of artificial

Manuscript received June 13, 2018. The research is funded by the Research Development Fund at Xi'an Jiaotong-Liverpool University, contract number RDF-16-01-34. We gratefully acknowledge the grant subsidy (XJTLU RIBDA2018-IRP1) provided by the Research Institute of Big Data Analytics (RIBDA), Xi'an Jiaotong-Liverpool University.

neural network used to process sequences of inputs. As traditional RNN suffers from the exploding and vanishing gradient problem during gradient descent training process, Long Short-Term Memory (LSTM) [16] was proposed to model long-term dependencies. The architecture of a standard LSTM unit is visualised in Fig. 1.



Fig. 1. Long short-term memory architecture

The LSTM architecture is composed of one input layer, one recurrent hidden layer and one output layer. A LSTM unit contains a cell *c*, an input gate *I*, an output gate *O* and a forget gate *F*. The cell is responsible for remembering values over time steps. The three gates allow LSTM memory cell to store and access information over long periods of time, thereby mitigating the vanishing gradient problem. Each of the three gates can be considered as a neuron that compute an activation using an activation function σ . Given a sequence of input $x = (x_1, x_2, ..., x_T)$, a standard LSTM computes a sequence of outputs $y = (y_1, y_2, ..., y_T)$ by iterating the following equations from (1) to (6):

$$F_t = \sigma(W_F \cdot [h_{t-1}, x_t] + b_F) \tag{1}$$

$$I_t = \sigma(W_I \cdot [h_{t-1}, x_t] + b_I) \tag{2}$$

$$c_t = F_t * c_{t-1} + I_t * \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$
(3)

$$O_{t} = \sigma(W_{0} \cdot [h_{t-1}, x_{t}] + b_{0})$$
(4)

$$h_t = O_t * \tanh(c_t) \tag{5}$$

$$y_t = W_f \cdot h_t + b_y \tag{6}$$

Where F_t represents the output of a forget gate, I_t denotes the output of an input gate, and O_t is the output of an output gate. The cell state and hidden state is respectively denoted as c_t and h_t . The weight matrices W, bias vectors b and sigmoid function σ are utilised to build connections between input layer, hidden layer and output layer.

B. LSTM Encoder-Decoder (LSTM-ED)

Since it is not clear how to apply an RNN to model problems whose input and output sequence might have different lengths, a RNN encoder-decoder architecture [14] (shown in Fig. 2) is proposed. The RNN Encoder-Decoder learns to encode a sequence into a context vector representation and to decode the context vector representation back into a sequence. Note that the input sequence length T and output sequence length T' may differ.



Fig. 2. Long short-term memory encoder-decoder architecture

To capture the long-term dependencies, LSTM units are used in the RNN encoder-decoder architecture. In the LSTM-ED model, the encoder is an LSTM that reads an input sequence $x = (x_1, x_2, ..., x_T)$ and generates a context vector C of the whole input sequence, where C is the final hidden state h_T of the encoder. Then, another LSTM is used as decoder to generate the output sequence $y = (y_1, y_2, ..., y_T')$. Output y_t at time step t is used as the input at the next time step t+1. Unlike the LSTM described in Section III.A, the decoder LSTM also includes the context vector C as input while updating the hidden state. Hence, the hidden state of the decoder at time t is calculated using Equation (7) below.

$$\hat{h}_{t} = g(\hat{h}_{t-1}, y_{t-1}, C)$$
 (7)

where g represents a LSTM unit, y_{t-1} denotes the output at time step *t*-1, and *C* is the context vector.

C. LSTM-ED for traffic flow prediction

The traffic flow of vehicle detector stations at the *t*th time step is denoted by f_t . At time T, the task is to predict traffic flow sequence $P = (f_{T+1}, f_{T+2}, ..., f_{T+T'})$ based on the historical traffic flow sequence $H = (f_1, f_2, ..., f_T)$. For each time step in the input and output sequence, f_t is a *d*-dimensional vector $\{f_t^1, f_t^2, ..., f_t^d\}$ that represents the traffic flow of *d* vehicle detector stations at time *t*. The input of the LSTM-ED encoder is the historical traffic flow sequence *H*, and the output of the LSTM-ED decoder is the predicted traffic flow sequence *P*. Since traffic flow may be affected by many other factors (e.g., time and weather), a simple solution is to append the information at each time step of input sequence *P*, in which case the input and output sequence dimensions *d'* and *d* are different.

IV. EXPERIMENTS & EVALUATION

A. Dataset

The Caltrans Performance Measurement System (PeMS) [17] is a widely used dataset for traffic flow prediction. The traffic data are collected every 30 seconds from various types of vehicle detector stations throughout the state of California. Then, the data is aggregated at 5-min interval for each detector station. We further aggregate the data into 15-min interval, as suggested by the Highway Capacity Manual [18]. In this paper, the traffic flow data is collected from 243 vehicle detector

stations in district 5 (including Monterey, San Benito, etc.) from May 1, 2017 to Feb 28, 2018. The data of the first nine months is used as the training set, and the data of the remaining one month is used as the testing set.

B. Experiments

Input and Output: The number of input time steps and output time steps may vary for different traffic prediction tasks. In this paper, the number of input time steps *T* is set to 96 and the number of output time steps *T'* is set to 8, which means the last 24 hours' traffic flow data is used to predict the next 2 hours' traffic flow data. Since the traffic flow is affected by time information, we further append three rows indicates the time of the day (from 0 to 95), the day of the week (from 0 to 6) and whether it is national holiday (0 or 1). Thus, we obtain the input of shape $R^{(243+3)\times96}$ and output of shape $R^{243\times8}$, where 243 is the number of vehicle detector stations. As traffic flow volume of different stations coupled with time information may have different scales, the input data is further normalised in the range of [0,1] by (8).

$$X' = \frac{X - X_{min}}{X_{max} - X_{min}} \tag{8}$$

Model parameters: With regard to the LSTM-ED network, we need to determine the number of hidden layers, the number of hidden units, the dropout rate, and etc. Because training LSTM is time consuming, the number of hidden layers is set to 1 in this study. We choose the number of hidden units from {32, 64, 128, 256} and the number of batch size from {64, 128, 256, 512}. After performing grid search runs, the best model parameters are shown in table I. LSTM-ED experiments were run using Tensorflow 1.3, python 3.6, and Windows 10 on a laptop with a i7-6700HQ CPU, 8GB RAM and GTX-970M GPU.

TABL	EI. MOI	DEL PARAMETER SETTIN	GS		
LSTM-ED model parameters					
Input length	96	Batch size	256		
Output length	8	Dropout	0.5		
Hidden layers	1	Epochs	1000		
Hidden units	128	Optimizer	Adam		

C. Evaluation

Fig. 3 presents the output of the proposed approach for the traffic flow prediction of three different vehicle detector stations at Freeway 1, Freeway 101 and Freeway 156. The actual traffic flow data is also included for comparison. The figure shows that the LSTM-ED model is able to learn the traffic flow patterns and provides reasonable traffic predictions. Note that every 2 hours' (8 timesteps) traffic flow prediction in Fig. 3 is calculated by the last 24 hours' (96 timesteps) traffic flow data.

To evaluate the effectiveness of the proposed model, we use the root mean square error (RMSE), which is calculated by (9).

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (|f_i - \hat{f}_i|)^2}$$
(9)





Fig. 3. Traffic flow prediction results of three different stations. (a) station at Freeway 1. (b) station at Freeway 101. (c) station at Freeway 156.

The performance of the proposed LSTM-ED method with the time information (LSTM-ED+T) is compared with the AVG method and LSTM-ED method. The AVG method is a simple baseline that calculates the average traffic flow of each station at specific time of week (e.g. 9 AM on Monday). As shown in Table II, as expected, simply using the AVG method leads to inaccurate prediction results. In contrast, the LSTM-ED method can provide much better prediction results. Since the time information (time of the day, day of the week, national holiday) may affect traffic flow patterns, the proposed LSTM-ED + T model further decrease the RMSE by 0.001.

 TABLE II.
 PERFORMANCE COMAPRISION OF DIFFERENT MODELS

RMSE
0.0686
0.0447
0.0437

D. Traffic Anomaly detection

As traffic incidents and events can significantly affect traffic flow, one interesting application is to use traffic anomalies to detect potential traffic incidents or events. In this section, a simple use case that applies the proposed method together with social media information to detect traffic incidents and traffic events is illustrated.



Since the proposed LSTM-ED method can capture traffic flow patterns, a high prediction error is likely to indicate a traffic anomaly. Fig. 4 shows the prediction error of a station at freeway 101 from 2018-02-05 to 2018-02-12. The prediction error is calculated from Fig. 3(b) by $e = |f_i - \hat{f_i}|$, where f_i is the observed traffic flow, and $\hat{f_i}$ is the predicted traffic flow. As shown in Fig. 3(b) and Fig, 4, the sudden drop of traffic flow at around 8 AM on 2018-02-09 yields high prediction error, which indicates a potential traffic incident or event. Because social media data can be used to complement and corroborate data from physical world, we further use the temporal and spatial information collected from the detected traffic anomaly to filter traffic related tweets posted by trusted organisations on Twitter. As shown in Fig. 5, a traffic accident caused by an overturned semi-truck is reported by CaltransD5 at the same time and the

same place.



Caltrans District 5 @CaltransD5 · Feb 9 SigAlert: left fast lane of northbound #Hwy101 is CLOSED near #Hwy58 due to overturned semi--truck so use caution & expect traffic delays this morning from SLO to north county.



Caltrans District 5 @CaltransD5 · Feb 9 SigAlert Update: both NB #Hwy101 lanes near #Hwy58 are CLOSED until 9am to remove truck.



Caltrans District 5 @CaltransD5 · Feb 9 SigAlert Update: both NB #Hwy101 lanes near #Hwy58 are now OPEN but expect residual traffic delays for awhile.

Fig. 5. traffic anomaly related tweets

This is a simple traffic anomaly detection example based on historical traffic data and tweets. Real-time traffic anomaly detection system can be implemented by collecting streaming tweets and comparing prediction results with the current traffic flow.

V. CONCLUSION AND FUTURE WORK

We proposed a deep neural network model, LSTM-ED, to generate a sequence of predictions. To improve prediction accuracy, time information is also added to the model input. The experiments show that the proposed model can capture traffic patterns and provide superior prediction results. A simple use case of the proposed method for traffic anomaly detection is also illustrated. For future work, more state-of-the-art deep learning models will be explored for time series prediction. Furthermore, data from other domains will be collected and used to further improve prediction performance.

REFERENCES

- N. Zhang, F. Y. Wang, F. Zhu, D. Zhao and S. Tang, "DynaCAS: Computational experiments and decision support for ITS," *IEEE Intelligent Systems*, vol. 23, no.6, 2008.
- [2] G. Hinton, L. Deng, D. Yu, et al, "Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp.82-97, 2012.
- [3] A. Krizhevsky, I. Sutskever, and G. Hinton, "Imagenet classication with deep convolutional neural networks," *International Conference on Neural Information Processing Systems*, pp. 1097-1105, 2012.
- [4] A. Vaswani, N. Shazeer, N. Parmar, et al, "Attention is all you need," Advances in Neural Information Processing Systems. pp. 6000-6010, 2017.
- [5] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no.7553, pp.436-444, 2015.
- [6] W. Huang, G. Song, H. Hong, and K. Xie, "Deep architecture for traffic flow prediction: deep belief networks with multitask learning," *IEEE Transactions on Intelligent Transportation Systems*, vol. 14, no. 5, pp.2191-2201, 2014.
- [7] Y. Lv, Y. Duan, W. Kang, Z. Li, and F. Y. Wang, "Traffic flow prediction with big data: a deep learning approach," *IEEE Transactions on Intelligent Transportation Systems*, vol. 16, no. 2, pp. 865-873, 2015.
- [8] X. Ma, Z. Tao, Y. Wang, H. Yu, and Y. Wang "Long short-term memory neural network for traffic speed prediction using remote microwave sensor data," *Transportation Research Part C: Emerging Technologies*, vol. 54, pp. 187-197, 2015.
- [9] Y. Tian, and L. Pan, "Predicting short-term traffic flow by long short-term memory recurrent neural network," *Smart City/SocialCom/SustainCom (SmartCity), 2015 IEEE International Conference on. IEEE*, pp. 153-156, 2015.
- [10] M. S. Ahmed and A. R. Cook, "Analysis of freeway traffic time-series data by using Box–Jenkins techniques," *Transportation Research Record*, no. 722, pp. 1–9, 1979.
- [11] H. Chang, Y. Lee, B. Yoon, and S. Baek, "Dynamic near-term traffic flow prediction: System oriented approach based on past experiences," *IET Intelligent Transportation Systems*, vol. 6, no. 3, pp. 292–305, 2012.
- [12] Y. S. Jeong, Y. J. Byon, M. M. Castro-Neto, and S. M. Easa, "Supervised weighting-online learning algorithm for short-term traffic flow prediction," *IEEE Transactions on Intelligent Transportation Systems*, vol. 14, no. 4, pp. 1700–1707, 2013.
- [13] K. Y. Chan, T. S. Dillon, J. Singh, and E. Chang, "Neural-network-based models for short-term traffic flow forecasting using a hybrid exponential smoothing and Levenberg–Marquardt algorithm," *IEEE Transactions on Intelligent Transportation Systems*, vol. 13, no. 2, pp. 644–654, Jun. 2012.
- [14] K. Cho, B. Van Merriënboer, C. Gulcehre, et al, "Learning phrase representations using RNN encoder-decoder for statistical machine translation," *Empirical Methods in Natural Language Processing* (*EMNLP*), pp.1724-1734, 2014.
- [15] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," *Advances in neural information processing systems*. pp. 3104-3112, 2014.
- [16] S. Hochreiter, and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp.1735-1780, 1997.
- [17] Caltrans, Performance Measurement System (PeMS), 2018. [Online]. Available: http://pems.dot.ca.gov.
- [18] "Highway Capacity Manual". Transportation Research Board, Washington, D.C. 2000.

Design of an Intelligent Temperature Control for the MIMO Thermal System

Yujia Zhai, Yuanye Fang, Zhejian Zhang, Sanghyuk Lee, Kejun Qian, Ka Lok Man

Abstract—An intelligent system designed for multi-input, multioutput (MIMO) temperature control of an exact substrate with the detailed control method is presented. Due to the problem of nonlinearity, large delay time and temperature inertia in the industry thermal control system, auto-tuning PID control with anti-windup technology are combined for more efficient and more flexible control effect. The controlled object was modeled by experimental data through MATLAB/Simulink for the aim of monitory and simulation with optimal tuning cost to study and verify the proposed control system. The PID algorithm with anti-windup has the merits of fewer tuning time, smaller overshoot, better precision and higher robustness approved by simulation, theoretical and experiment analysis.

Index Terms-MIMO; PID control; anti-windup; MATLAB

I. INTRODUCTION

TEMPERATURE control with high precision in industry is **L** an important part especially in some special areas such as oil refining, plastic product manufacturing and domestic air conditioner etc. PID control is the most conventional method applied to deal with the problems with characteristics of nonlinearity, time-varying and big lag [1]. However, the nonlinear system is hard to be controlled by the conventional PID to maintain the process at the desired environment efficiently. To deal with the multi-input, multi-output (MIMO) process control, a variety of other strategies have been imposed, such as Fuzzy [2], Artificial Intelligence [3], Fuzzy Self-tuning PID [4], Model Predictive Control (MPC) [5], etc. The self-tuning PID algorithm with anti-windup combined with MPC, which is a way that is easy to be implemented and enable the temperature control process maintained at the desired operating conditions efficiently and safely.

In the MIMO temperature system, the multi-variable control based on the experimental data of all the controlled variables in the multi-loop is considered. Meanwhile, rather than traditional approach that is on account of knowledge of process, experience and insight, this kind of control is modelbased which means it depends heavily on the accuracy of the process model and can be incorporated-directly with the control law like MPC [5]. Furthermore, this approach ensures the simulation in the computer which allow the test of alternative control strategies and the determination of the controller settings' preliminary values.

The dynamic model built in this paper was developed empirically from measurements combined with the theoretical basis. The multi-zone control is achieved by applying autotuning PID control, standard feedback and feed-forward technologies while multi-variable and constraint control uses MPC. Realtime optimization (RTO) is also considered to be implemented to make some constraints on the input variables. Moreover, the data reconciliation and parameter estimation technologies with the help of recent plant data are planned to be applied to update the process model. Process dynamics and control issues were attempted to be solved with planning and scheduling operations at the beginning. The computer simulation is developed with MATLAB/Simulink programming and the hardware during the research is based on B&R Industrial Automation Temperature Control Device.

The rest of this paper is organized as follows: a brief introduction and related works about the research background is presented in Section II. The model of the temperature control system and its mathematical representation are described in Section III. Detailed multi-zone control approaches are provided in Section IV. Section V focuses on results and analysis. Finally, conclusion remarks and future work proposal are made in Section VI.

II. BACKGROUNDS AND RELATED WORKS

A. Structure of B&R PID temperature control device

The object been controlled is the B&R PID temperature control device. There are three metal rods in the entire controlled object, which means that this is a three-zone temperature control system. One heater is installed in the bottom of the metal rod to raise the temperature in each metal rod. Meanwhile, three metal rods are arranged side by side with a metal baffle above them. Three temperature sensors are placed on the left side of each metal rod. In addition, a fan under the middle of the second

Yujia Zhai, Yuanye Fang, Zhejian Zhang, Sanghyuk Lee, are with the Department of Electrical and Electronic Engineering, Xi'an Jiaotong-Liverpool University, China.

Email: yujia.zhai@xjtlu.edu; yuanye.fang15@student.xjtlu.edu.cn; zhejian.zhang15@student.xjtlu.edu.cn; sanghyuk.lee@xjtlu.edu.cn

Kejun Qian is with Suzhou Power Grid Comparny, China. Email: <u>kejun.qian@xjtlu.edu.cn</u>

Ka Lok Man is with Xi'an Jiaotong-Liverpool University, China and Swinburne University of Technology Sarawak, Malaysia Email: ka.man@xjtlu.edu.cn

metal rod is used to cool the system.

The electrical part includes mounting rails, relays and 24V switching power supply. The digital output of the control system is connected to the solid-state relay. The solid-state relay output controls the voltage of the heater (the PWM uses the PWM mode to control the output), and the relay 220V powers (in the PWM 100% output when it is the biggest). The digital output signal of the control system supplies power directly to the cooling fan.

The heaters and the fan are designed to be controlled by equivalent PWM because of the feature of the large inertia hysteresis in temperature control. The input and output modules of the B&R control system used in the entire system are as follows: X20CP1584 is used as the CPU controller; the bus controllers use X20BC0083 and X20PS9400; X20DO4332 is applied to control the heater while X20DS1119 is to drive the fan. Finally, the temperature sensors are controlled by X20AT6402 and X20AT6402a.

B. B&R Programmable Computer Controller

The intelligent temperature control for the MIMO thermal system is designed based on Programmable Computer Controller (PCC) which is a B&R industrial automation product that combined the advantages of PLC and industrial computer. The system controlled by PCC is implemented pragmatically with its proprietary high-level language, Automation Studio which supports various programming languages specified in IEC61131-3 such as the standard C that used during this design. Additionally, as far as the use of MATLAB/Simulink during the design of the process control and simulations, the function of the Automation Studio that enables the translation between the process model and C code in PLC is utilized to maximize the efficiency of process design and control.

III. THE MODELING OF PROCESS DYNAMICS AND CONTROL

The empirical dynamic model directly from experimental data was applied because of the limitation of the rigorous theoretical model. To develop the empirical dynamic models, the activity is called the process or system identification [6]. The empirical dynamic model is simpler than the theoretical model and can also be solved in the actual problem while its response time is much shorter than the actual process response time. Furthermore, the discrete-time model was also implemented by digital computer through sampling the continuous time process variables at regular intervals. The firstorder model with general model-fitting techniques based on nonlinear regression [7] is used to calculate model parameters in this research. The first-order model was chosen not only because of the minimal measurement difference between model prediction s and data, but also the criteria to choose the simplest model structure to ensure the physical applicability. Overall, the empirical dynamic model is developed followed by the systematic model building procedure presented from Ljung (1999) [8].

Consider about the mathematical representation of the first

order model, the transfer function can be represented as:

$$\frac{Y(s)}{U(s)} = \frac{K}{\tau s + 1} \tag{1}$$

while the step response of it is:

$$y(t) = KM(1 - e^{-\frac{t}{\tau}})$$
(2)

then it can be found that the initial slope of the normalized step response after rearranging and evaluating the limit at t = 0 is:

$$\frac{d}{dt} \left(\frac{y}{KM} \right)_{t=0} = \frac{1}{\tau} \tag{3}$$

which shows that the tangent at t = 0 with the horizontal line, $\frac{y}{KM} = 1$, occurs at $t = \tau$ while the response y(t) reaches 63.2% of its final value at $t = \tau$. The graphical analysis can be used to do the process identification with these theoretical bases. However, departures from the ideal first-order model are common. A time-delay term is included to improve the agreement between model and experimental responses. The fitting of a first-order-plus-time-delay model (FOPTD) with transfer function of:

$$G(s) = \frac{Ke^{-\theta s}}{\tau s + 1}$$
(4)

The method of Sundaresan and Krishnaswamy (1978) [9] can do parameter estimation in this condition.

Rather than analyzing the step response data graphically, nonlinear regression and standard software such as Excel and MATLAB can also be used to obtain the first-order transfer function model. The optimization method used by Excel is due to the generalized reduced gradient technique (Edgar et al., 2001) and the default algorithm in MATLAB is the BFGS quasiNewton method (Ljung, 2007) [10]. The regression method is effective with sufficient response data measured from experiments and does not depend on graphical correlations when has fewer differences between predictions and data compared with the graphical analysis approach. Therefore, the system identification and the discretization of process model are chosen to be completed on MATLAB during the research. The final model considered about coupling in the system can be seen in Figure 1.



Fig. 1: Multi-temperature zone coupling modeling

IV. CONTROL APPROACHES

PID is important in the field of industrial control with its advantages of convenience, stability and reliability. In order to achieve the control requirements, a parallel form PID controller is adopted in this temperature control project. In a close loop system, the PID controller can sum the proportional, integrate, and derivative actions of a control error, which is the difference between a reference and the system output, as an input of the controlled plant to control the system. A discrete PID controller is implemented in B&R PLC and the compensator formula is defined by:

$$K_p + K_i \cdot T_s \frac{1}{z-1} + K_d \cdot \frac{1}{T_s} \frac{z-1}{z}$$
 (5)

This is a parallel controller with the Forward Euler Integrator method, where Ts is the sampling period and the sum weights K_p , K_i , and K_d are the proportional, integral, and derivative gain parameters, respectively.

The integral action is harnessed to eliminate the steadystate error and a relatively large sum weights K_i can expedite the elimination. Nevertheless, a large K_i may leads the integration wind-up if the controller output is saturated, which causes the system recovers slowly from the non-linear region so that the controller response is delayed. The great inertia and slow cooling of the temperature control system may aggravate the situation and reduce the system controllability. Therefore, an integrator clamping based anti-windup method is used to improve the saturated PID system performance by decreasing the excessive accumulation of the integral part.

The derivative part can restrain the error in advance to depress the overshoot. A large K_d may over advance the control response, prolong the adjustment time, and increase the sensitivity of the system to disturbance. For the temperature control system with large delay and inertia, suppress the overshoot can shorten the time of system reach the steady state efficiently.

To obtain a satisfactory result balancing the performance and robustness, an appropriate controller tuning is required. The one-quarter decay ratio tuning relations such as Ziegler-Nichols and Cohen-Coon methods are time-consuming for the temperature control project meanwhile the tuning relation based on integral error criteria results may miss some robust. Under the premise of high modeling accuracy, the PID tuner is applied to obtain an intuition and accuracy results with a fast speed. Approaching the physical response speed with a smaller overshoot is the criteria of the tuning.

V. RESULTS AND ANALYSIS

The demonstration of research design about both simulations and experiments then their results can be seen in the following. Firstly, the single PID control to zone 3 was applied while the sensor placed at the bottom of the zone was used to analyze the results. Meanwhile, the results of simulations through MATLAB will also be shown, where the values of parameters in PID control were found with some adjustments in actual conditions. The anti-windup technology [11] was also implemented to optimize PID control in this thermal system when overshoot occurs.



Fig. 2: The simulated result of single PID control in zone 3

Seen from the results in Figure 2 and Figure 3, the temperature aimed to be achieved was 80 °C and because the



Fig. 3: The actual result of single PID control in zone 3

accuracy of the sensor is 0.1 $^{\circ}$ C, the steady state value is 800. In Figure 3, the overshoot is limited to 1% and the peak is 81 $^{\circ}$ C while the fluctuation is less than 0.5 $^{\circ}$ C. The experimental result almost met the simulation result and the performance of overall control is acceptable.

Secondly, multiple PID control based on the coupling model was tested. The results are stated from Figure 4 to Figure 9.



Fig. 4: The simulated result of multiple PID control in zone 1



Fig. 5: The actual result of multiple PID control in zone 1



Fig. 6: The simulated result of multiple PID control in zone 2

The results of the multi-variable control were monitored by the commands that the temperature of zone 1 was respected to achieve 80 °C with the input at 0 second, the temperature of zone 2 was respected to achieve 120 $^{\circ}$ C with the input at 600 seconds and the temperature of zone 3 was respected to achieve 150 $^{\circ}$ C with the input at 1200 seconds. Due to the phenomena of coupling, the auto-tuning PID control cannot control the system effectively although can make the results get the targets in some situations like what have been described above.





Fig. 8: The simulated result of multiple PID control in zone 3



Fig. 9: The actual result of multiple PID control in zone 3

VI. CONCLUDING REMARKS

This research has introduced and studied the control of the MIMO thermal system based on B&R PID temperature control device and B&R PCC using auto-tuning PID control approach with the technology of anti-windup. Furthermore, MPC was considered to be applied to optimize the control result and reduce the resource cost of the whole system. The control effect has been studied and analyzed by theory, simulations and experiments, which indicates that the adjustment time can be reduced, almost no overshoot occurs while real-time tracking and robustness have also been improved. For the future work, the realization of MPC in the system will be studied. Moreover, the optimization of process models is decided to be implemented.

ACKNOWLEGEMENT

This research is partially supported by the XJTLU Key Programme Special Fund (KSP-P-02). The authors would like to thank all the parties concerned.

REFERENCES

- [1] G. K. Venayagamoorthy, "Dynamic, stochastic, computational, and scalable technologies for smart grids," IEEE Computational Intelligence Magazine, vol. 6, no. 3, pp. 22-35, 2011.
- [2] D. Shi, G. Gao, Z. Gao, and P. Xiao, "Application of expert fuzzy pid method for temperature control of heating furnace," Procedia Engineering, vol. 29, no. 4, pp. 257-261, 2012.
- Y. Tunckaya and E. Koklukaya, "Comparative performance evaluation of [3] blast furnace flame temperature prediction using artificial intelligence and statistical methods," Turkish Journal of Electrical Engineering & Computer Sciences, vol. 24, no. 3, 2016.
- [4] H. Q. Wang, J. I. Chang-Yinga, T. Z. Liu, F. Gao, and J. Y. Xian, "Modeling and simulation of fuzzy self-tuning pid temperature control system," Computer Engineering, vol. 38, no. 7, pp. 233-235, 2012.
- G. Mantovani and L. Ferrarini, "Temperature control of a commercial [5] building with model predictive control techniques," IEEE Transactions on Industrial Electronics, vol. 62, no. 4, pp. 2651-2660, 2015.
- J. Chen, "System identification and control of fopdt or sopdt processes," [6] Advanced Materials Research, vol. 317-319, pp. 2393-2397, 2011.
- [7] B. Kestenbaum, Non-Linear Regression. Springer New York, 2009.
- [8] A. Simpkins, "System identification: Theory for the user, 2nd edition (ljung, l.; 1999) [on the shelf]," IEEE Robotics & Automation Magazine, vol. 19, no. 2, pp. 95-96, 2012.
- [9] K. R. Sundaresan and P. R. Krishnaswamy, "Estimation of time delay time constant parameters in time, frequency, and laplace domains, Canadian Journal of Chemical Engineering, vol. 56, no. 2, pp. 257-262, 2010
- [10] D. E. Seborg, T. F. Edgar, D. A. Mellichamp, and F. J. D. Iii, Process Dynamics and Control, no. 33, 2010.
- [11] C. Bohn and D. P. Atherton, "An analysis package comparing pid antiwindup strategies," IEEE Control Systems, vol. 15, no. 2, pp. 34-40,

Scale Invariant Privacy Preserving Video via Wavelet Decomposition

Chengkai Yu, Charles Fleming and Hai-Ning Liang

Abstract—Video surveillance has become ubiquitous in the modern world. Mobile devices, surveillance cameras, IoT devices, all can record video that can violate our privacy. One proposed solution for this is privacy-preserving video, which removes identifying information from video as it is produced. Several algorithms for this have been proposed, but all of them suffer from scale issues: in order to sufficiently anonymize near camera objects, distant objects become unidentifiable. In this paper we propose a scale invariant method, based on wavelet decomposition.

Index Terms-Privacy, anonymization, video

I. INTRODUCTION

Cameras and camera-embedded devices have become pervasive in our daily life, not only wearable devices such as GoPro cameras and Google Glasses, but also surveillance systems, IoT devices, or even drones threaten our privacy. Firstperson videos are especially becoming very popular among YouTubers and video bloggers. Wearable cameras are also widely equipped by the police for security and evidence gathering purposes. Many such videos are eventually being uploaded to the Internet and processing techniques are often used to remove sensitive information such as faces. However, studies have shown that video blurring techniques are unable to balance privacy with awareness of risky situations by the person being recorded [1]. As the privacy issues of such videos recorded by wearable cameras are attracting more attention by the public [2], people are paying more attention to the potential threats to their privacy, and methods to preserve privacy in video need to be developed.

The usual video processing technique for privacy protection is to anonymize the video by applying blurring effects on sensitive regions, so that the information is not observable by the viewers. However, common methods for anonymization processing often require human selection of the sensitive regions. Faces are the only thing that are blurred out in most videos. The effect might not be actually anonymous due to other information revealed in the videos, e.g. the person's body shape or information from the background and nearby objects. Other

All authors are with the Department of Computer Science and Software Engineering, Xi'an Jiaotong-Liverpool University, Suzhou, Jiangsu, China.(email:chengkai.yu14@student.xjtlu.edu.cn,

charles.fleming@xjtlu.edu.cn, haining.liang@xjtlu.edu.cn).

methods blur the whole video frame, but this has the side effect of making distant objects too indistinct to separate from the background.

II. RESEARCH AIM AND OBJECTIVES

As minor details can threaten people's privacy, real anonymity using an existing algorithm would produce a video that is so anonymous that it destroys most of the visual information, including the objects and background. However, this would generally make the videos themselves unusable for whatever purpose they were recorded for. Detecting and performing the blurring effect on different information in the video based on scale could potentially blur out the sensitive information but maintain usability.

III. LITERATURE REVIEW

Privacy issues in wearable cameras have been widely discussed with the growth of popularity in recording videos by portable devices. Wearable cameras are more portable than conventional video recording devices and have emerged as a popular way to capture a wide variety of experiences which threaten people's privacy more aggressively and pervasively. Nguyen et al. conducted an extensive study on how individuals perceive and react to being recorded in first-person videos [3]. Findings suggest that most people would like to be asked for permission, in case the recordings are shared with others, and most people would mind if they are being recorded without being notified.

A study on addressing privacy concerns from videos taken in first-person point-of-view evaluated the effectiveness of four techniques (face detection, image cropping, location filtering and motion filtering) at reducing privacy infringing content [4]. Minor information that could be linked back to an individual as a non-obvious threat to privacy was pointed out in the study. All four methods were not particularly effective and could still pose privacy concerns. Glasses-style wearable devices were investigated with respect to recording and privacy and how these devices differ from other classes of cameras [5]. The qualitative study found that for such subtle devices, reactions to recording can be affected by the perception of the recorder and whether or not they could be identified in the recording. It was also pointed out that people frequently change their perceptions with repeated exposure or change their views as they become active users of such devices.

Ethical considerations of wearable cameras were investigated by [6]. Guidelines for all involved and best practices for third

Manuscript received July 4, 2018.

parties were proposed to address the ethical considerations of wearable cameras. Considerations in pervasive recording technologies were assessed by [7] for an insight into how design, technology, and policy can work together for appropriate usage of such technologies. An interesting finding suggested that a lack of control over recordings could actually make recording more tolerable. In an inescapable situation everything has been decided, and people could rationalize it better than controllable situations. Wearable cameras are demonstrated as an important emerging method to provide personalized feedback and support in public health interventions [2]. Ethical approval and privacy concerns stand as the most significant barriers and require more research. An important goal is to create interventions that explain, control, and notify people about the technologies.

The issue of privacy protection has been widely discussed in computer vision especially when it comes to human recognition from sources such as robots [8], wearable cameras [9], [10], first-person video [11], or surveillances videos [12]. Human activity recognition has received a great amount of attention and recognition algorithms for various environments and activities were introduced [13], [14], [15] which highlight the importance of privacy preserving videos. However, studies that cover privacy issues were focused on approaches for protecting peoples' privacy but not on analyzing the actual privacy protection of the current anonymous algorithms.

IV. ANONYMIZATIONS ALGORITHM

As baseline algorithms, we compare with three commonly used blurring algorithms including Gaussian blur, downsampling, and super pixel, applying these to our test videos. The three algorithms tend to produce different blurry effects and thus were chosen for determining and assessing their degree of anonymization in surveillance and wearable settings. We compare our algorithm with these standard algorithms on both up close and distant objects in videos and show that our algorithm outperforms all three.

V. WAVELET TRANSFORMATION

A wavelet transformation decomposes the signal into several bands to capture and separate different characteristics of the original signal. The signal information of an object contains detectable differences which could be captured by one or more bands during wavelet decomposition. The edges and surfaces are separated by different sets of bands due to the different characteristics of their signal. Different levels tend to capture the changes in color or textures in original images.

The wavelet transformation technique used in the study is the discrete wavelet transformation (DWT). The discrete wavelet transformation of an image signal of a level is calculated by passing it into a filter bank where a series of filters perform levels of decomposition. The signal is transformed by a highpass filter which gives the output of the detail coefficients while the low-pass filter returns an approximation coefficient. As the image signal has been decomposed at the current level, the output of the low-pass filter will then be subsampled by 2 in the next level.

The Wavelet transformation anonymization algorithm (WTAA) demonstrates the potential of balancing anonymity and usage by performing scalable blurring effect based on the scale of the objects. We anonymize video by decomposing the video using the wavelet transform, selectively destroying certain wavelet coefficients, then reconstructing the image via the inverse transform.

We compare the performance of our algorithm by considering a series of videos with both near and far objects. The wavelet transformation anonymization algorithm shows a great improvement when objects are far away. Gaussian blur has good anonymization but preserves insufficient information for usage. Downsampling has the lowest level of anonymization and does not keep the outline of the object at a distance. Superpixel has a reasonable anonymization level and preserves part of the shape; however, would not be sufficient to be considered useful. Wavelet transformation can keep both the shape and the color at a comparable level to the original video and performs extremely well in preserving distance objects due to the inherent scale in the wavelet decomposition.

In a zoomed-in shot, the person is about 15 meters away from the camera. Gaussian blur does not preserve any shape nor color at this distance in order to maintain the level of anonymity. Downsampling preserves a degree of color when the color contrast to background is large but the shape is left out. Superpixel merges the person with adjacent superpixels and the color would be destroyed as well as part of the shape. WTAA preserves both the shape and the color enough for the figure to recognizable as a person. This effect is more pronounced when viewed as a video.



Fig. 1. Comparison of anonymization algorithms for distant figures



Fig. 2. Comparison in level of anonymity in the last frame.

ACKNOWLEGEMENT

This research is partially supported by the XJTLU Key Programme Special Fund (KSP-P-02). The authors would like to thank all the parties concerned.

REFERENCES

- C. Neustaedter, S. Greenberg, and M. Boyle, "Blur filtration fails to preserve privacy for home-based video conferencing," ACM Transactions on Computer-human Interaction, vol. 13, no. 1, pp. 1—36, 2005.
- [2] A. R. A. Doherty, S. E. S. Hodges, A. C. A. King, E. E. Berry, C. J. A. C. Moulin, S. S. Lindley, P. P. Kelly, and C. C. Foster, "Wearable cameras in health: the state of the art and future possibilities," *American journal of preventive medicine*, vol. 44, no. 3, pp. 320–323, 2013.
- [3] D. H. Nguyen, G. Marcu, K. N. Truong, J. Scott, M. Langheinrich, and C. Roduner, "Encountering sensecam: personal recording technologies in everyday life," *Proceedings of the 11th international conference on Ubiquitous computing*, pp. 165—174, 2009.
- [4] E. Thomaz, A. Parnami, and J. Bidwell, "Technological approaches for addressing privacy concerns when recognizing eating behaviors with wearable cameras," ACM UbiComp'13, pp. 739–748, 2013.
- [5] T. Denning, Z. Dehlawi, and Z. Kohno, "In situ with bystanders of augmented reality glasses: Perspectives on recording and privacymediating technologies," ACM CHI 2014, 2014.

- [6] P. Kelly, S. J. Marshall, H. Badland, J. Kerr, M. Oliver, A. R. Doherty, and C. Foster, "An ethical framework for automated, wearable cameras in health behavior research," *American journal of preventive medicine*, vol. 44, no. 3, pp. 314–319, 2013.
- [7] M. Massimi, K. Truong, D. Dearman, and G. Hayes, "Understanding recording technologies in everyday life," *IEEE Pervasive Computing*, vol. 9, no. 3, 2010.
- [8] D. J. Butlerand, J. Huang, F. Roesner, and M. Cakmak, "The privacyutility tradeoff for remotely teleoperated robots," *In ACM/IEEE International Conference on Human-Robot Interaction*, pp. 27–34, 2015.
- [9] R. Templeman, M. Korayem, D. Crandall, and A. Kapadia, "Placeavoider: Steering first-person cameras away from sensitive spaces," *In Network and Distributed System Security Symposium (NDSS)*, 2014.
- [10] Y. Poleg, C. Arora, and S. Peleg, "Temporal segmentation of egocentric videos," CVPR, 2014.
- [11] S. Narayan, M. S. Kankanhalli, and K. R. Ramakrishnan, "Action and interaction recognition in first-person videos," *CVPR 2014 Workshp*, 2014.
- [12] W. Niu, J. Long, D. Han, and Y. F. Wang, "Human activity detection and recognition for video surveillance," *IEEE International Conference* on Multimedia and Expo (ICME), 2004.
- [13] M. S. Ryoo, B. Rothrock, C. Fleming, and H. J. Yang, "Privacypreserving human activity recognition from extreme low resolution." in AAAI, 2017, pp. 4255–4262.
- [14] J. Dai, B. Saghafi, J. Wu, J. Konrad, and P. Ishwar, "Towards "privacypreserving recognition of human activities," *ICIP 2015*, 2015.
- [15] J. K. Aggarwal and M. S. Ryoo, "Human activity analysis: A review," ACM Computing Surveys 43:16:1–16:43., 2011.

Comparative Studies of Segmentation Algorithms

Yanda Zhu, Yuxuan Zhao and Ka Lok Man

Abstract—Digital Image Processing is a transforming process from analog signal to digital signal. It is an essential part of **Computer Science. In Image Processing filed, Image Segmentation** addresses issues about how to extract some regions that people are interested in, which is a challenging task for many scientists and scholars. In general, image segmentation algorithms could be divided into two categories: region-based and edge-based. This paper investigates and compares four common Image Segmentation algorithms include Thresholding, Region Growing Labeling, Edge Detection and Mean Shift algorithms. Due to image segmentation problems are similar with classification problems in Machine Learning, plotting Receiver Operating characteristic curve and Precision Recall curve to in algorithm evaluation phase. Results of experiments demonstrate that each algorithm has its advantage and disadvantage. The significance of comparative studies for different segmentation algorithms tries to improve existing algorithms or find new algorithms to achieve better performance.

Index Terms—Algorithm, Computer Vision, Image Segmentation, Machine Learning.

I. INTRODUCTION

IMAGE obtained from cameras is consists of continuous simulated signal. The images are converted into discrete 2-d digital matrix so that computer is able to process them. *Computer Vision* is an important subfield in *Artificial Intelligence* that deals with how computer understand highlevel information from images or videos [1]. For more complex image recognition problems, computers usually need to extract features that task needs. Image Segmentation techniques solve many object extraction problems, which become a fundamental step in computer vision problems. In general, Image Segmentation partitions an image into several disjoint regions or region of interest (ROI), which keeps similarities to some extent according to their gray level, color, texture, shape and other characteristics.

This work is partially supported by the XJTLU Key Programme Special Fund (KSP-P-02) and the Grant Subsidy (XJTLU RIBDA2018-IRP1) provided by the Research Institute of Big Data Analytics (RIBDA). Yanda Zhu and Yuxuan Zhao are with the Xi'an Jiaotong- Liverpool University, China (e-mails: yanda.zhu16@student.xjtlu.edu.cn and yuxuan.zhao@xjtlu.edu.cn). Ka Lok Man are with the Xi'an Jiaotong- Liverpool University, China and Swinburne University of Technology Sarawak, Malaysia (e-mails: ka.man@xjtlu.edu.cn).

This papers reviews four common Image Segmentation algorithms: threshold-based, region-based, edge-based and a kind of feature space analysis method-mean shift algorithm [2]. Different Segmentation algorithms are suitable for processing different kinds of images according to their gray level value, texture, color and boundary characteristics.

The aim of Image Segmentation is in order to extract meaningful information from the original images. The accuracy of segmentation and algorithm time complexity becomes two important criteria to determine whether algorithm is good or not. In general, evaluation methods can be mainly divided into two categories- Subjective Analysis and Objective Experiment Method [3]. In this paper, we mainly focus on experimental methods to evaluate the accuracy of segmentation. Once regarding Image Segmentation problem as a classification problem in Machine Learning so that we can evaluate algorithms quantitatively using ROC curve and Precision-Recall (PR) curve to evaluate algorithms.

II. BACKGROUND

Over the years, image processing techniques have been widely applied into various fields such as industry application, medical images and traffic surveillance systems. Image Segmentation is the pre-stage for some further object recognition and feature extraction of high-level computer vision techniques [2].

From the perspective of mathematics, the segmentation of image \mathbf{R} is a finite set composed of $R_1,...R_s$, submits to no intersection between any two subsets. Supposing that $R_1, ..., R_s$ is the finite set is the complete segmentation of Image \mathbf{R} . The segmentation process could be denoted as the equation 1 shown.

$$\mathbf{R} = \bigcup_{i=1}^{S} R_i, R_i \cap R_i = \emptyset, \ i \neq j$$
(1)

From the past decades, many researchers have explored various algorithms based on the two criteria. One is homogeneity inside one region. The other is discontinuity between different segmented regions. Based on above two criteria, segmentation can be divided into edge detection methods and region growing or merging methods.

Except for edge detection and region merging or growing methods, Image Segmentation algorithms can be implemented by explicit theory. Around 2000, a new Image Segmentation method based on 'super-pixel' issued by Xiaofeng Ren attracts many researchers attention [3]. Super-pixel refers to irregular pixel patches with similar texture, color and brightness. It groups pixel according to their similarity of pixel, replacing a large number of pixels with less number of super-pixels. Another Segmentation algorithm combined with Graph-theory has been extensively used issued by Felzenszwalb in 2004 [4]. For this algorithm, images are regarded as undirected weight graph G=(V,E), E is limited edge set, V refers to nodes. Transferring Image Segmentation problem into minimum cut problem in graph. The advantage of this algorithm lies on each pixel owns different weight so that it is insensitive for shape of object. However, too long computation time is a huge problem.

Although many novel methods have been issued, exploring and improving Segmentation problem is still ongoing.

III. SEGMENTATION ALGORITHMS

In this paper, for a comparative studies purpose, four popular segmentation algorithms are compared. The comparison and simulation are carried out by MATLAB [6].

For threshold based segmentation, the general idea is to determine a suitable threshold according image gray level distribution histogram. If the gray value great than threshold, it becomes 1. Otherwise, it is 0. The procedure can be represented by the equation 2.

$$g(i,j) = \begin{cases} 1, \ f(i,j) > T \\ 0, \ f(i,j) \le T \end{cases}$$
(2)

In general, determining threshold between the peak and valley of histogram is simple but effective. For instance, the histogram of *Lenna.jpg* and *Girlface.png* is shown in Fig-1. The Girlface.png owns more distinct peak and valley. If we segment *Girlface.png* by T=50 or T=150, the outcome is excellent. But the *lenna.jpg* owns multiple irregularly peaks so it is hard to use peak-valley method to determine a suitable threshold. For Lenna, the other method to select a threshold was called *OTSU*, which proposed by a Japanese Scientist N, OTSU [7]. He thought the variance is a measurement of gray distribution. When the variance between foreground and background get larger, the performance of segmentation is better.



Fig-1 Gray Level Histogram for Girlface.png and Lenna.jpg

Edge based segmentation utilizes discontinuity of gray level, texture, color or other characteristics to mark different regions.

Edge describes positions where brightness changes dramatically. Edge is a vector variable with magnitude and direction. Due to the sensitivity of the edge detector towards noise, we always need to other preprocess techniques to prevent the influence of noise or distortion caused by edge detectors like contrast enhancement, smoothing. Edge detectors can be divided into two main categories: the gradient detector and Laplace detector [8]. Gradient Laplace detector owns rotational invariance.

Region based algorithms include Region Growing and Region Merging algorithm. The critical issue in the region growing algorithm is seed point selection and similar region determining rule. In this project, the seed point was selected automatically in the foreground by a function 'getpts' in MATLAB to interactively. The growing rule can be decided by humans. Comparing the gray value difference between seed point and neighboring pixel is the first stage. If difference less than a predefined threshold, the neighboring point can be inserted into the growing region. Until there are no pixel can be added in this region, growing processes stops.

These three algorithms are sufficient to process gray image. For color image, because multiple channels problem, the computation time and algorithm complexity will be increased. Mean Shift is a clustering algorithm which is suitable for segmenting color images. The principle is mapping each pixel in the original image into RGB or HSI feature space can obtain better performance than just process it in gay level [9].

IV. EVALUATION BY MACHINE LEARNING METHODS

For most of image processing problems, evaluating performance can be done by the subjective and objective perspective [13]. People observe experimental results are subjective and arbitrary. Here, combined with machine learning can evaluate quantitatively. First, importing reference image and turning it into binary image which foreground labeled as '1' in white and background labeled as '0' in black. Observing the Receiver Operating characteristic (ROC) curve and analyze area under the ROC curve and slope for threshold based and region growing algorithm. ROC curve is always used for evaluating a classifier in Statistics and Machine Learning [14]. In this condition, the ROC curve analyzes a binary classification problem. We call that true if the pixel in the ground truth image is "1". Otherwise is false. At the same time, when the gray level value great than threshold, the pixel is predicted positive. Otherwise is negative. So, we can generate a Confusion Matrix as table-1 shown. The equation (3) and (4) shows the formulation of True Positive Rate (TPR) and False Positive Rate (FPR). TPR is the probability of true examples were predicted correctly. FPR is the probability of false examples were incorrectly predicted as true.

	Actual True	Actual False
Predicted Positive	TP	FP
Predicted	FN	TN
Negative		

Table-1 Confusion Matrix

$$TPR = \frac{TP}{TP + FN} \tag{3}$$

$$FPR = \frac{TN}{TN + FP} \tag{4}$$

By calculating the TPR and FPR to draw Receiver Operation Curve (ROC) curve at different threshold so that transforming Image Segmentation problem into a classification problem in Machine Learning is easy [14]. The curve which approaches (0,1) is better. Fig-2 shows results of the ROC curve for region growing and threshold algorithms, so we can conclude that threshold performs better the than region growing algorithm in *Lenna.jpg*.



As for edge based segmentation algorithm, mainly focus on extract edge or boundary feature from the original image. This task is more like information retrieval rather than information division. Particularly, if true negative is not much valuable to a problem, or negative examples are abundant. Then, Precision Recall (PR) curve is typically more appropriate to evaluate this type of segmentation algorithm. 'P' refers to Precision, 'R' is Recall [14]. Equations as follows shown:

$$p = \frac{TP}{TP + FP} \tag{4}$$

$$r = \frac{TP}{TP + FN}$$
(5)

According the performance of different edge detectors in one same image, table-2 can be concluded.

Туре	Characteristics	Performance	Robustness
Sobel	Weighted	Normal	Depend on
	Average		size of mask;
	U		sensitive
			towards
			noise.
Prewitt	Approximate 1 st	Normal	Depend on
	derivative		size of mask
Canny	Two thresholds	Best	High
			robustness
			towards
			noise
LoG	2 nd derivative	Bad	Sensitive to
	Zero Crossing		noise;
			isotropy
	TT 1 1 0 C	C 11 CC . 1 1	

Table-2 performance of different edge detectors

V. CONCLUSIONS

The findings of this paper are the following:

- Threshold based algorithm is a fast but simple algorithm. For non-complex tasks, threshold based algorithm is effective to distinguish background and foreground.
- Region growing algorithm can overcome existing segmentation discontinuous disadvantage. However, region growing algorithm also has many disadvantages. For instance, unsuitable seed point selection directly leads to different segmentation outcomes. In addition, for some inhomogeneous grey level distribution or blurring edge images may cause under segmentation or over segmentation.
- Edge based segmentation might lose much details but essential structural properties can still be remained. As mentioned before, distinct edge detectors own different characteristics and performance. Different from the above two algorithms, edge detection extracts edge information from image rather than partition background such as region based segmentation algorithms. For some complex images, edge detection is easily to lose some relevant information.
- Mean shift is a non-parameter Probability Density Function (PDF) estimation algorithm, which seeks local maximum point along with the PDF gradient [8]. In addition, this algorithm does not need too much preprocess, and the iteration times are limited.
- As for evaluating segmentation algorithms, transferring segmentation problem into classification or retrieval problem in Machine Learning. By drawing ROC or PR curve to evaluate different algorithms is easier than from complex theory perspective.

REFERENCES

[1] R. C. Gonzalez, and R. E. Woods, *Digital Image Processing*: Upper Saddle River, N.J. : Pearson/Prentice Hall, c2010. 3rd ed., 2010, pp105-106 [2] M. Sonka, V. Hlavac, and R. Boyle, *Image Processing, Analysis, and Machine Vision*: Cengage Learning, 2014.

[3] L. Alphonsa and R. Resmi, "Performance metrics comparison of various image segmentation methods," in *International Conference on Control, Instrumentation, Communication and Computational Technologies*, 2016, pp. 410-414.

[4] X. Ren, C. C. Fowlkes, and J. Malik, "Learning Probabilistic Models for Contour Completion in Natural Images," *International Journal of Computer Vision*, vol. 77, no. 1-3, pp. 47-63, 2008.

[5] P. F. Felzenszwalb, and D. P. Huttenlocher, "Efficient Graph-Based Image Segmentation," *International Journal of Computer Vision*, vol. 59, no. 2, pp. 167-181, 2004.

[6] G. Rafael, and P. Wintz. *Digital Image Processing Using MATLAB*. Publishing House of Electronics Industry, 2013, pp55-58.

[7] N. Otsu, "A thresholding selection method from grey-level histogram," *IEEE Transactions on Systems Man & Cybernetics*, vol. 9, no. 1, pp. 62-66, 1979.

[8] L. Xiao, W. Hong, and L. Jian, "Sub-pixel edge detect technique based on LOG operator," *Journal of Baotou University of Iron & Steel Technology*, 2002.

[9] A. K. Jain, "The fundamentals of digital image processing prentice-hall," *Chemical Society Reviews*, vol. 24, no. 4, pp. 243-250, 1989.

[10] T. Yang, P. Zhikang, T. Min, W. Pingan, and X. Dechen, "Image Segmentation algorithm based on hierarchical mean shift," *Computer Research and development*, vol. 46, no. 9, pp. 1424-1431, 2009.

[11] Z. Liubing, and H. Shuo, "Image dimensionality reduction based on HSI color model," *Modern Electronics Technique*, no. 14, pp. 79, 2013.

[12] D. Comaniciu, and P. Meer, "Mean Shift: A Robust Approach Toward Feature Space Analysis," *IEEE Transactions on Pattern Analysis & Machine Intelligence*, vol. 24, no. 5, pp. 603-619,2002.

[13] G. K. Harrup, "ROC analysis of IR segmentation techniques," M.S. Thesis Air Force Inst. of Tech., Wright-Patterson AFB, OH. School of Engineering. 1994.

[14] Z. Zhihua, *Machine learning* Beijing: Tsinghua University press, 2016.pp 63-65.

Classification of Genetic Mutations for Cancer Treatment with Machine Learning Approaches

Gangmin Li and Bei Yao

Abstract— Cancer treatment is one of the most important and difficult tasks in medical treatment. Its success will influence millions of people's life. It is extremely important to develop an understanding of genetic mutations in cancer tumors so that the precise therapies can be applied for patients. However, due to the large amount of manual work required in the interpretation of genetic mutations and their effects it is difficult to classify genetic mutations based on clinical text. This paper reports our research on classification of genetic mutations with an approach of constructing a machine learning model with a hope to improve the performance of classification of genetic mutations.

Index Terms— Cancer Treatment, Machine Learning, SVM, Text Classification, XGBoost.

I. INTRODUCTION

ANCER is a leading cause of death all over the world and it is still increasing in alarming rate. The progress on cancer treatment has been quite slow. This is due to the complexity nature of cancer and its appropriate treatment. Personalized medicine and treatment are studied that takes genetic makeup into account to attempt to maximum efficiency and minimize toxicity for each patient by using the right drug with right dosage [1]. However, personalized medicine and treatment require analyzing the cause of cancer and its trail treatment. These involve large amount of data to be analyzed which cause manual analysis difficult and even impossible [2]. In addition, mutations have different consequences to the cancer development. Some of them contribute to tumor growing while others might be neural mutations [3]. It requires background information and nuanced decision-making for analyses the causes. It further increases the difficulties to classify available clinical datasets. Currently, the process of distinguishing genetic mutations is being done manually which is really laborious work since scientists need to review and classify each mutation manually according to clinical evidence obtained from text-based literature [3].

Until recently, it has gradually become a broad consensus that machine learning can bring significant improvements to genome sequencing datasets analysis [4]. But there is still lot of work to be done to exploiting different classification methods based on their performance on classification [5].

This paper reports a project aims to build a machine learning model for efficient classification of genetic mutations to further enable personized medicine for cancer treatment. The paper is organized as follows. Section II describes methods used including experimental dataset and approaches. Results of different approaches are presented in Section III. Section IV discusses and analyzes the results in detail. Section V concludes the paper and further works are identified.

The paper is organized as follows. Section II describes methods used including experimental dataset and approaches. Results of different approaches are presented in Section III. Section IV discusses and analyzes the results in detail. Section V concludes the paper.

II. RESEARCH METHODS

Machine learning uses artificial intelligence and statistical techniques to let computers to learn progressively to improve performance on a specific task with data with or without out being supervised [6]. Within the field of data analytics, machine learning is a method used to devise complex models and algorithms that lend to prediction. It is also known as predictive analytics. These analytical models allow us to produce reliable, repeatable decisions and results and uncover hidden insights through learning from historical relationships and trends in the data [7]. Depends on whether there is a learning feedback available, machine learning are typically classified into two broad categories supervised and unsupervised learning.

In our research the task is to classify genetic mutations to enable personized medicine for cancer treatment. The data comes from public available source provided by the Kaggle competition named "Personalized Medicine: Redefining Cancer Treatment" [3]. Due to the nature of given dataset is text-based, a conversion is needed for any classification algorithms can adopted. We used TF-IDF feature extraction and conversion; we used XGBoost and SVM for classification to compare the performance. The evaluation of two classification methods was done based on multi class log loss, which is specified by the Kaggle website.

This work was supported by the RIBDA and Jiangsu DOE for Branding Programme of IMS.

Gangmin Li is with the Department of Computer Science and Software Engineering, Xi'an Jiaotong-Liverpool University. Suzhou, Jiangsu, PR. China. 215123, phone: +86 521 88161510; fax: +86 88161500; e-mail: Gangmin.li@xjtlu.edu.cn.

Bei Yao is with the Department of Computer Science and Software Engineering, Xi'an Jiaotong - Liverpool University. Suzhou, Jiangsu, PR. China. 215123, phone: +86 521 88161510; fax: +86 88161500; e-mail: bei.yao@xjtlu.edu.cn.

A. Dataset

The dataset used in our research came from the Kaggle competition and Memorial Sloan Kettering Cancer Center (MSKCC). The dataset has three parameters: *genes*, *variations* and *clinical text*. In the training set, 9 classes of mutations are given. Our goal is to use these three provided variables to predict mutation classes. It can be seen from the Table 1 that the number of data entries in the training set which combined by the original training set and some released solutions of test set. The training set has 3689 data entries which are about four times as as much as the test set. A similar distribution can also be seen with the *types of variations*. However, there are more types of *genes* in the test set compare with the number of types in the training set.

TABLE I. DATA STATISTICS



Fig. 1. Imbalanced distribution of different genetic mutation classes

We look into the distribution of different genetic mutation classes in train set. According to the figure 1, an obvious imbalanced distribution can be seen. Class 3, 8 and 9 are underpresented with less than 100 samples, this will cause problems and affect the accuracy of a prediction model.

B. Approach

Our approach is illustrated in Figure 2. It involves two major steps: firstly, the given training dataset in a form of text needs be cleaned and useful features are extracted for training. Then, classification methods are adopted for the model construction, which given data are enriched to overcome data imbalance problem.



Fig. 2. The model used in our classification of genetic mutations

Since the classification of genetic mutations was on the basis of clinical evidence which came from related medical literatures, features were required to be extracted from text for further use. Text files are actually sets of words while genes and variants are composed of characters. In order to apply machine learning algorithms, we need to convert the text fields into numerical feature vectors. We applied label encoder to convert gene and variation into feature vectors for training, which can help us to use non-linear features by normalizing them. In our study, TF-IDF was used to transform text into feature vectors which could analyze frequency of words in literature and further used as input to estimator. Before the transformation, we have cleaned the text by removing all punctuations and keeping only regular characters in the literature. Since the literature was composed of meaningful words, we set analyzer to 'word' in order to extract features in word format.

Given training set is imbalanced and for certain classes there is no enough samples to support, we implemented an oversampling method to fix the under-presented problem. Specifically, we used Synthetic Minority Oversampling Technique (SMOTE) [8], which can randomly pick a point from the minority class and calculate k-nearest neighbors for the point, then add those points somewhere between the selected point and each of its neighbors.

Regarding classification model, Both XGBoost and SVM were used. They are two popular classification algorithms for small-scale datasets without overfitting.

XGBoost is an open-source software library which provides the gradient boosting framework. Gradient boosting is a machine learning technique for regression and classification problems, which produces a prediction model in the form of an ensemble of weak prediction models. It builds the model in a stage-wise fashion like other boosting methods do, and it generalizes them by allowing optimization of an arbitrary differentiable loss function.

In our classification of genetic mutations problem we have an output variable y (mutation classes) and a vector of input variables x (genes, variations and features vector form clinical text) that y = F(x). Our training set { $(x_1, y_1), ..., (x_n, y_n)$ } of known values of x and corresponding class of y, the goal is to find an approximation $\hat{F}(x)$ to a function F(x) that minimizes the expected value of a log loss function L(y, F(x)).

$$\hat{F} = \underset{F}{\arg\min} \mathbb{E}_{x,y}[L(y, F(x))].$$

Clearly it is an optimization problem. In many cases it is computationally infeasible to find the best function \hat{F} for an arbitrary loss function. We used log loss because it is simple and computationally inexpensive. The log loss is actually Logarithmic loss. It measures the performance of the classification model where the prediction input is a probability value between 0 and 1. The goal of our model is then becomes to minimize this value. A perfect model would have a log loss of 0. When log loss increases as the predicted probability diverges from the actual label.

The good thing is we are not worried about the mathematical complexity. Many integrated Data analysis tools provide build in XGBoost library. We used Weka, a Waikato Environment for Knowledge Analysis (Weka) is a suite of machine learning software written in Java, developed at the University Of Waikato, New Zealand. It is free software licensed under the
GNU General Public License.

SVM, support vector machines are supervised learning models with associated learning algorithms that analyze data used for classification and regression analysis. Traditionally, given a set of training examples, each marked as belonging to one or the other of two classes, an SVM training algorithm builds a model that assigns new examples to one class or the other, making it a non-probabilistic binary linear classifier. However in our problem we used Platt scaling and then the SVM can be models in a probabilistic classification setting. Basically, an SVM model is a representation of the examples as points in space, mapped so that the examples of the separate classes are divided by a clear gap that is as wide as possible. New examples are then mapped into that same space and predicted to belong to a class based on which side of the gap they fall. Again we used Weka software to test SVM classifier. We have experimented three different approaches to split given training set into training and validation subsets and compare their performances.

The first was to split original training dataset in sequence, which meant that we used the first 80% of data for training and the rest 20% for validation. Secondly, simple random sampling was adopted and finally, stratified random sampling was applied. For stratified random sampling, genetic mutation classes were used as labels to split data in a stratified fashion. To avoid contingency, for each random sampling method, validating process was carried out 3 times to obtain an average value.

III. RESULTS AND EVALUATION

The results of different approaches to select subsets were evaluated on both of public and private leaderboard. On the Kaggle website, public board is based on 76% of the test set while private board is based on other 24% of data which presents the final result of the competition. Also, validation sets' performances were also judged based on their log loss and accuracy.

The comparison is provided in the table 2. It can be seen the improvements on the validation set resulted from simple random sampling compared to sequential sampling. It is also reflected on the test set which means that the performance of models was actually improved instead of overfitting. However, the simple random sampling seems to obtain an even better result on validation set compared to stratified random sampling. Stratified random sampling performed better on test set, which means that validation subset chosen by simple random sampling might not be representative enough. We will further discuss the representativeness of validation sets generated by different methods in the next section.

 TABLE II.
 COMPARISON
 AMONG
 DIFFERENT
 APPROACHES
 TO
 SPLIT

 SUBSETS

 </t

Models	Public Board	Private Board	Log Loss	Accuracy
TF-IDF + Sequential sampling (8:2) + XGBoost	1.80	2.48	1.08	0.60

TF-IDF + Simple Random Sampling (8:2)+ XGBoost	1.74	2.52	0.88	0.68
TF-IDF + Stratified Random Sampling (8:2)	1.71	2.46	0.93	0.66
+ XGBoost				

After implementation of SMOTE method, the log loss and accuracy of the validation set were significantly improved owing to oversampling as indicated in the table 3. The log loss decreased from 0.91 to 0.37 and the accuracy increased from 0.67 to 0.87. However, this improvement did not reflect on the test set while there was no great change of the scores on the leaderboard and the score on the public board was even a little bit worse compared to the result before oversampling.

TABLE III. RESULT OF EMBEDDING SMOTE

Method	Public Board	Private Board	Log Loss	Accuracy
TF-IDF + SMOTE + Stratified Random Sampling (8:2) + XGBoost	1.73	2.21	0.37	0.87

For the two classification methods we used. It can be easily seen from the table 4, where the XGBoost model outperformed the SVM model both on the validation set and the test set when the other methods were all the same.

TABLE IV. COMPARISON BETWEEN PERFORMANCES OF DIFFERENT MODELS

Method	Public Board	Private Board	Log Loss	Accuracy
TF-IDF + Stratified Random Sampling (8:2) + XGBoost	1.71	2.46	0.93	0.66
TF-IDF + Stratified Random Sampling (8:2) + SVM	1.87	2.55	0.99	0.62

To better understand the hidden details of our predications, confusion metrics were applied and visualized for the two classification models' performances. Figure 3 and 4 are the results. It could be seen that both of the models XGBoost and SVM successfully distinguished more than half of those instances that belonged to class 1, 4 and 7, which were classes that occur most frequently in the training set. And XGBoost performed better in distinguishing those classes compared to SVM does. Also, it is worth mentioning that although class 3, 8 and 9 were all under-presented, more than half of data entries of class 9 were correctly predicted in both models especially in the SVM model. There existed some classes that often were mistakenly predicted to be another one, including class 2 and class 7, class 4 and class 1. Also, it is obvious on the figures that the XGBoost model has a better predictive ability than the SVM model has.



Fig. 3. Confusion Matrix of The SVM Model



Fig. 4. Confusion Matrix of The XGBoost Model

IV. DISCUSSION

The preliminary results from our study has demonstrated the potential of embedding machine learning algorithms at least XGBoot and SVM into classification of Genetic Mutations, which is text-based and probabilistic classification problem. This shed lights on personalized medicine to reduce tedious manual work and expedite the process in analyzing cancer causes. In this section, we will further interpret this project based on encountered problems in the modeling process and analysis of our results.

Firstly, TF-IDF, the method we used in this problem, is a commonly used word embedding method for feature extraction. However, it does not take semantical similarities and orders of words into consideration; instead, it focuses only on the frequency of each word. In our study, clinical text came from scientific literatures and a certain number of words or phrases should be informative, it should be meaningful to analyze the semantics and word order for better learning of features.

Secondly, in our model construction, we treated *genes*, *variations* and *clinic text* as three independent attributes since we did not have any domain knowledge in biomedical field and genetic mutations. It is the easiest way to extract features to solve this problem. However, in real application, mutations should not be annotated in this way. It can also be seen from confusion metrics that there were certain pairs of classes often

mistakenly distinguished which might be one of drawbacks of simply concatenate three parameters. It can be supposed that experts should firstly check the types of *genes* and their *variants* to find related literatures, which could help them identify mutation classes.

Thirdly, it can be seen from the comparison table that simple random sampling could lead to a better performance of the model than sequential sampling, which suggests that the distribution of genetic mutations probably exists a certain relationship with their orders in the dataset. Nevertheless, random could be a choice that still not good enough in this case due to the imbalanced data. Although 80/20 should be a good splitting for this problem, it is still possible that some classes occur only in training or validation subsets since the dataset is high skewed which might cause the training subset not representative enough. Thus, stratified random sampling also helped slightly improve the performance of model compared to simple random sampling on the test set. It reveals that it could select a more representative training subset from the whole training set.

Fourthly, oversampling was initially expected to reduce a bias due to highly-skewed distribution of the given data set. However, the improvement of the validation set did not reflect on the whole test set, instead, only the score on the private board got an obvious improvement which presented as the best performance with a score of 2.21, ranking about top 6% on the leaderboard. Nevertheless, since the score on the private board is based on a smaller subset of test set (24%), it is possible that this relatively good performance is due to bias caused by limited size and unbalanced distribution of data. In general, oversampling method did not lead to a significant change on the predictive ability for the test set. It seems that SMOTE solved under-presented problem of certain classes but caused other problems that could lead to overfitting in the same time. Therefore the overall performance of models on the whole test set was basically the same. There are two main possible reasons for this: one is that it might mistakenly amplifies and strengthens some noise in the training set by KNN algorithm since in biomedical data noise is unavoidable, thereby the risk of overfitting increased. Another possible reason is that, as we mentioned in the earlier sample data statistics, that there are more types of genes in the test set, it will cause certain difficulties in the prediction since the model was not able to learn about some genes in the training process.

Lastly, regarding the performances of the two classification models, it is clear that XGBoost outperformed SVM in this case. From the confusion metrics, it could be seen that XGBoost had a better predictive capacity especially for those sufficient presented classes. It may be caused by its features' regularization. However, for the parameter tuning, we only attempted several parameters of SVM model and left other parameters as default values since it requires a large amount of knowledge and testament for parameter tuning. It might cause the relatively poor performance of SVM.

V. CONCLUSION AND FUTURE WORK

This paper reports our ongoing study of machine learning for Genetic Mutations Classification. We firstly explored the data by applying transformation on data entries and features. Visualization was also utilized to display the data structure in order to understand give dataset. Then we managed to implement a commonly used word embedding method, TF-IDF, for feature extraction. The efficiency of different methods of subsets choosing was compared and their performances were discussed in this paper. Alternative ways to sample was conducted and its performance was discussed with possible leading causes. In addition, we implemented two different classification models, XGBoost and SVM, which were expected to be better employed for small-scale datasets. The results demonstrated the feasibility of embedding machine learning algorithms into Genetic Mutations classification.

There is still much work to do for the further study. Firstly, for word embedding, except for TF-IDF, word2vec or doc2vec, which take semantic similarities and orders into consideration, is expected to be implemented for the improvement of proposed models for further development. Secondly, it is meaningful to explore the relationship between *genes* and their *variations* and combine genes and variations together to find related literatures. Thirdly, for the parameter tuning, instead of only several parameters, other parameters could be further explored with the usage of grid search or random search to simplify the process of parameter tuning and optimize the model for a better performance.

REFERENCES

- M. Verma, "Personalized Medicine and Cancer," *Journal of Personalized Medicine*, vol. 2, no. 1, pp. 1-14, 2012.
- [2] A. Holzinger, "Trends in interactive knowledge discovery for personalized medicine: cognitive science meets machine learning," IEEE Intell Inform Bull, vol. 15, no.1, pp. 6-14, 2014.
- Kaggle, "Personalized Medicine: Redefining Cancer Treatment," 2017. [Online]. Available: https://www.kaggle.com/c/msk-redefining-cancertreatment.
- [4] M. W. Libbrecht and W. S. Noble. (2015, Mar.). Machine learning applications in genetics and genomics. Nature Reviews Genetics. [Online]. 16(6), pp. 321-332. Available: https://doi.org/10.1038/nrg3920
- [5] A. Bhola and A. K. Tiwari, "Machine Learning Based Approaches for Cancer Classification Using Gene Expression Data", Machine Learning and Applications: An International Journal (MLAIJ), Vol. 2, No. 3/4, Dec. 2015.
- [6] Samuel, Arthur (1959). "Some Studies in Machine Learning Using the Game of Checkers". IBM Journal of Research and Development. 3 (3): 210–229. doi:10.1147/rd.33.0210.
- [7] "Machine Learning: What it is and why it matters". www.sas.com.
- [8] N. V. Chawla, K. W. Bowyer, L. O. Hall and W. P. Kegelmeyer. (2002, June). "Smote: Synthetic minority over-sampling technique," Journal of Artificial Intelligence Research. [Online]. 16, pp. 321-357. Available: https://doi.org/10.1613/jair.953

Evaluate Cancer Patients Quality of Life after Receiving TIVAPDS Treatment through Big Data Analytics

Gangmin Li, Shiyang Zhang, and Xuming Bai

Abstract—Cancer plays a leading role in causing morbidity and mortality worldwide. Several treatments have been developed and practiced for fighting against cancer. TIVAPDS treatment is a new method utilizing TIVAP drag delivery method, which is one kind of Intrathecal Drug Delivery method (IDD) with lower side effects and increased patient's quality of life (QOL). This paper reports our study aiming to evaluate the TIVAPDS treatment especially on patients' QOL in order to make contributions to generalize the treatment in China. Our data samples come from The Second Affiliated Hospital of Suzhou University, a forerunner of TIVAPDS practices in China and with patients' agreement. The methods adopted are BDA techniques. By creating of a prediction model based on historical data, we could predict the complications caused by the implant of TIVAP therefore taking necessary measures in advance to improve QOL for a patient.

Index Terms— Data analysis, Logistic regression, Prediction model, Quality of life (QoL), TIVAPDS.

I. INTRODUCTION

ANCER plays a leading role in causing morbidity and mortality worldwide. According to the statistics from WHO (World Health Organization), in 2015, the number of people who died in cancer is about 8.8 million, and it is nearly one-sixth of global deaths. In order to fight with cancer, several methods have been developed and practiced, including surgery [1], chemotherapy [2, 3], radiation therapy [4, 5], hormonal therapy [3], targeted therapy (including immunotherapy) [6] and synthetic lethality [7]. These methods deliver drugs into the intrathecal space utilizing oral or intravenous infusion associated with long-term and discontinuing venous transfusion. Recently, the Intrathecal drug delivery system (IDDS), a new method with lower side effects hence improving patient's quality of life (QOL) become popular globally [8]. It contains three delivery methods, including Central Venous Catheter (CVC) [9], Peripherally Inserted Central Catheter (PICC) [10], and a Totally Implantable Venous Access Port (TIVAP) [11].

TIVAPDS (Totally Implantable Venous Access Port Drug Supply) is a treatment utilizing TIVAP delivery method. Applying this therapy, it is necessary to insert one long hollow tube into one of the large veins in patient's body. One end of this tube is located in the vein, and it usually just above the heart. Meanwhile, the other end is connected to the injection port. It is under patient's skin on the chest. Medicines could be delivered through this certain port rather than injecting into veins. Simultaneously, the port inserted under skin is not easy to be noticed. It is helpful for reducing the chances of infection, low maintenance, consequently improves patient's quality of life (QoL).



Fig. 1. Diagram showing poortacath in place

Comparative Effectiveness Research (CER), through a comprehensive analysis of the therapeutic data and feature data of patients, compared effective of the TIVAPDS method to other cancer therapies such as, CVC (Central venous catheters) and PICC (peripherally inserted central catheters). It finds that the TIVAPDS is the best treatment method for a wide range of specific cancer patients.

However, because of the short history and lack of research on TIVAPDS itself, both researchers and practitioners still hold split views on this method. Especially in opinions about the patients QOL after treatment in China, there are few hospitals and organizations applying TIVAPDS in treating cancer patients so far. As a forerunner, the 2nd Affiliated Hospital of Suzhou University plays an influential role in doing TIVAPDS researches and practices. In order to generalize this treatment and improve cancer patients' QOL, it is necessary to study its

This work was supported by the RIBDA and Jiangsu DOE for Branding Programme of IMS.

Gangmin Li and Shiyangzhang are with the Department of Computer Science and Software Engineering, Xi'an Jiaotong-Liverpool University.

Suzhou, Jiangsu, PR. China. 215123, phone: +86 521 88161510; fax: +86 88161500; e-mail: Gangmin.li, Shiyang.zhang@xjtlu.edu.cn.

Xuming Bai is with the Department of Interventional Radiology, The Second Affiliated Hospital of Suzhou University, Suzhou, China. Suzhou, Jiangsu, PR. China. 215123, e-mail: 2005baixuming@163.com

effectiveness and specially in terms of improving patients' QOL since a complete cure is hard to achieve. This paper reports our study on the effectiveness of improving patients' QOL.

The paper is organized as follows: Section II presents our research methodology, which literally follows the process of BDA. It includes Acquisition, which defines the goal of analysis, Understanding of data, preprocess of data, analyses of data and evaluation. Second III is evaluations where not only the model adopted are evaluated against data samples to assess the accuracy of the models, furthermore, evaluation of the approaches against ultimate goal of study are reported. Finally section IV we draw our conclusions and important issues are discussed.

II. METHODOLOGY

This study aims to study the effectiveness of TIVPDS specifically on improving patients' treatment QOL. Theoretically, the quality of life is a non-consensus concept. people have different understanding Deferent and interpretations. Since TIVPDS treatment implants an alien device into patients' body. It will inevitably affect patients' normal movement and normal living life. However compare with other treatments, it is believed that it can minimize the discomfort and maintain quality of life. This opinion is challenged by many people including patients. Obviously perceptions of QOL are affected by various other factors including medical condition, physiology condition and patient's personal feeling and psychology factor etc. The medical and physiological parameters can be measured and compared with the standard index but other factors cannot be easily expressed and measured. Due to the complexity of the measurement, they may provide bias results for any analysis. This research used a simple but important indicator, after a discussion with doctors. We are suggested specifically study the complications of implant operation. This is because it is simple and easy to measure. The historic data provides different complications with each specific case are used to build a prediction model to forecast the potential risks of each complication therefore can be avoided in the treatment, therefore to increase the comfort and QOL of the patients.

A. Data Samples

Our data comes from the The Second Affiliated Hospital of Suzhou University. To protect patients' privacy all data were desensitized in advance. This hospital is one of the earliest and most influential hospitals conducting TIVPDS research and practice in China. The hospital receives not only natives of Suzhou but also patients across the country. Moreover, acting as a national training center for TIVPDS, this hospital has advanced treatment facilities supporting TIVPDS treatments and real-time clinic data, both non-structured and structured collection.

From 2014 to 2016, a total 1055 patients with cancer had received systemic TIVPDS treatment in the hospital. Among of 1055 cases, there are 366 cases had 10 different complications. They are *Pinch-off syndrome*, *Catheter Fibrin Sheath and DVT*,

Catheter Heterotopia, Catheter Rupture and Fracture, Port Device Related Infection, Local Skin Problems: Ulceration Defects, Pneumothorax, Incision local infection, Incision dehiscence, and Inadvertent Arterial Puncture. Among the patients who had complications, apart from age, gender and the diseases they suffer, the blood sample shows in normal range: White Blood Cell (WBC) count within (3.50-9.50*109/L), within (223-341.5*109/L), Platelet (PLT)count and Prothrombin Time (PT) within (11.5-14.0/S), and the international normalized ration (INR) within (0.9-1.15), therefore can be taken as predictors. However the data samples shows the patients who suffer from complications had 14 different diseases: Malignant thymoma, Postoperative recurrence of hepatocellular carcinoma, Obstructive jaundice, Colon cancer, Prostate cancer with bone metastasis, Breast cancer, Esophageal cancer, Gastric cancer, Postoperative gastric cancer, Pancreatic cancer, Rectal cancer, Postoperative rectal cancer, Left lung cancer, and, Left breast cancer.

The sample data summary is presented in Table I. Among these 365 data samples, 55.6% are female while the rest 44.5% are male. Most of these data sample are aged between 41 and 70 counted as 75.7%. The top two cancers are breast and Gastric Cancer counted as 40.8%.

 TABLE I.
 Sample data Summary

Factor	Number	-
Gender		
Female	203 (55.6%)	
Male	162 (44.4%)	
Age (years)		
10-20	2 (0.54%)	
21-30	8 (2.2%)	
31-40	28 (7.7%)	
41-50	81 (22.2%)	
51-60	101 (27.7.%)	
61-70	94 (25.8%)	
71-80	44 (12%)	
81-90	7 (2%)	
Diseases		
Rectal Cancer	36 (9.9%)	
Gastric Cancer	57 (15.6%)	
Breast Cancer	92 (25.2%)	
Ovarian Cancer	13 (3.6%)	
Colon Cancer	34 (9.3%)	
Liver Cancer	15 (4.1%)	
Lung Cancer	38 (10.4%)	
Others	80 (21.9%)	

B. Prediction model

Build a prediction model to predict the implications are a simple task. However, an accurate prediction needs variables and their historical value to be accurate. One of concern is the independence or correlation to be evaluated. Another study [12] indicated there is no obvious correlation among the individual predictors and the complications. Furthermore, no linear relationship between predictors and the dependent variable can be found. Therefore a straight forward logistic regression model was adopted to build our predictor.

A commercial package SPSS (Statistical Product and Service Solutions) for Windows (Version 19.0), was used in our analyses. SPSS supports building logistic regression equations, which can forecast the probability of a new cancer patient received TIVAPDS who suffer from each complication. All 10 complications were individually analyzed with 10 different logistic regression equations.

Constructed and test our model, data samples are randomly split into two groups 300 samples used as training sets and the rest 65 samples used as test set. The coefficients of the logistic regression equation for complication were obtained through training set. The coefficient of gender (x_1) is -0.225, the coefficient of age (x_2) is -0.019, and the coefficient of disease (x_3) is -0.027. The coefficient of constant is -0.983. Therefore, the probability of the logistic regression equation for complication is shown in (1).

$$\ln\left(\frac{f(x)}{1-f(x)}\right) = -0.983 - 0.225x_1 - 0.019x_2 - 0.027x_3.$$
(1)

Which $t = -0.983 - 0.225x_1 - 0.019x_2 - 0.027x_3$.

Converted to

$$\delta(t) = \frac{e^t}{e^t + 1} = \frac{1}{1 + e^{-t'}}$$

and $\delta(t)$ is the probability of occurrence of complications.

With these models all the remaining 65 cases as the testing set. The models on all the test sets show 92.4% of accuracy rate.

For each individual complications similar approach were conducted. We have the following results as shown in Table II.

TABLE II. MODEL AND ACCURACY OF 10 COMPLICATIONS ON TEST SET

Complication	Model	Accuracy
Pinch-off	$t = -6.271 + 1.162\chi_1 - 0.009\chi_2$	98.48%
Syndrome	$+ 0.108\chi_3$	
Catheter	t = -2.723 = 0.249 x = 0.006 x	100%
Rupture and	$t = 2.725 0.219_{\lambda_1} 0.000_{\lambda_2} = 0.287 v_{-}$	
Fracture	0.207 /3	
Catheter Fibrin	$t = -13.904 + 0.357\chi_1$	100%
Sheath and DVT	$+ 0.057 \chi_2$	
Catherton	$+0.721\chi_{3}$	1000/
Catheter	$t = -4.181 - 0.306\chi_1 - 0.01/\chi_2$	100%
Reterotopia	$+ 0.087 \chi_3$	06 70/
Port Device	$l = 18.285 - 17.620\chi_1 - 0.089\chi_2$	96.7%
Infaction	$+ 0.013\chi_3$	
Illegration	$t = -38374 \pm 16137$ v	100%
Defects	$t = 30.374 + 10.137 \chi_1 + 0.039 \gamma_2$	100 /0
Dereets	$-0.389 \chi_2$	
Pneumothorax	$t = 19.528 - 17.206 \chi_1$	98.48%
	$-0.096\chi_2$	
	$-0.522\chi_{3}$	
Incision Local	$t = 11.695 + 0.615 \chi_1$	98.48%
Infection	$-0.021 \chi_2$	
	$-14.195 \chi_3$	
Incision	$t = -32.975 + 16.032 \chi_1$	100%
Dehiscence	$-0.023 \chi_2$	
Turanda and a set	$-1.143 \chi_3$	1000/
	$\iota = -5.098 \pm 0.540\chi_1$	100%
Anterial	+ 4.033 * $10^{-4} - 4x$	
Functure	+ 0.145v	
	1 0.115/23	

III. EVALUATIONS AND CONCLUSION

Our goal was analyze patients QOL after receiving TIVAPDS treatment. We were suggested to use complications caused by the treatment. We adopted logistic regression model since no leaner relation can be found among predictors and dependent variable. Our primary results after test were proving to be pretty accurate.

However, we understand the limitations of our study. First of all the sample size is small. It is hardly confident to be used for future accurate predictions. It can be further improved if more data becomes available. Secondly, the logistic model should be further evaluated by AUC (area under curve). Again due to sample size, AUC does not provide convincing results. Thirdly, some diseases has multiple implications which we don't explicitly take into consideration.

Nevertheless, our study is useful in that we attempted to predict patients QOL quantitatively using BDA methods after a treatment. It does give medical practitioners some hints that the potential implication may occur based on the historical data.

ACKNOWLEDGMENT

This research project was supported by the fund provided by the Research Institute of Big Data Analytics (RIBDA), Xi'an Jiaotong-Liverpool University and the 2nd Affiliated Hospital of Suzhou University.

REFERENCES

- Subotic, S., S. Wyler, and A. Bachmann, Surgical Treatment of Localised Renal Cancer. European Urology Supplements, 2012. 11(3): p. 60-65.
- [2] Mieog, J.S.D., J.A. van der Hage, and C. de Velde, Neoadjuvant chemotherapy for operable breast cancer. British Journal Of Surgery, 2007. 94(10): p. 1189-1200.
- [3] Abe, O., et al., Effects of chemotherapy and hormonal therapy for early breast cancer on recurrence and 15-year survival: an overview of the randomised trials. Lancet, 2005. 365(9472): p. 1687-1717.
- [4] Bao, S.D., et al., Glioma stem cells promote radioresistance by preferential activation of the DNA damage response. Nature, 2006. 444(7120): p. 756-760.
- [5] Stupp, R., et al., Radiotherapy plus concomitant and adjuvant temozolomide for glioblastoma. New England Journal Of Medicine, 2005. 352(10): p. 987-996.
- [6] Topp, M.S., et al., Targeted Therapy With the T-Cell-Engaging Antibody Blinatumomab of Chemotherapy-Refractory Minimal Residual Disease in B-Lineage Acute Lymphoblastic Leukemia Patients Results in High Response Rate and Prolonged Leukemia-Free Survival. Journal Of Clinical Oncology, 2011. 29(18): p. 2493-2498.
- [7] Kaelin, W.G., The concept of synthetic lethality in the context of anticancer therapy. Nature Reviews Cancer, 2005. 5(9): p. 689-698.
- [8] Lynch, L., Intrathecal drug delivery systems. Continuing Education in Anaesthesia Critical Care & Pain, 2014. 14(1): p. 27-31.
- [9] Klek, S., et al., Enteral and Parenteral Nutrition in the Conservative Treatment of Pancreatic Fistula: A Randomized Clinical Trial. Gastroenterology, 2011. 141(1): p. 157-163.e1.
- [10] Hoshal, V.L. and Jr, Total intravenous nutrition with peripherally inserted silicone elastomer central venous catheters. Archives of Surgery, 1975. 110(5): p. 644-646.
- [11] Shankar, G., et al., Totally Implantable Venous Access Devices in Children Requiring Long-Term Chemotherapy: Analysis of Outcome in 122 Children from a Single Institution. Indian Journal Of Surgical Oncology, 2016. 7(3): p. 326-331.
- [12] Gangmin Li, Jianze Gu and Xuming Bai, Study of the Effectiveness of TIVAPDS Treatment on Cancer through BDA, CICET 2018, Taiwan, China.

Chinese Microblog Sentiment Analysis by Adding Emoticons to Attention-Based CNN

Yi-Jen Su, Member, IEEE, Chao-Ho Chen, Tsong-Yi Chen, Cheng-Chan Cheng

Abstract-Nowadays people are used to sharing their views and ideas on social media service (SMS) platforms and generating enormous amounts of data every day. Sentiment analysis was adopted in this research to discover embedded information in Chinese short texts, serving as an integral part of Social Media Monitoring and Analytics. The research proposed a deep learning method. Attention-of-Emoticons Based Convolutional Neural Network (AEB-CNN), by integrating emoticons and attention-based mechanisms with CNN to enhance the accuracy of sentiment analysis without external knowledge. Implementation was performed by TensorFlow and the accuracy of sentiment polarity of Chinese microblogs reached as high as 85.8%.

Index Terms—Sentiment Analysis, Attention-Based, CNN, Emoticon.

I. INTRODUCTION

S ENTIMENT analysis (or Opinion Discovery), a key process of Social Media Monitoring and Analytics, detects the sentiment polarity of public opinions by collecting enormous amounts of microblogs from social media platforms. Most related researches adopted the natural language processing (NLP) technique, n-gram, sentiment dictionary, bag-of-mouth and so on, to parse and retrieve potentially useful information from collected short texts, but the accuracy of sentiment analysis has room for improvement.

Since Web 2.0 was published, social media services have attracted Internet users to leave their views and ideas on SMS platforms and interact with others. The interactions between users easily lead to concrete subgroups with high-density relationships between members.

Neural Network was first successfully applied to handwriting recognition by Hinton et al., who proposed Deep Belief Network (DBN) in 2006 [1]. In these years, with the rapid advances of computer technology, Deep Learning applications have been widely adopted in a great variety of research domains. For example, the hidden layer of Recurrent Neural Network (RNN) was proposed by Rumelhart et al. [2], in which the state store unit can be used to merge the vectors of the current term and its previous state to make prediction of the next term. RNN however has the weakness of vanishing gradients. Therefore, Hochreiter and Schmidhuber [3] proposed the Long Short-Term Memory Network (LSTM) to solve the disadvantage of RNN. Xinjie Zhou et al. [4] used LSTM to perform cross-language sentiment analysis without external knowledge in order to enhance accuracy. Yequan Wang et al. [5] combined the attention mechanism and LSTM to improve the accuracy rate of sentiment analysis.

Convolutional Neural Network (CNN) is the core technique of deep learning for image processing in these years. In 2012, Krizhevsky et al. [6] proposed AlexNet to lead a new trend of CNN. When CNN is processing the input data, the attention mechanism can assist CNN to focus only on important features while ignoring the rest. For example, Wenpeng Yin et al. [7] proposed combining single attention and CNN to solve the text entailment problem; Linlin Wang et al. [8] merged multi-attention mechanisms and CNN to classify sentence relations.

This research proposed an attention-based mechanism with emoticons that can significantly reduce the manually tagging time on each aspect word derived from the output of the segmentation process. In addition, the Convolutional Neural Network (CNN) was also adopted to capture the local features of each sentence, which makes it possible to pay attention on terms with higher implications of sentiment polarity and ignore the others for effective identification of sentiment. In the experiment, the proposed method for sentiment analysis could reach significantly higher accuracy.

II. RELATED WORK

A. Word Vector or Word Embedding

In the NLP domain, one-shot representation is one of simplest word vector methods, uses a long n-dimension vector to show a term. The length of a vector is the total number of terms in which 1 represents current term and others are 0. The drawback of this method easily has the dimensionality problem causing the size of space growing too fast and sparse data distribution. The second problem, each term is independent and hard to find the relation between two terms.

The representation of word vector maps a term to a real number vector with lower dimension. All vectors are forming a term-vector space. In the problem space, related terms or similar terms can be identified by Euclidean Distance or Cosine Similarity to conquer the weakness of one-hot representation

Yi-Jen Su is with the Department of Computer Science and Information Engineering, Shu-Te University, Kaohsiung 82445, Yanchao, Taiwan (e-mail: iansu@stu.edu.tw).

Chao-Ho Chen, Tsong-Yi Chen and Cheng-Chan Cheng are with the Department of Electronic Engineering, National Kaohsiung University of Science and Technology, Kaohsiung 80778, Sanmin, Taiwan (e-mail: thouho@kuas.edu.tw, chentso@kuas.edu.tw, 1105305117@gm.kuas.edu.tw).

and solve the problem of curse of dimensionality for promoting the accuracy of sentiment polarity judgement.

B. Convolutional Neural Network (CNN)

CNN is one of important neural network architectures in deep learning, especially in the application domain of image recognition. There are two major reasons causing RNN be replaced by CNN in the application domain NLP. First, CNN owns excellent feature retrieving function. Secondary, CNN is better on classification result and training speed. Kalchbrenner adopted CNN to capture the sentiment polarity of people opinions from movie comments and Twitter owned better accuracy [9]. Yoon [10] proposed the CNN-Text model with word-vector to solve the judgement of sentiment polarity, as shown in Figure 1. The architecture contains 5 layers, Input Layer, Convolutional Layer, Pooling Layer, Fully Connected Layer, and Output Layer.

Input Layer: The segmentation process for dividing input



Fig. 1. The architecture of CNN-text.

Chinese sentence to terms, each term will be transferred to a word vector to lower down the dimension of distributed space. A term, however, maps to a *k*-dimension vector $x_i \in \mathbb{R}^k$. Therefore, a sentence with fixed length *n* will be represented as the Equation (1):

$$x_{1:n} = x_1 \oplus x_2 \oplus x_3 \oplus \dots \oplus x_n \tag{1}$$

Convolutional Layer: Using different size of kernels, length is d and width is k, to proceed convolutional computing with the input layer for deriving the local features of input layer. The derived local features cu as shown in Equation (2):

$$c_u = f_c (W_c \cdot x_{uu+d-1} + b_c) \tag{2}$$

Where W_c is the weight of convolutional kernel weight, d stands for the maximum length of convolutional kernel, and b_c is the bias value, f_c stands for the Activation Function. $X_{u:u+d-1}$ shows the matrix of the kernel. After convolution computation, a sentence with length l is adopted to derive the feature map as shown in Equation (3):

$$c = [c_1, c_2, c_3, \dots, c_{n-d+1}]$$
(3)

Where *c*, the feature map, stands for a 1-dimension matrix with the length n-d+1.

Pooling Layer: In general, Max-over-time pooling is adopted to derived convolutional feature map by convolutional computation to output the feature with the maximum influence, as shown in Equation (4):

$$\overline{c} = \max(c) \tag{4}$$

Fully Connected Layer: After convolutional layer and pooling layer proceeded iteration learning, an input feature enters the fully connected layer will merge distributed feature together as the classification feature to be used for prediction.

C. Attention Mechanism

Attention mechanism was firstly proposed in the image recognition domain. When people watching a picture, their attention only focus on some interested areas. In 2014, the team of Google Mind [11] applied RNN with attention mechanism to classify images. Following the order of input images to observe the interested areas, based on the observation result to adjust the observation areas to prediction the final classification result.

Bahdanau et al. [12] was the first adopted attention mechanism on Neural Machine Translation (NMT). The proposed model is sequence to sequence to learn the target term and the attention weight of terms. There are three major Equations, (5), (6), and (7), are used to calculate the focus areas, as shown in followed.

$$f(x_i, x_k) = x_i^T \cdot x_k \tag{5}$$

$$f(x_i, x_k) = x_i^T \cdot W \cdot x_k \tag{6}$$

$$f(x_i, x_k) = W \cdot [x_i; x_k] \tag{7}$$

Where x_j and x_k is the input, Equation (5) shows the calculation of inner product. Equation (6) denotes to increase the Weight *W* by iteration learning to derive the value of the best attention weight. Equation (7) uses concatenation-based attention and iteration learning the weight *W* for deriving the best value of the attention weight.

III. RESEARCH METHOD

To promote the accuracy of sentiment analysis for Chinese microblogs, the research proposed the Attention of Emoticons Based-Convolutional Neural Network (AEB-CNN) model. The basic concept of AEB-CNN model is to point out the local feature of a sentence by CNN and to recognize the important part of a sentence with highly strength sentiment degree by merging emoticon and attention mechanism, as shown in Figure 2.







$$inp_{i,j} = f(x_{i,j}, aw_i) \tag{9}$$

$$Att_{i,j} = \frac{\exp(inp_{i,j})}{\sum_{j=1}^{l} \exp(inp_{i,j})}$$
(10)

Where x^{pad} is filled for the wild-convolution, $x_{i,j}$ and aw_i are adopted for the inner product $inp_{i,j}$. The Equation (10) shows the softmax algorithm to derive $Att_{i,j}$, represents the strength of attention. The example of attention mechanism computation as shown in Figure 4.

	今天	玩	的	很	開心	,	期待	下次	的	出遊
A	0.1221	0.4894	0.0019	0.0019	0.1252	0.0019	0.0554	0.0724	0.0019	0.0140
	2336	375	3216	3216	6338	3216	7103	6048	3216	3722
	早安	好友	們	一起	加油	pad	pad	pad	pad	pad
В	0.2400	0.1410	0.0032	0.1735	0.2376	0.0027	0.0027	0.0027	0.0027	0.0027
	2655	9854	4814	9838	3435	2524	2524	2524	2524	2524
	氣	到	心臟	好痛	,	到底	在	想	뇬	什麼
С	1.77865	3.61750	3.45407	2.76380	3.61750	3.61750	3.61750	8.10462	3.70973	3.61750
	74e-01	43e-05	12e-03	25e-04	43e-05	43e-05	43e-05	30e-01	86e-03	43e-05
	可愿	的	下雨天	破壞	我	的	計畫	pad	pad	pad
D	9.34169	2.29353	4.13832	4.25091	2.29353	2.29353	6.82664	1.73746	1.73746	1.73746
	11e-01	10e-04	14e-02	29e-03	10e-04	10e-04	43e-03	88e-04	88e-04	88e-04

Fig. 4. Examples of attention mechanism computation

The feature extraction process products the attention weight matrix and the convolution feature map matrix and uses the activation function to generate the hidden layer as shown in the Equation (11)

$$h_t = f_h(W_h \cdot (c_t \times Att_i) + b_h) \tag{11}$$

Where c_t stands for the convolution feature map matrix, and Att_i represents the attention weight matrix, are applied throug

Fig. 2. The architecture of AEB-CNN.

In data preprocessing process, all training sentences with emoticons were collected from Plurk, a famous Chinese SMS platform. A sentence with emoticon is easily to judge its sentiment polarity either positive or negative, for example :-D stands for positive and X-(stands for negative. To avoid the training result of word vector is not qualified, all collected sentence are removed no effective terms or symbols, e.g. URL, digit, nonsense symbol, and so on. In addition, the Chinese segmentation work for all collected sentences is adopted THULAC (THU Lexical Analyzer for Chinese) [16] to complete based on the consider of the segmentation time. Then, the Word2vec uses the segmentation result as the input to generate term vectors, as shown in Equation (8),

$$s_i = \{x_{i,1}, x_{i,2}, x_{i,3}, \dots, x_{i,n}\}$$
(8)

Where s_i denotes for the *i*-th sentence and $x_{i,j}$ stands for the *j*-th term of the *i*-th sentence.

In the attention weight computation process, the sentence s_i with attention emotion term aw_i are transferred to term vector. Both s_i and aw_i are used to proceed the inner product for generating the attention weights matrix, as shown in Figure 3 and Equation (9) and (10) weight of hidden layer W_h , bias of hidden layer b_h , and activation function of hidden layer f_h to generate hidden layer.

The research adopted max-over-time pooling method sampling as show in Equation (12)

$$p_t^h = \max(h_{t,1} \ h_{t,2} \ h_{t,3} \dots \ h_{t,p})$$
(12)

Where using length *h* convolution kernel to output the *t*-th feature map. From $h_{t,1}$ to $h_{t,p}$ proceeds the max-over-time pooling method sample to derive the fully connected layer matrix $_{p_{1}^{p_{1}}}$.

IV. EXPERIMENT RESULT

To prove the effectiveness of the AEB-CNN model on Chinese microblog sentiment analysis, the research randomly collected Chinese short text with emoticon from the SMS platform Plurk. All collected microblog are divided into two datasets, training set and testing dataset. In experiment, there are 10,000 sentences in the training set including 5,000 positive sentences and 5,000 negative sentences. The testing dataset, including 3,000 sentences, has two subgroups, positive and negative are the same 1,500 sentences.

Based on the research observation, most users of PTT, a famous BBS service platform, do not like to leave Chinese short texts with emoticons. Therefore, the research collected 9,607,232 sentences during 2018/03/01~2018/06/01 from PTT as the word vector training set. After the training, there are 122,218 terms are retrieved in the corpus.

As shown in Table 1, the experiment has 10,000 sentences in the training dataset and another 1,000 sentences in the testing dataset. To compare with other 5 previous research method, AEB-CNN owns the better accuracy of sentiment analysis.

As shown in Figure 5, to observe the change of accuracy in TABLE I

AC	ACCURACT COMPARISON OF SENTIMENT ANALYSIS							
Model	Accuracy	Precision	Recall	F1				
CCLM [13]	62.6%	97.8%	26.8%	42.1%				
JCCLM [14]	59.8%	97.3%	21.3%	48.8%				
SVM [15]	77.3%	83.0%	68.6%	75.1%				
LSTM [3]	76.8%	83.1%	67.2%	74.3%				
CNN [10]	80.5%	78.9%	83.2%	81.0%				
AEB-CNN	85.8%	81.7%	92.2%	86.6%				

different size of training dataset, there are five datasets, 10,000, 15,000, 20,000, 25,000, and 30,000 sentences. The proposed SEB-CNN model also has better accuracy performance.

V. CONCLUSION

The research proposed AEB-CNN model by merging emoticon and attention mechanism with CNN. Using emoticon in the calculation of attention weight, the important key terms



Fig. 5. The accuracy of 6 sentiment analysis methods.

will be set high attention, but ignore the terms with low emotion strength for promoting the accuracy of sentiment polarity of short texts either positive or negative. Especially, the proposed method does not need to manually tag aspect word or target word for saving lot of time. The experiment shows the proposed method owns high accuracy of sentiment polarity with Chinese microblog.

REFERENCES

- G. E. Hinton, S. Osindero and Y. W. Teh, "A fast learning algorithm for deep belief nets," *Neural Computation*, vol. 18, no. 7, pp. 1527-1554, 2006.
- [2] D. E. Rumelhart, G. E. Hinton and R. J. Williams, "Learning internal representations by error propagation," in *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, vol. 1, pp. 318-362, 1986.
- [3] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735-1780, 1997.
- [4] X. Zhou, X. Wan and J. Xiao, "Attention-based LSTM Network for Cross-Lingual Sentiment Classification," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2016.
- [5] Y. Wang, M. Huang, X. Zhu and L. Zhao, "Attention-based LSTM for Aspect-level Sentiment Classification," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2016.
- [6] A. Krizhevsky, I. Sutskever and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in *Proceedings of the 25th International Conference on Neural Information Processing Systems*, 2012.
- [7] W. Yin, H. Schutze, B. Xiang and B. Zhou, "ABCNN: Attention-Based Convolutional Neural Network for Modeling Sentence Pairs," *Transactions of the Association for Computational Linguistics*, vol. 4, pp. 259-272, 2016.
- [8] L. Wang, Z. Cao, G. Melo and Z. Liu, "Relation Classification via Multi-Level Attention CNNs," in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, 2016.
- [9] N. Kalchbrenner, E. Grefenstette and P. Blunsom, "A Convolutional Neural Network for Modelling Sentences," in *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, 2014.

- [11] V. Mnih, N. Heess, A. Graves and K. Kavukcuoglu, "Recurrent Models of Visual Attention," in *Proceedings of the 27th International Conference* on Neural Information Processing Systems, 2014.
- [12] D. Bahdanau, K. Cho and Y. Bengio, "Neural machine translation by jointly learning to align and translate," in *Proceedings of the 3rd International Conference on Learning Representations (ICLR-2015)*, 2015.
- [13] Y. J. Su, R. C. Chen, C. M. Hsiung, Y. Q. Chen, S. W. Yu and H. W. Huang, "Using Prefix Tree to improve the performance of Chinese Sentiment Analysis," in *Proceedings of the 30th IEEE International Conference on Advanced Information Networking and Applications* (AINA-2016), Crans-Montana, Switzerland, 2016.

Yi-Jen Su received his Ph.D. degree in electrical engineering from Cheng Kung University in Taiwan. He is an Associate Professor in Computer Science and Information Engineering at Shu-Te University. He has published more than 40 papers in journals and conference proceedings. His current research interests include social network analysis, sentiment analysis, data mining, artificial intelligence, e-Learning and image processing.

Chao-Ho (Thou-Ho) Chen received his Ph.D. degree in electrical engineering from National Taiwan University, Taiwan, in 1992. He has been the consultants of Eastern Graphics Co., HuperLab. Co. and Vivotek Inc.. From 2002, he worked at National Kaohsiung University of Applied Sciences and has been the chairman (also the founder) of Department of Computer Sciences and Information and the director of computer and network center. He obtained the Best Paper Awards of IIHMSP2007, NCWIA-2011, -2013, and -2014. His research interests mainly involve image/video processing and computer vision, and there are about 180 papers published and 27 ROC/USA patents.

Tsong-Yi Chen received his M.S. degree and Ph.D. degree in computer science from Illinois Institute of Technology, Chicago, U.S. A, in 1999. From 2000, he is an Assistant Professor in the Department of Electronic Engineering, National Kaohsiung University of Applied Sciences. His technical interests include digital image/video processing, computer vision and expert systems.

Cheng-Chan Cheng received his MS degree in Department of Electronic Engineering from National Kaohsiung University of Applied Sciences, Taiwan, in 2018. His recent research interests include image/video processing and deep learning, and sentiment analysis.

Visible Light Optics for Indoor Toxic Gas Detection and Positioning Systems

¹Shih-Hao Chang*, ²Chih-Chieh Hung, ³Mao-Sheng Hung

^{1,2,3}Department of Computer Science and Information Engineering, Tamkang University, New Taipei City, Taiwan 25137 shhchang@mail.tku.edu.tw, oshin@mail.tku.edu.tw, erichung1974@gmail.com

Abstract— LED light has many advantages over incandescent lighting such as power efficiency, improvements to quality, longer lifespan, and focus emission. LED lighting has been employed in our daily life in various ways, such as car lamps and lights in vehicles, roadside signal lights, indoor lighting systems, etc. These LED light bulbs can not only provide light but also can be employed for indoor positioning, namely visible light positioning (VLP), which is a fast-growing technology developed to provide necessary output such as data transmission and object orientation. Due to its attractive functionalities, it has drawn significant attention to the wireless communication and navigation systems that also reduce cost after the system employed. VLP can utilize received signal strength (RSS) from the use of the signal from the receiver with respect to the strength of the relationship between the distance attenuation. This approach also can be applied to create a measurement of optical power absorption and scattering across an air gap that is to be exposed by toxic gas or smoke. The paper presented focuses on a model-based understanding of particulates scattering established. Subsequent studies analyzed the impact of visible light scattering on toxic gas and the potential role of dust accumulation. This highlights some of the challenges associated with visible light scattering measurements of particulates in the presence of real-world surfaces.

Keywords—Visible Light Positioning, Toxic Gas Detection, Model-Based, Received Signal Strength, Light Emitting Diodes

I. INTRODUCTION

In last few years, air pollution due to noxious gases, such as volatile organic compounds (VOCs), CO, SO2, H2S, NH3, and NO2, has become a dangerous problem that causes harmful effects on plants, aquatic animals, and human health [1]. Among others, NO2 is one of the most highly toxic gases. Inhalation of low-level concentrations of NO2 (50 ppm) can cause damage to the lungs, cardiovascular system, and upper respiratory tract in humans. Therefore, the Occupational Safety and Health Administration (US) announced that the permissible exposure for general industries is 5 ppm, and 1 ppm for short-term exposure (15 min). The development of gas sensors for monitoring NO2 is highly important and necessary to protect people from over exposure to such dangerous gases and to improve environmental quality. In the event of an indoor air pollution emergency, people inside the building should be evacuated from the structure as soon as possible. However, it can be very difficult to organize and execute a quick evacuation procedure due to the complexity of modern buildings and the great numbers of people that can be inside. Critical problems such as airlessness, huddling, smothering, trampling and inaccessibility of exits may arise during the evacuation procedure. The effects of deposition, cracks and insulation in ventilation channels is also treated, as well as the effects of internal filtration. The values of available indoor deposition velocities for some toxic gases are reviewed.

Typical toxic gas sampling methods employ manual grab samples that are collected on site and then transported to a laboratory for analysis. These sampling methods can be very costly, time-consuming and can compromise the integrity of the sample during sample collection, transport, storage and analysis. Portable, robust, accurate methods of analysis are needed to achieve monitoring such that the samples can be analyzed in the field. These enable results to be available faster, at low cost and they minimize the risk of contamination by eliminating the transportation of the samples [2]. Using LED lighting has many advantages such as power efficiency, better quality, longer lifespan, and focus emission. LED lighting has been employed in our daily life in various ways, (car lamps and lights in vehicles, indoor lighting systems, etc.) [6-12]. The advancement in LED sources and photo-detector technologies provide a solution to these issues as composed by compact, low power and low-cost detectors designed for incorporated colorimetric analytical methods.

These analytical methods can be used to study toxic gas from accidents or attacks with chemical warfare agents. The influence of indoor pollutant location, total exposure, and maximum concentration is studied. Visible Light Positioning (VLP) [3-5] is a promising complementary and/or alternative technology to its radio frequency (RF) counterparts in indoor environments. In contrast, Global Positioning System (GPS), which can barely penetrate building walls and tunnels; cannot be used efficiently for indoor areas or inside tunnels. VLP is now a fast-growing technology which provides indoor positioning using low-cost and energy efficient light-emitting diodes (LED). The visible light source can be modulated at a high-speed data rate; therefore, it's able to efficiently transmit positioning data. Since VLP technology has been proved as a pertinent positioning system, we can apply a certain computational algorithm in the literation to obtain an accurate

position of the receiver. As described in [3-5], VLP technology can utilize receiver signal strength (RSS) [13] and angle of arrival (AOA) [14], or time-difference-of-arrival (TODA) [15] to estimate location information.

In the event of an indoor air pollution emergency, people inside the building should be evacuated from the structure as soon as possible. The problem of evacuating buildings through the shortest path safely during a case of indoor toxic gas and air pollution incident has become more important than ever due to the nature of today's complex and elaborate building schemes. To address this problem, we propose an indoor visible light toxic air detection and positioning system. This system will model toxic gas positioning by applying received signal strength (RSS) [16]. This RSS also can be used as a measure of visible light absorption and scattering power across an air gap that is exposed by smoke or toxic gas. This paper is organized as followed: In Section II, background and related works are briefly described. In Section III, the RSS approach that are applied to the toxic gas-positioning algorithm has been modeled. The emulation has been presented in section IV. And finally, our conclusions and future work are explored in Section V.

II. BACKGROUND AND RELATED WORKS

Optical remote sensing technologies are expected to play a role, possibly a significant one, in the development of residual risk requirements in facility emissions compliance. However, a more immediate and urgent application for optical remote sensing technologies lies in the requirements for the detection of accidental releases. This discussion is limited to the capabilities of electro-optical technology to measure toxic gases and vapors using open-path methods. (The gases covered under previous legislation— sulfur dioxide (SO2), carbon monoxide (CO), nitrogen dioxide (NO2), and ozone (O3)-are mentioned only briefly in this article.) General reviews of electro-optical remote sensing instruments are found in a number of references [1]. Some techniques, such as spectral mask correlation spectroscopy, will not be discussed here. Mask correlation spectrometers are generally limited to the ultraviolet (UV) and visible spectral regions where sunlight is available and to gases, such as SO2 and NO2, that have spectra with periodic repetition of absorption features, and do not seem to apply to any of the new toxic gases [2].

On the other hand, with regard to indoor positioning, radiolocation used in the past generally use a fixed position to know the source of the signal development, depending on the intensity of the received signals [3, 4], azimuth and transmission time, supplemented by algorithms calculate the target's indoor location. But because variation in radio waves are due to changes in the indoor environment, arising masking signal or signal drift problem, consider locating in a visible way. The positioning calculation method currently used basically can be summarized into the following categories: Moreover, Visible Light Positioning [5] can be highly directional, thus it is difficult

to re-establish a link that has been lost due to movement or rotation of one of the devices in the link. Therefore, the proposed VLP link protection mechanism does not influence the design or functionality of the presented solution. Obviously, real-time and reliability adaptation can be, and have to be, implemented at each system level.

Arrival time method (Time of Arrival, TOA)--a method using three or more of the known locations from the signal source, the emission of the time information with the same content--spend more time after the signal reaches the receiver to receive all the differences [13]. The Cramer-Rao Bound (CRB) is a widely used approach to evaluate the theoretical accuracy in estimation theory in which a lower bound on the mean square error of an unbiased estimate can calculate the relative distance between a known points, then locate and triangulate its position. However, this method requires the use of a high-precision signal source calculator to achieve precise positioning.

Signal strength method (Received Signal Strength, RSS): RSS approach is the measurement of the signal from the receiver with respect to the strength of the relationship between the distance attenuation, with three or more known fixed positions using source signal strength, measured before the construction of the transmission loss pattern, and then again after the comparison operation can be learned from each other the relative distance. With the three sets of data, one is able to draw three circumferential distances and the focus can be positioned, i.e. three circumferential own self-position [14].

Recipient angle method (Angle of Arrival, AOA): AOA uses a directional antenna connected to the source of wave reception signals to the signal to determine the angle of arrival of the signal source position, and therefore at least two or more antennas receiving this signal. The incidental angle is extended to the intersection points to obtain position signal source. This method must use directional antennae, but manipulation of the surrounding environment is easy using reflection multipath propagation problems [15]. However, these methodologies only consider normalized received signal without people and goods moving through in indoor environment. Therefore, we propose a positioning control system and visible signal strength detection module which will send a specific visible identification, coded in a LED lighting drive circuit. When the illuminating light is emitted, the signal through silicon photodiode becomes the receiving analog signal, and then the analog-digital processing of legislation can be used by software to interpret the information content received.

III. PROBLEM FORMULATION

RSS approach uses signal from the receiver with respect to the strength of the relationship between the distance attenuation using three or more known fixed position signals source signal strength. The signal is recorded once before the construction of the transmission loss patterns are recorded, and then again after the comparison operation can be learned from each other the relative distance. The three sets of data are used in order to be able to draw three circumferential distances and the focus can be positioned, i.e. three circumferential own self-position [14]. In a static optical channel the total received power Pr at the receiver is given by:

$$P_r = P_t H_d(0) + \int P_t H_{\text{ref}}(0) \tag{1}$$

where Hd (0) and Href(0) are the DC channel gain of the direct and reflected paths, respectively and Pt is total optical power transmitted by LEDs. The DC channel gain of the direct path is given as:

$$H_{d}(0) = \begin{cases} \frac{(m+1)}{2\pi d^{2}} A_{\det} \cos^{m}(\theta_{r}) T_{s}(\vartheta) g(\vartheta) \cos(\vartheta), \\ \text{for } 0 \leq \vartheta \leq \vartheta_{\text{FOV}} \\ 0, \quad \text{for } \vartheta > \vartheta_{\text{FOV}} \end{cases}$$
(2)

where m is the order of the Lambertian radiant, d is the distance between transmitter and receiver, θ r is the angle of irradiance, ϑ is the angle of incidence, Ts (ϑ) is the optical filter gain, g(ϑ) is the optical concentrator gain and ϑ FOV is the FOV of the receiver. The DC gain of the reflected path can be written as:

$$H_{d}(0) = \begin{cases} \frac{(m+1)}{2(\pi d_{1}d_{2})^{2}} \rho A_{\det} dA_{w} \cos^{m}(\theta_{r}) \cos(\alpha) \times \\ \cos(\beta) T_{s}(\vartheta) g(\vartheta) \cos(\vartheta), \\ \text{for } 0 \leq \vartheta \leq \vartheta_{\text{FOV}} \\ 0, \quad \text{for } \vartheta > \vartheta_{\text{FOV}} \end{cases}$$
(3)

where d1 is the distance between transmitter and reflective point, d2 is the distance between the reflected point and receiver, ρ is the reflectance coefficient, dAw is a small reflective area on the wall, α is the angle of incidence from the transmitter and β is the angle of irradiance from a reflected point. The multipath character of the optical link can be described as the delay profile of the channel [4]. The delay spread estimates the upper bound of the maximum channel capacity without the need for equalization [5]. To avoid very complicated calculations of electronic transitions, numerous measurements of the absorption cross-sections of the atmospheric atoms and molecules absorbing in the visible light have been performed. In general, the absorption cross section varies with temperature.

IV. SIGNAL STRENGTH METHOD (RECEIVED SIGNAL STRENGTH)

In order to obtain a visible light positioning accuracy algorithm (although most visible positioning algorithms defect in the article and the positioning algorithm in the channel model) the emphasis is on certainty of the distance, different angles and different algorithms that produce different signal strength. However, this does not consider the different environments and physical obstructions caused by shadows, so practicality is limited. In order to improve this situation, it is proposed that a positioning algorithm and channel model, this algorithm will assist us in limited circumstances continue to make the most precise positioning. Its method is as follows: We can assume that under the transfer position, the three transmission ends T0, T1 and T2. It is at the position (X0,Y0,Z0)t0(X1,Y1,Z1)X2,Y2,Z2)t, the square of the distance between the pair and T0 and T1, called , is used in the formation of a field to shows in terms of equal and fields as follows:

$$d_{0,1}^{2} = (X_{0} - X_{1})^{2} + (Y_{0} - Y_{1})^{2} + (Z_{0} - Z_{1})^{2}$$
(4)

Suppose the measurement point in the entire environment is the location coordinates K0, K1, K2. In accordance with these variables, the measuring point receiver reference point to calculate where. However, the measured value it receives signals because of the size and shape of the obstacle in the lighting configuration which create problems in the lighting configuration. [8] In the same plane, according to the signal value of the energy, we can expect to produce different estimates, as follows:

$$F1 = \sum_{i=1}^{3} |EIPiRF - PiRF|$$
(5)

Assuming F1 is in the same plane and has the ability to form a triangular region using light irradiation from three sources, each class and each of the receiver estimates the actual calculated values, minus the value produced from the received difference value calculating accurate. However, in order to measure the signal value, we used a mathematical model to calculate the probability density function (PDF) of the measured values. [9] Set X as a measure of value, σ as the variance in the following:

$$f(x; \sigma) = \frac{x}{\sigma} e^{-x^2/2\sigma^2}, x \ge 0$$
 (6)

When the two components of a two-dimensional random vector are independent, they have the same normal variance through this algorithm to prove our accuracy. Based on the RMS delay spread channel model, we can estimate the range of LED lamps by using measured optical power at the receivers, which is referred to as received signal strength (RSS). Therefore, we can discuss the RSS approaches in more detail where usually the LS method is applied to obtain the position of the receiver in the literature. On top of the LS algorithm, we can develop an approach by applying an ML estimator to highly enhance the accuracy and performance of the positioning systems in this article. Besides the described approaches, a grid method can be employed where reference points are densely specified for a small area. If there is strong confidence that a targeted receiver is within the area, the grid method can be a

good choice for the problem; the accuracy and performance depends on the resolution of the grid.

The fingerprinting method is similar to the grid method. In the Wi-Fi fingerprinting technique, a radio map of an area of interest is created based on RSS from several access points, and generates a probability distribution of RSS values for any location within the area. This approach can be adopted in LED-based VLC positioning systems, and the accuracy and performance can be enhanced by taking advantage of VLC technology. Nonetheless, these two methods may not be of great interest due to limited applications in a relatively small area. Due to visible light, we will depend on the measurement of optical power absorption and scattering across an air gap that is to be exposed by smoke or toxic gas. As light propagates through a heterogeneous medium such as smoke aerosol in air, it is both absorbed and scattered. Rays that enter a material with non-zero absorption coefficient are attenuated according to Beers'-Lambert's Law of transmission:

$$P_t = P_i e^{-\alpha l} \text{ or } P_a = P_i (1 - e^{-\alpha l})_{(7)}$$

where *Pa*, *Pt* and *Pi* are absorbed, transmitted and initial optical flux, respectively, *l* represents an optical path length through the material, and \Box is the absorption coefficient related to the imaginary part of the material's refractive index, *k*:

$$\alpha = \frac{4\pi k}{\lambda}$$
 and $m = n + ik$

Here, *m* is the material's complex refractive index, while λ represents wavelength. However, this requires that scattering is negligible compared to absorption. When both absorption and scattering contribute to total extinction of a propagating light beam, Bougher's modified law should be used instead.

$$P_t = P_i e^{-\tau_{ext} l}$$
 and $\tau_{ext} = N(\sigma_{abs} + \sigma_s)_{(9)}$

Here, τ ext is the extinction coefficient combining absorption and scattering through the corresponding cross-sections, σ abs and σ s are expressed as 1/cm, while N represents the number of aerosol particles per unit volume A measurement of beam attenuation within a smoke test box or smoke room can provide the experimental τ ext value. However, in order to quantify light attenuation due to absorption and scattering, the appropriate cross-sections, σ abs and σ s need to be supplied as well. Typically, these two parameters are lumped together and reported as mass a scattering cross-section. They can be found in the literature as several common test fires and nuisance aerosols. Note that the measurement of aerosol concentration is still required before an experimental τ ext value may be used within the ray-tracing software to estimate Pt. The correlation of results from the model and experiments requires accurate determination of particle number density, size distribution and

fire type during typical test fires, such as a cotton wick smoke box test.

Most recently proposed approaches based on RSS methods employ the LS method to compute the position of the receiver by using estimated distances. Although the LS method obtains a solution easily, optimality cannot be guaranteed. With the assumption of Gaussian distributed measurement noise, the solution that gives zero to the first-order derivative of the likelihood function based on the measurement can be obtained iteratively. We can then obtain an IMLE estimate from the iterative solution.



Fig. 2 – Iterative Maximum Likelihood Estimator

The iterative method (e.g., the Newton-Raphson method) is needed when the solution cannot be obtained in a closed form. The iteration step finds the estimate that approaches the solution; however, beyond a certain number of iterations, it may converge to a value. The initial guess has high impact on the accuracy of this iterative method. We propose an IMLE approach to enhance the accuracy and performance of the positioning system beyond the LS method with the Gaussian assumption for the measurement noise. We adopt the LS solution (i.e., LS) as the initial value in this approach. The diagram for the iterative IMLE approach is shown in Figure 2, where details of the steps for the approach can be seen.

V. EVALUATION

To cover any area of a museum environment, multiple luminaries can be installed. Figure 3 shows a simulation map where we placed different quantities of luminaries in different 8 showrooms. The luminaries are all identical to the single luminary simulation and their locations are indicated in figure 3. The same pattern can be seen around the outside of the room except for the middle lobby, where there appears to be many more regions than expected. The additional regions are from the overlap between adjacent luminaries, which can also be used for positioning.

The simulations on going, when the simulation completed, we will supplementation in the paper.

VI. CONCLUSION

In recent years, indoor positioning has drawn significant research attention. Normal wireless communication has many well-known delivery problems, such as signal fading, multipath propagation, signal obscured and interference problems which will affect the indoor positioning accuracy. Visible light positioning (VLP) utilizing received signal strength (RSS) approach is the use of the signal from the receiver with respect to the strength of the relationship between the distance attenuation. This approach also can be applied to the measurement of optical power absorption and scattering or light across an air gap that is exposed by toxic gas or smoke. The presented paper focuses on a model-based understanding of particulates scattering established. Subsequent studies analyzed the impact of visible light scattering on toxic gas and the potential role of dust accumulation. However, future challenges will also need to be highlighted that associate visible light scattering measurements of particulates in the presence of real-world surfaces.

REFERENCES

- Finlayson-Pitts, B. J.; Pitts, J. N. "Monitoring Techniques for Gaseous Criteria and Non-Criteria Pollutants," in Atmospheric Chemistry: Fundamentals and Experimental Techniques, J. Wiley & Sons, NY, pp 305-379, 1986.
- [2] Hall, F. F., Ed., "Remote Sensing of Hydrocarbons and Toxic Pollutants (Workshop Minutes)," U.S. EPA Report No. EPA 600/9-90/009, EMSL, Office of R & D, U.S. EPA, Las Vegas, NV, 1990.
- [3] Zahid Farid, Rosdiadee Nordin, and Mahamod Ismail, "Recent Advances in Wireless Indoor Localization Techniques and System," Journal of Computer Networks and Communications, vol. 2013, pp.12, 2013.
- [4] R. Zurawski, Ed. Boca Raton , " Indoor tutoring communication systems," in The Indoor tutoring Information Technology Handbook, CRC Press, 2005, Sec. 3, pp. 37.1–47.16.
- [5] J. R. Moyne and D. M. Tilbury, "The emergence of indoor tutoring control networks for manufacturing control diagnostics and safety data," In Proc. of IEEE, vol. 95, no. 1, pp. 29–47, Jan. 2007.
- [6] N. Farr, A. Bowen, J. Ware, C. Pontbriand, M. Tivey, An integrated, underwater optical/acoustic communications system,in: Proceedings of IEEE OCEANS, 2010, 1-6.
- [7] Xiao-Wei Ng, Wan-Young Chung, VLC-based medical healthcare information system, Biomed. Eng.: Appl. Basis Commun. 24 (2) (2012) 155–163.
- [8] R. Murai, T. Sakai, H. Kawano, Y. Matsukawa, Y. Honda, K. Campbell, A novel visible light communication system

for enhanced control of autonomous delivery robots in a hospital, in: Proceedings of the IEEE/SICE International Symposium on System Integration (SII), 2012, pp 510-516.

- [9] S.-B. Park, D. K. Jung, H.S. Shin, D.J. Shin, Y.-J. Hyun, K. Lee and Y.J. Oh, Information broadcasting system based on visible light signboard, Presented at Wireless and Optical Communications 2007, Montreal, Canada, 2007.
- [10] http://www.vlcc.net/?ml_lang=en (27.06.15)
- [11] Y. Wang, N. Chi, Y. Wang, L. Tao, J. Shi, Network architecture of a high-speed visible light communication local area network, IEEE Photonics Technol. Lett. 27 (2) (2015) 197–200.
- [12] (https://mentor.ieee.org/802.15/dcn/08/15-08-0171-00-0vl c-10mbps-visiblelight-transmission-system.pdf),
- [13] Aomumpai, Supattra, et al. "Optimal Placement of Reference Nodes for Wireless Indoor Positioning Systems." Electrical Engineering Electronics, Computer, in Proc. of 11th IEEE Telecommunications and Information Technology (ECTI-CON) International Conference.
- [14] Ganti, Divya, Weizhi Zhang, and Mohsen Kavehrad. "VLC-based indoor positioning system with tracking capability using Kalman and particle filters.", in Proc. 2014 IEEE Consumer Electronics (ICCE) International Conference.
- [15] Hyun-Seung Kim; Deok-Rae Kim; Se-Hoon Yang; Yong-Hwan Son; Sang-Kook Han, "An Indoor Visible Light Communication Positioning System Using an RF Carrier Allocation Technique," Proc. in Journal of Lightwave Technology, vol.31, no.1, 2013, pp.134,144.
- [16] T. J. Fagan, "Letter: nomogram for bayes theorem," The New England journal of medicine, vol. 293, no. 5, pp. 257–257, 1975.
- [17] Zvanovec, S., Haigh, P.A., Ghassemlooy, Z, "Channel Characteristics of Visible Light Communication Within Dynamic Indoor Environment.", Journal of Lightwave Technology, Vol.33, issue. 9, pp.1719 – 1725, March , 2015
- [18] IEEE 802.15 WPAN[™] Task Group 7 (TG7) Visible Light Communication http://www.ieee802.org/15/pub/TG7.html.
- [19] J. M. Kahn and J. R. Barry, "Wireless infrared communications," Proc. IEEE, vol. 85, no. 2, pp. 265–298, Feb. 1997.
- [20] Z. Ghassemlooy, W. Popoola, and S. Rajbhandari, Optical Wireless Communications. Boca Raton, FL, USA: Taylor & Francis, 2012.
- [21] A. Burton, H. Le-Minh, Z. Ghassemlooy, S. Rajbhandari, and P. A. Haigh, "Smart receiver for visible light communications: Design and analysis," in Proc. 8th Int. Symp. Commun. Syst., Netw. Digit. Signal Process., Jul. 18–20, 2012, pp. 1–5

Detection of Intentional and Unintentional Financial Restatements using Data Mining Techniques

Tengfei Qian and Ou Liu

Abstract— this study proposed an idea of multi-class classifier to predict intentional fraudulent financial restatement, unintentional financial restatement and normal financial statement. Most prior studies on detection of financial fraudulent restatement use balanced dataset. They only focused on the fraudulent case and ignored the unintentional case. This study employs a large imbalanced dataset that includes 70781 financial statement from 5962 companies. It includes 596 fraudulent cases, 14751 unintentional cases and 55434 normal cases. In this study, dataset is preprocessed and different feature selection algorithms, resample algorithms and different data mining techniques will be applied in the future study.

Index Terms-Financial restatement, multi-class classifiers, large imbalanced dataset.

I. INTRODUCTION

FINANCIAL restatement, which can erodes investors' confidence on the companies, has receiving increased attention from both companies and investors. Financial restatements happen when firms make errors on their financial statements. There are two types of financial restatement: unintentional restatement and intentional restatement. The unintentional restatement is caused by the unintentional errors in financial statements. The intentional fraudulent restatement is caused by the intentional error like altering financial data to misguide market participants. Most prior researches focused on predictive model of fraudulent financial restatement and ignored the unintentional restatement. There is few research concentrate on unintentional restatement [3,7]. Dutta et al. [7] developed a predictive model to classify restatement (both fraudulent and unintentional) and normal statement. To the best of our knowledge, there is only one literature aims to classify fraudulent restatement, unintentional restatement and normal statement. In this research, they implemented a three-class financial statement fraud detection model which detects intentional misstatement, unintentional misstatement and non-fraud statement [3]. It used Accounting and Auditing Enforcement Releases (AAER) database. The restatement data from this database is reliable as all the data is investigated by Security Exchange Commission (SEC). However, this leads to an issue that many restatement data are excluded in this database as SEC review one-third of public companies' financial statements per year [11]. In our research, we uses Audit Analytics database, a commercial database that collects all the financial restatements from 1994 to current period. We use Audit Analytics database to get all the financial restatement data from 1994 to 2018. Meanwhile, we use Compustat database to get normal statement. The number of intentional fraudulent restatement is relatively small, compared to the number of unintentional restatement and the number of normal statement. The total number of statement in the dataset of this study employed is 70781. The of intentional fraudulent percentage restatement. unintentional restatement and normal statement in the dataset are 0.8% (596), 20.84% (14751) and 78.36% (55434) respectively. As most prior studies focus on predicting fraudulent restatement though using balanced data (see table 1), this study contributes to i) develop a three-class classifier to predict intentional fraudulent restatement, unintentional restatement and normal statement; ii) employ extremely imbalanced dataset to develop the predictive model with different data mining technique.

II. LITERATURE REVIEW

Financial restatement includes fraudulent restatements and unintentional restatements. A fraudulent restatement is revising a firm's previous fraud financial statement, while unintentional restatement is revising firm's previous financial statements where occurs material errors. The majority of data mining predictive models of detecting financial statement fraud concentrated only on the fraudulent restatements but ignored the intentional restatements.

There has been a widely use of machine learning technologies applied in detecting financial statement fraud. However, there is no consistent conclusion on the best performance of predictive models. The previous study shows that different data structure and time period affect the performance of detecting models. Logistic regression, SVM, artificial neural network and decision tree are algorithms that are most commonly used in detection of financial statement fraud. A summary of related financial statement fraud detection research is listed in table 1.

Ravisankar et al.[1] uses six data mining techniques Probabilistic Neural Network, Logistic Regression,

This work is partially supported by the Grant Subsidy (XJTLU RIBDA2018-IRP1) provided by the Research Institute of Big Data Analytics (RIBDA).

Tengfei Qian is with the International Business School Suzhou, Xi'an Jiaotong-Liverpool China University. Suzhou (e-mail: Tengfei.Qian@xjtlu.edu.cn).

Ou Liu is with the International Business School Suzhou, Research Institute of Big Data Analytics, Xi'an Jiaotong-Liverpool University, Suzhou, China. (corresponding author, phone:+86 512 8188 3264 e-mail: Owen.Liu@xjtlu.edu.cn).

Multilayer Feed Forward Neural Network, Support Vector Machines and Genetic Programming to predict financial statement fraud. The authors concluded that PNN has the best performance in detecting financial statement fraud without feature selection and GP and PNN have the best performance if feature selection is applied. The dataset of this research includes 101 fraudulent Chinese listed companies and 101 non-fraudulent Chinese companies.

Liu et al.[5] applied random forest, Logistic regression, Knearest Neighborhood, Decision Tree and Support Vector Machines with feature selection in financial statement fraud with Chinese listed company data. The authors found that the ratio of debt to equity is the most import feature in detection financial statement fraud. RF and SVM obtain much higher accuracy than the other three other algorithms LR, KNN and DT. It suggested that random forest shows the highest accuracy 88% and SVM got 80.18. The accuracy of LR is only 42.91%, the KNN is only 60.11% and DT is 66.43%.

Yeonkook et al. [3] divided financial misstatement into two types that are errors (unintentional) and irregularities (intentional). They implement a three-class financial statement fraud detection model which detects intentional fraud, unintentional misstatement and non-fraud statement. Hennes et al.[4] suggested that it is more important to distinguish intentional fraud than simply distinguish fraud. Intentional financial fraud causes series consequence like shaking the investor's confidence. Intentional financial fraud will lead to more frequent securities class-action lawsuits follow than unintentional misstatement.

Throckmorton et al.[6] used financial statement data as nonverbal features and public company quarterly earnings conference call audio record as verbal features to detect financial statement fraud. In this study, Bayesian-based GLRT is used as classification algorithm. The number of nonfraud samples is 1531 which is significantly larger than the size of fraud samples which is only 41. The authors concluded that verbal features provide complementary information to detect financial statement fraud.

Dutta et al. [7] used data mining technique to develop a predictive model to determine financial restatements with 3513 financial restatement data and 60720 normal statements over a period from 2001 to 2014. Five classifiers are evaluated (ANN, DT, NB, SVM, and BBN) and SMOTE is used to address the issue of class imbalance and cost imbalance. ANN and DT obtain very close high accuracy which is 77.7% and 79.9% respectively.

Zhou et al. [2] considered that executives who commit financial statement fraud have knowledge on the fraud detection mechanism or application, which are usually not difficult to obtain. Executives try to adapt the mechanism while they commit fraud. In such case it is difficult to detect the same mechanism. Therefore, Zhou et al. proposed an adaptive financial statement fraud detection framework to detect evolutionary financial statement fraud. This framework attempts to form domain knowledge on fraud statement detection and analyzes fraud statement data based on this domain knowledge. However, Zhou et al. only proposed such model and had not yet developed and evaluated such model. Therefore, the performance of this framework is still unknown.

Lin et al. [8] used stepwise regression to rank the importance of financial features used in the data mining models (Logistic Regression, Decision Trees and Artificial Neural Networks). The ranking results are compared with the ranking that obtained from an experts' questionnaire survey, which evaluates the fraud financial features by using Lawshe's approach [9]. The research shows that the empirical results are very consistent with the decisions of the experts.

Research	Number of feature	Data set	Algorithm	Accuracy
Ravisankar (2010)	18 financial measures	101 fraud samples and 101 non-fraud	Multi-layer feedforward neural network	75.32%
		samples	Support Vector Machines	72.36%
			Genetic programming	89.27%
			Group method of data handling	88.14%
			Logistic regression	70.86%
			Probabilistic neural network	90.77%
Huang(2014)	9 financial measures	72 fraud samples and 72 non-fraud samples	Dual Dynamic multilayer hierarchical network structures	82.47%
Lin(2015)	9 financial measures	129 fraud samples and 447 non-fraud	Logistic Regression	88.5%
	8 financial measures	samples	Decision Trees	90.3%
	32 financial measures		Artificial Neural Networks	92.8%
Throckmorton (2015)	4 financial measures + 14 Acoustic and Linguistic measures	41 fraud samples and 1531 non- fraud samples	Bayesian-based GLRT	81%

 TABLE I

 PREVIOUS RESEARCHES ON DETECTION OF FRAUDULENT FINANCIAL STATEMENT

Liu (2015)	8 financial measures	138 fraud samples and 160 non-fraud samples	Logistic Regression K-Nearest Neighbor Decision Tree	42.91% 60.11% 66.43%
			Support Vector machines	80.18%
			Random forest	88%
Yeonkook (2016)	49 financial measures	355 error samples, 214 fraud sample and 2156 non-	Multinomial logistic regression	86.9%
	fraud samples		Support vector machines	85.4%
			Bayesian networks	82.5%
Dutta (2017)	15 financial	3513 restatement	Decision tree	79.9%
	measures	and 60720 normal statement	Artificial Neural Network	77.7%
	Satement		Naïve Byes	60.1%
			Support Vector Machines	73.4%
			Bayesian Belief network	75.1%

III. DATA PRE-PROCESSING

The dataset of this study is obtained from two databases: Audit Analytics and Compustat. We firstly collected financial restatement data from Audit Analytics. The data range of financial restatements are from 1994 to 2018. There are 17532 financial restatements from 10091 companies. 106297 Financial statements data is then obtained from Compustat database queried by company codes from Audit Analytics. With the fraud records in Audit Analytics, we then classified these 106297 financial statements into 733 intentional fraudulent restatements, 21438 unintentional restatements and 84126 normal statements. After classification, we calculated 9 financial statement features based on the research of Perols [12]. The features are listed in table 2. We also removed null value records in the dataset in order to get meaningful data of features. In addition, we used data winsorizing method to remove the outliers. Finally the dataset contains 596 intentional fraudulent restatements, 14751 unintentional restatements and 55434 normal statements from 5962 companies.

	T.	AE	3L	E	Π		
 				_		 	_

NINE FINA	NCIAL FEATURES
rect_to_sale	Accounts Receivable to Sales
rect_to_ta	Accounts Receivable to Total
	Assets
ppegt_at	Fixed Assets to Total Assets
gross_margin	Gross Margin
invt_sale	Inventory to Sales
ppent_at	Property Plant and Equipment to
· · ·	Total Assets
sale_at	Sales to Total Assets
total_accruals_to_at	Total Accruals to Total Assets
lt_to_at	Total Debt to Total Assets

IV. EXPECTED RESEARCH OUTCOMES

This research expects the following outcomes:

1. Overcome the highly imbalanced dataset with resample algorithms like ENN,SMOTE and NCL. Evaluate these

resample algorithms to find the most suitable algorithm for this highly imbalanced financial statement dataset.

- 2. Add more financial features into the dataset and use feature selection algorithms to rank the importance of the financial features to find most related financial features to train the predictive model.
- 3. Try different machine learning algorithms like Support Vector Machines, Random Forest and Artificial Neural Network. Evaluate the performance of these algorithms to find the most efficient and accurate predictive model.

V.CONCLUSION

We proposed an idea of developing a multi-class classifier to predict intentional fraudulent financial restatement. unintentional financial restatement and normal financial statement. We employed a very large dataset which includes 70781 statements from 5962 companies. In the future study, we will resample the dataset using algorithms like ENN.SMOTE and NCL to overcome data imbalance issue. More financial features will be added to the dataset and feature selection, and feature selection algorithms will be used to rank the importance the financial features. We will also test different data mining techniques on this dataset and find out the bst performance of data mining technique to predict both intentional fraudulent restatement and unintentional restatement.

REFERENCES

- Ravisankar, Pediredla, et al. "Detection of financial statement fraud and feature selection using data mining techniques." Decision Support Systems 50.2 (2011): 491-500.
- [2] Zhou, Wei, and Gaurav Kapoor. "Detecting evolutionary financial statement fraud." Decision Support Systems 50.3 (2011): 570-575.
- [3] Kim, Yeonkook J., Bok Baik, and Sungzoon Cho. "Detecting financial misstatements with fraud intention using multi-class cost-sensitive learning." Expert Systems with Applications 62 (2016): 32-43.
- [4] Hennes, Karen M., Andrew J. Leone, and Brian P.

Miller. "The importance of distinguishing errors from irregularities in restatement research: The case of restatements and CEO/CFO turnover." The Accounting Review 83.6 (2008): 1487-1519.

- [5] Liu, Chengwei, et al. "Financial fraud detection model: based on random forest." (2015).
- [6] Throckmorton, Chandra S., et al. "Financial fraud detection using vocal, linguistic and financial cues." Decision Support Systems 74 (2015): 78-87.
- [7] Dutta, Ila, Shantanu Dutta, and Bijan Raahemi.
 "Detecting financial restatements using data mining techniques." Expert Systems with Applications 90 (2017): 374-393.
- [8] Lin, Chi-Chen, et al. "Detecting the financial statement fraud: The analysis of the differences between data mining techniques and experts' judgments." Knowledge-Based Systems 89 (2015): 459-470.
- [9] Huang, Shaio Yan, et al. "Fraud detection using fraud triangle risk factors." Information Systems Frontiers 19.6 (2017): 1343-1356.
- [10] Huang, Shin-Ying, Rua-Huan Tsaih, and Fang Yu. "Topological pattern discovery and feature extraction for fraudulent financial reporting." Expert Systems with Applications 41.9 (2014): 4360-4372.
- [11] Dechow, Patricia M., et al. "Predicting material accounting misstatements." Contemporary accounting research 28.1 (2011): 17-82.
- [12] Perols, Johan. "Financial statement fraud detection: An analysis of statistical and machine learning algorithms." Auditing: A Journal of Practice & Theory 30.2 (2011): 19-50.

Improved GMDA based DOA Technique using Pre-training Phase Unwrapping for Source Localization

S.-I. Kang, S. B. Kim and S. M. Lee

Abstract— In this paper, we propose a novel technique to improve the performance of generalized mixture decomposition algorithm (GDMA) based on a pre-training phase unwrapping. From an investigation of the GDMA scheme, it is discovered that the conventional GDMA method cannot take full consideration of phase unwrapping since the estimated inter-channel phase difference (IPD) slop is initialized by the random. To avoid this phenomenon, the proposed GDMA approach initialize IPD slope from the data of low-frequency bins. Experimental results show that compared to the conventional GMDA method the proposed GMDA technique based on pre-training phase unwrapping obtains low estimation error and when integrated into a source localization system achieves improved results.

Index Terms—Source localization, Inter-channel phase difference, Generalized mixture decomposition algorithm, Pre-training

I. INTRODUCTION

COURCE localization is an important tool used in many a Dimultichannel signal processing system and may include other functions such as source tracking, signal separation, speech enhancement, and noise suppression. A number of source localization algorithms have been proposed, e.g., adaptive eigenvalue decomposition algorithm associated with blind channel identification [1], LMS-type adaptive time delay estimation (TDE) algorithm [2] and the generalized crosscorrelation (GCC) method [3]. Recently, one of the successful source localization techniques is the generalized mixed decomposition algorithm (GDMA) which obtains an estimate the direction of arrival (DOA) of the sources by utilizing interchannel phase difference (IPD) between the dual channels on the sinusoidal tracking [4]. The GDMA based source localization technique is refined using the sinusoidal track method in the existing IPD distribution and shows robust performance at the white noise environment, but the result of the DOA estimation is sensitive to the accuracy of the phase unwrapping.

In this paper, we propose a GDMA-based DOA estimation incorporating the pre-training phase unwrapping for source localization. In order to efficiently determine the phase unwrapping of the GMDA method, pre-training is performed in a low-frequency bins where phase unwrapping does not occur. The proposed technique is substantially adopted for a source localization technique and is evaluated under various conditions.

II. REVIEW OF GENERALIZED MIXTURE DECOMPOSITION ALGORITHM (GMDA) BASED DOA TECHNIQUE

Let $x_1[k]$ received at one microphone denote a noisy signal that is the sum of a desired source signal s[k] and an uncorrelated additive noise signal $n_1 = [k]$; $x_1[k] = s[k] +$ $n_1[k]$. Another signal $x_2[k]$ is the sum of a delayed version of s[k] and noise signal $n_2[k]$. Applying a short-time discrete Fourier transform (DFT), we have in the time-frequency domain

$$X_1(\omega) = S(\omega) + N_1(\omega) \tag{1}$$

$$X_2(\omega) = S(\omega)e^{-j\omega\tau} + N_2(\omega)$$
⁽²⁾

where ω represents angular frequency.

The inter-channel phase difference (IPD) $\psi_X(\omega)$ that can be used for estimating the direction of arrival (DOA) between the two channel signals and is calculated as

$$\psi_X(\omega) = \angle X_1(\omega) - \angle X_2(\omega) \tag{3}$$

where $\angle X_1(\omega)$ and $\angle X_1(\omega)$ are the phase of $X_1(\omega)$ and $X_2(\omega)$, respectively. Considering the noise components, equation (3) is expressed as follows:

$$\psi_X(\omega) = \omega\tau + 2\pi n + \upsilon(\omega) \tag{4}$$

where τ is the time delay of the desired source, $2\pi n$ (n = ..., -1, 0, 1, ...) represents possible phase unwrapping and

This research was supported by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education (2010-0020163).

S.-I. Kang, S. B. Kim., and S. M. Lee are with the Department of Electronic Engineering, Inha University, Incheon 22212 Republic of Korea (e-mail: rkdtkddlr@gmail.com; zelabean@naver.com; sanglee@inha.ac.kr).

 $v(\omega)$ is the IPD error. The DOA of the desired source can be derived from τ using the following equation:

$$\tau = d\sin\theta / c \tag{5}$$

where d is inter-microphone distance, θ denotes the DOA of the desired source, and c is the sound speed.

In contrast to the white signal, the speech has a sparsity in the time-frequency domain. The speech signals are many short pauses and silence segments in the time domain. In the frequency domain, the power of the signal is concentrated on the harmonics of the pitch frequency in the voiced speech. The white noise signal can be effectively removed by applying the sinusoidal track method considering the sparsity of the speech signal [5].

According to equation (4), the points on the IPD versus frequency plot are spread over several lines based on the DOA information of the sources. The GMDA is adopted for clustering and estimation of the directions of multiple sources. The parameters of the mixture model are to be trained from the data. By employing the maximum-likelihood approach and using the expectation-maximization (EM) algorithm [6], the conditional expectation of complete data log-likelihood given the observed data under the previous parameter value is

$$Q(\mathbf{\Theta}; \mathbf{\Theta}(t)) = \sum_{i=1}^{N} \sum_{j=1}^{m} P(C_j \mid \mathbf{y}_i; \mathbf{\Theta}(t)) \ln(p(\mathbf{y}_i \mid C_j; \mathbf{\theta}) P_j)$$
(6)

where $\boldsymbol{\Theta} = [\boldsymbol{\Theta}^T, \boldsymbol{P}^T]^T$ is the parameters for the whole mixture model, $\boldsymbol{\Theta}$ is the parameter vectors for all clusters, and \mathbf{P} is *a priori* probability vectors. $P(C_j | \mathbf{y}_i; \boldsymbol{\Theta}(t))$ is the *a posterior* probability for class C_j given the previous parameter value $\boldsymbol{\Theta}(t)$ and data point \mathbf{y}_i . *t* is the iteration number.

The GMDA-based DOA technique proposed to select the phase unwrapping that yields the highest probability for the all observed data point. To look for a proper phase unwrapping factor for the IPD, this pdf is reversed to be

$$p(\mathbf{y}_{i} \mid C_{j}; \mathbf{\theta}_{j}) = \max_{n} \frac{1}{\sqrt{2\pi\sigma_{j}}} \times \exp\left\{-\frac{(\psi_{x,i}(\omega_{i}) + 2\pi n_{i} - \alpha_{j}\omega_{i})^{2}}{2\sigma_{j}^{2}}\right\}$$
(7)

where α_j is slope of the line, σ_j^2 is the variance of the model. Denote

$$n_i^j = \arg \max_{n_i} \frac{1}{\sqrt{2\pi\sigma_j}} \times \exp\left\{-\frac{(\psi_{x,i}(\omega_i) + 2\pi n_i - \alpha_j \omega_i)^2}{2\sigma_j^2}\right\} (8)$$

Let

$$J = \arg\max_{i} p(\mathbf{y}_{i} \mid C_{j}; \boldsymbol{\theta}_{j})$$
(9)

Then \mathbf{y}_i is chosen as

$$\mathbf{y}_i = [\omega_i, \psi_{x,i}(\omega_i) + 2\pi n_i^j]^T \tag{10}$$

III. ENHANCED GDMA BASED ON PRE-TRAINING PHASE UNWRAPPING

In the previous section, we note that the important parameter in the GDMA method is α_j , represent the slope of the line iteratively calculated based on the Gaussian model from the data point \mathbf{y}_i . The conventional GMDA approach proposed three ways to initialize α_j such as an ITD histogram method, a random initialization method, and uniformly distributed initialization method to determine the initial α_j . However, when the slope is randomly initialized, it takes a long time to converge when selected in the other direction, and it is not accurate when the distribution is symmetric. When histogram method is used, a large error is expected if the noise



Fig. 1. Inter-channel phase difference versus frequency. (DOA : 60°) (a) Original. (b) After adjusting phase using pre-training approach.

concentrated at a certain frequency is added.

Fig. 1 shows the IPD versus frequency plot original distribution and after phase unwrapping using pre-training approach. Since it is known that possible phase unwrapping parameter $2\pi n$ is zero in low frequency, we first consider the data point \mathbf{y}_i where distributed in low frequency as follows:

	Noise	-60?		-30?		0?		30?		60?	
		random	proposed								
Anechoic	clean	8.41	6.22	8.15	4.72	6.84	4.81	7.41	3.31	7.46	1.47
	white	9.59	4.12	8.79	4.32	8.27	3.73	9.73	4.61	8.94	1.48
	babble	10.16	4.28	9.59	4.55	8.16	4.39	9.15	5.38	8.17	2.68
Echoic	clean	9.17	8.93	8.91	6.01	8.06	4.57	8.61	5.04	8.96	4.06
	white	10.71	10.51	10.35	6.83	9.94	3.92	9.16	5.41	9.51	3.70
	babble	13.51	13.23	11.48	6.28	11.10	4.78	11.46	5.49	11.58	4.72

 TABLE I

 Relative RMSE for DOA estimation obtained from randomly initialized GMDA and proposed method.

$$\mathbf{y}_{i}^{T} = [\omega_{i}, \psi_{x,i}(\omega_{i})]^{T}, \omega_{i} < \eta$$
(11)

where η is the threshold of angular frequency. The initial slope $\alpha_i(0)$ is given by

$$a_{j}(0) = \arg\max_{a_{j}} \frac{1}{\sqrt{2\pi\sigma_{j}}} \times \exp\left\{-\frac{\left(\psi_{x,i}(\omega_{i}) - \alpha_{j}\omega_{i}\right)^{2}}{2\sigma_{j}^{2}}\right\}$$
(12)

IV. EXPERIMENTS

Several experiments are performed to research the performance of the proposed algorithm in various conditions. All experiments use two microphone channels and the source signal is the speech sound that constitutes the TIMIT database [7] with a sampling rate of 16 kHz. The speech average sound pressure level (SPL) is 65 dB SPL and noise (white, babble) average SPL is 55 dB SPL (SNR : 10 dB). The direction of the speech signal is tested at intervals of 30° from -60° to 60° . The noise location is fixed at -30°. Experiments are conducted in the anechoic chamber and the echo room to evaluate the performance depending on the presence or absence of echoes. The size of the echo room is $6m \times 6m \times 2m$. The window size of 25ms is used with frames overlapping frames, where the frames are shifted by 10ms. The distance between the microphones is 10cm. The mean absolute error (MAE) in the DOA estimates is used to evaluate the performance of the various algorithms.

Table 1 presenting the average DOA MAE performance shows that the proposed pre-training based GMDA approach outperformed original GDMA-based scheme which is randomly chosen initial slope. The conventional GMDA-based algorithm and proposed algorithm use the average of 10 experimental results to compensate for the random variation of performance. These results confirm that the proposed pretraining GDMA scheme substantially improves the GDMA method in source localization.

V. CONCLUSION

In this paper, we have proposed a robust approach to incorporate the pre-training into the conventional GDMAbased phase unwrapping estimation scheme for source localization. The initial slope for phase unwrapping is determined by the estimation method using the Gaussian distribution in the low frequency bins. Compared to the conventional GDMA method, it has been demonstrated that the proposed technique provides better phase unwrapping estimates in the source localization systems.

REFERENCES

- J. Benesty, "Adaptive eigenvalue decomposition algorithm for passive acoustic source localization," *J. Acoust. Soc. Amer.*, vol. 107, pp. 384– 391, 2000.
- [2] F. Reed, P. Feintuch, and N. Bershad, "Time delay estimation using the LMS adaptive filter—Static behavior," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. ASSP-29, no. 3, pp. 561–571, Jun. 1981.
- [3] D. Hertz, "Time delay estimation by combining efficient algorithms and generalized cross-correlation methods," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. ASSP-34, no. 1, pp. 1–7, Feb. 1986.
- [4] Wenyi Zhang and B. D. Rao, "A Two Microphone-Based Approach for Source Localization of Multiple Speech Sources," *IEEE Trans. Audio Speech Lang. Process.*, vol. 18, no. 8, pp. 1913–1928, Oct. 2010.
- [5] R. McAulay and T. Quatieri, "Speech analysis/ synthesis based on a sinusoidal representation," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. ASSP-34, no. 4, pp. 744–754, Aug. 1986.
- [6] N. M. Laird, A. P. Dempster and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Ann. Roy. Stat. Soc.*, pp. 1-38, Dec. 1977.
- [7] Defense, The DARAPA TIMIT Acoustic-Phonetic Continuous Speech Corpus (TIMIT) Speech Disc Cd1-1.1 ed., DARPA-ISTO, 1990.

HyDM: Data Migration Methodology for Hybrid Memories targeting Intel Knights Landing Processor

Jongmin Lee, Kwangho Lee, and Kidong Yun, Mucheol Kim, Geunchul Park, and Chan Yeol Park

Abstract-High-performance computing (HPC) systems provide huge computational resources and large memories. The hybrid memory is a promising memory technology that contains different types of memory devices, which have different characteristics regarding access time, retention time, and capacity. However, the increasing performance and employing hybrid memories induce more complexity as well. In this paper, we propose a data migration methodology called HyDM to effectively use hybrid memories targeting at Intel Knight Landing (KNL) processor. HyDM monitors status of applications running on a system and migrates pages of selected applications to the High Bandwidth Memory (HBM). To select appropriate applications on system runtime, we adopt the roofline performance model, a visually intuitive method. HyDM also employs a feedback mechanism to change the target application dynamically. Experimental results show that our HyDM improves over the baseline execution the execution time by up to 22%.

Index Terms—performance, data migration, roofline model.

I. INTRODUCTION

W ITH the ever-shrinking feature size in the CMOS process technology and increasing performance demands, modern processors typically integrate multiple cores and the number of cores in the same chip area has grown significantly. Continuous technology scaling realizes a many-core processor with hundreds of cores on a single chip [1]–[3]. These trends necessitate larger DRAMs to accommodate more and bigger programs in the main memory. DRAMs have been popularly used to implement the main memory because of their high densities and low prices. Due to the scaling limitation of DRAMs and the high bandwidth demands, hybrid storage architectures, which contain heterogeneous memories, are likely to be the future memory systems in high-performance computing (HPC) systems [4]–[7].

Knights Landing (KNL) is the code name for the secondgeneration Intel Xeon Phi product family [1], [8]. The KNL processor contains tens of cores and it provides the HBM 3D-stacked memory as a Muti-Channel DRAM (MCDRAM). DRAM and MCDRAM differ significantly in terms of access time, bandwidth and capacity. Because of those differences between DRAM and MCDRAM, performance will vary depending on the application characteristics and the usage of memory resources. The switch to multi/many-core processors and hybrid memories means that microprocessors will become more diverse. The growing complexity in HPC environments makes difficult for users to determine the performance of applications quantitatively. The roofline performance model is a simple and visual model that offers insights for performance analysis [9]. To evaluate performance, the roofline model ties floating-point performance (GFlops/second), arithmetic intensity (Flops/Byte), and memory bandwidth (GB/second) together. The peak floating-point performance and the peak memory bandwidth represent the attainable performance on a system and the arithmetic intensity shows a ratio of computations to memory accesses.

In this paper, we propose a data migration strategy for hybrid memories (HyDM). HyDM periodically monitors the application's execution status and it selects appropriate applications, which require more memory bandwidth. By migrating pages of the applications to the high bandwidth memory (i.e. MCDRAM), HyDM improves the memory usages on hybrid memories. In order to trace performance changes of applications during their executions, HyDM employs a feedback mechanism to change target applications dynamically. Our experimental results demonstrate that HyDM significantly improves the performance of mixed application sets on Intel KNL processor. HyDM enhances performance by up to 22% compared to the baseline execution time.

The rest of this paper is organized as follows. We provide background in the next section. Section III presents our proposed data migration strategy using the roofline model. Experimental results are given in Section IV. Section V presents related prior works. Section VI concludes this paper.

II. BACKGROUND

A. Intel Knight Landing (KNL) Processor

In this section, we briefly summarize the main features of the Intel KNL processor, especially we focus on its memory system.

Fig. 1 illustrates the KNL processor and its connection to the hybrid memories. The KNL processor integrates up to 72 cores together with eight Multi-Channel DRAM (MCDRAM) memories, which support 16GB of memory and they provide the peak bandwidth of 400GB/second. The processor also integrates six DDR4 channels supporting up to 384GB of memory with the peak bandwidth of 100GB/second. The MCDRAMs are positioned on-chip while DRAMs are offchip. Fig. 1 shows 36 tiles in the KNL processor and each tile consists of the two cores sharing 1MB L2 cache. Tiles are connected through a 2D-mesh network on-chip and they can

Jongmin Lee, Kwangho Lee, Kidong Yun are with the Department of Computer Engineering, Won-Kwang University, Iksan, Korea. E-mail: {square55, lkh002, s2mos2}@wku.ac.kr

Mucheol Kim is with School of Software, Chung-Ang University, Seoul, Korea. E-mail: mucheol.kim@gmail.com

Geunchul Park and Chan Yeol Park are with National Institute of Supercomputing and Networking, KISTI, Daejeon, Korea. E-mail: {gcpark, chan}@kisti.re.kr



Fig. 1. A structure of Intel Knight Landing processor



Fig. 2. Roofline performance model

be clustered in several NUMA configurations. In this paper, we only use the Quadrant cluster configuration where the tiles are partitioned in four quadrants as it reduces the latency of L2 cache misses because the worst-case path is shorter. This configuration is the one recommended by Intel as a symmetric multi-processor [10]. MCDRAM can be configured at boot time in three modes: cache, flat or hybrid mode. The Flat mode configures MCDRAMs to the same address space with DRAMs, Cache mode configures MCDRAMs as a last-level cache. The Hybrid mode separates MCDRAMs as two parts and one is used for an additional addressable memory with DRAMs and another is used for a last-level cache. In this work, we consider the Flat mode. For more details on KNL processor can be found in [1], [11].

B. Roofline Performance Model

The roofline performance model is a visually intuitive method used to bound the performance of floating-point programs running on multi/many-core processors [9]. Rather than simply using percent-of-peak estimates, the model can be used to evaluate the quality of attainable performance including locality, bandwidth, and computational throughput.

Fig. 2 shows the roofline model of Intel KNL processor



Fig. 3. HyDM Methodology Overview

with NAS parallel benchmark suites [12]. We periodically record the position of each benchmark to see performance changes of each benchmark over time. Detailed experimental methodologies will be shown in Section IV-A. The blackcolored lines show the peak performance of KNL processor with DRAM, MCDRAM, and floating-point units, respectively. The x-axis shows the arithmetic intensity that is the ratio of total floating-point operations to total data movement. The y-axis represents performance that is the number of floatingpoint operations completed by the cores. As shown in Fig. 2, most benchmarks are changing their positions over time and they are located under the memory-bound area with small arithmetic intensities. Arithmetic intensity with a small number means there are more memory requests, and the opposite case means more computations. Thus, one of the straightforward approaches to enhance the performance is moving data of the applications, which require more memory bandwidth, to the high bandwidth memory.

III. PROPOSED TECHNIQUES

In this section, we introduce a data migration methodology for hybrid memories called HyDM.

A. Overview

Fig. 3 shows an overview of HyDM method. HyDM employs three stages to enhance the performance of applications on the KNL processor. We first monitor the applications during system runtime using hardware monitoring tools. Then, based on the historical data, we select a candidate application, which requires more memory bandwidth, using the roofline model. Next, we migrate data stored in both MCDRAM and DRAM dynamically. By managing application data on MCDRAM and DRAM, HyDM effectively uses hybrid memories.

Algorithm 1 shows the implementation of HyDM. If running applications exist, HyDM makes the list of running application L (line 1). L stores unique PIDs for each application. The three stages of HyDM are repeatedly performed in a time window p. At each time point, historical monitored data from each application are stored in the list $W = \{W_0, W_1, ..., W_n\}$ (line 2). Note that W_0 is the current time window and the time window W_1 is the previous time window. W_n represents the

Algorithm	1	HyDMA lgorithm	(p))
-----------	---	----------------	-----	---

Input: $p \leftarrow$ time window period
/*
L : the list of running application
W: the list of windows for monitoring data
b: the selected application id
*/
initiate L
while $length(L) > 0$ do
wait p
1. $L \leftarrow getCurrentApplication()$
2. $Monitoring(L, W)$
3. $b \leftarrow Selection(L, W)$
4. $Migration(b, L, W)$
end while

Algorithm 2 Monitoring(L, W)

Input: $L \leftarrow$ the list of running application $W \leftarrow$ the list of windows for monitoring data /* M: the list of monitoring data for applications *i* : the index of application */ for $i \leftarrow L_0$ to length(L) do 1. $M_i \leftarrow getMonitoringData(i)$ 2. /* $M_i.fp: #$ of floating point operations */ 3. /* $M_i.pref$: # of page references */ 4. /* $M_i.pf$: # of page faults */ end for 5. insert M to W_0

n-th previous window. After the *Selection* procedure with the lists L and W, HyDM returns the candidate application b for migration (line 3). The pages of selected application b are migrated by the *Migration* procedure (line 4). We present the details of our method in the following subsections.

B. Monitoring

During the system runs, HyDM monitors applications using hardware monitoring tools. Most processors now include hardware support for performance monitoring such as *perf_event* [13] and LIKWID [14]. In this paper, we use perf_event. In Algorithm 2, the inputs include the list of running application L and the list of windows for monitoring data W. Let Mdenote the list of monitored data for running applications in the current time window. Let M_i denote the *i*th application. The Monitoring procedure collects the number of floatingpoint operations $(M_i.fp)$, page references $(M_i.pref)$ and page faults $(M_i.pf)$, and it stores those values into the entry corresponding to each type in M_i (line 1-4). $M_i.fp$ and $M_i.pref$ will be used to compute the arithmetic intensity of each application. The for loop stops when i is equal to the size of length(L). Then, M is inserted to the W_0 to prepare the next stage (line 5). Because HyDM only stores a few types of monitoring data, the storage overhead is very small compared to the total memory.

Algorithm 3 $Selection(L, W)$
Input: $L \leftarrow$ the list of running application
$W \leftarrow$ the list of windows for monitoring data
Output: $b \leftarrow$ the selected application id
/*
S: the sorted list of applications
V: the miss predicted list of applications
i: the index of applications
pf_{avg} : the averaged page faults
*/
1. $S \leftarrow regressionAndSort(W)$
2. $pf_{avg} \leftarrow getAvgPageFault(W_0)$
for $i \leftarrow 0$ to $length(S)$ do
3. if S_i in V then continue
4. if $S_i.pf > pf_{avg}$ then continue
5. $b \leftarrow getAppId(S_i, L)$
6. return b
end for

C. Selection

Algorithm 3 shows the *Selection* procedure that chooses an application as a candidate for migration. In order to select an application, which requires more memory bandwidth, HyDM uses the roofline model. When the execution status of applications is mapped to the roofline model, HyDM chooses an application with the lowest arithmetic intensity in the memory-bound area. The strategy in HyDM is to give more chances to the application that shows the highest ratio of memory references to computations.

The regressionAndSort procedure first computes the arithmetic intensity of each application using historical floatingpoint operations and page references stored in W (i.g. $W_{time.appid.fp}$ and $W_{time.appid.pref}$). After that, we perform the linear regression to predict the next arithmetic intensity value for each application. Finally, the application list (S)are sorted in ascending order according to the next arithmetic intensity values (line 1). Since all candidate applications are sorted in the list S, the first application of the list is considered for migration. We first check that the candidate application has a historical record of miss prediction (V). The list V generated in *Migration* procedure stores applications that did not show performance improvement after migration (line 3). In order to filter applications with low memory locality, HyDM employs a simple technique using a number of page faults monitored in the Monitoring procedure. HyDM compares the number of page faults from the first application with the averaged number of page faults (line 4). If all operations are finished, the *Selection* procedure returns the candidate application b for migration (line 6).

D. Migration

Algorithm 4 shows the *Migration* procedure. Because of the limited capacity of MCDRAM (16GB), we identify the possibility before performing the data migration. We check that the total memory usage, including the current memory usage of the selected application, does not exceed the threshold

Algorithm 4 $Migration(b, L, W)$	
Input: $b \leftarrow$ the selected application id	
$L \leftarrow$ the list of running application	
$W \leftarrow$ the list of windows for monitoring	data
/*	
V: the list of miss predicted applications	
r: the index of application for rollback	
t: the threshold ratio of MCDRAM use	
*/	
if $isMigrationPossible(t, b)$ then	
1. migrationToMCDRAM(b, W)	
else	
2. $r \leftarrow checkFlops(W)$	
3. if r exists then	
4. insert r to V	
5. $migrationToDRAM(r, W)$	
6. end if	
end if	

parameter t (e.g. 90%). If the usage of MCDRAM is less than t, the entire page of the selected application is migrated to MCDRAM (line 1). Although we migrate the entire page to maintain the simplicity of the HyDM, the page grouping techniques for selecting the critical pages of the entire page in the application are applicable to our proposed scheme [15], [16].

When the usage of MCDRAM is larger than t, we check the changes in flops for applications, which have migrated to MCDRAM (line 2). If the flops of an application have improved at least once during the time windows compared to the previous flops, we give more chances to the application to be in MCDRAM. Since migrating pages frequently induces additional overheads in terms of performance and energy, we employ a strict methodology for rollback. If we found the application that flops does not change, the *checkFlops* returns the application id r. After that, the application is inserted to the list V to record the miss prediction and pages of the application are migrated to the DRAM again (line 3-6). By employing the feedback mechanism above, HyDM effectively uses hybrid memories when many applications are running on a system.

IV. EXPERIMENTAL RESULTS

In this section, we present the methodologies for evaluations and their results with discussion.

A. Methodology

The experimental system is equipped with the Intel Xeon Phi(TM) CPU 7250@1.40GHz, 68 cores per socket, 4 threads per core, and a total of 272 threads available with the hyper-threading technology. The system includes 96GB DDR4 (DRAM) and 16GB HBM (MCDRAM).

We evaluated NAS Parallel Benchmark (NPB) related to computational fluid dynamics [12]. The NPB consists of five kernel benchmarks (IS, EP, CG, MG, FT) and three pseudo benchmarks (BT, SP, LU). For all experiments, we used standard test problems (CLASS-C). Table I shows the benchmark execution results when they are run on the system alone including averaged execution times (twenty times), floatingpoint operations, memory accesses, and the amount of peak memory use, respectively.

Table II shows the design parameters of HyDM and their values set in evaluations. The minimum time unit that can be monitored through *perf_event* is 1ms. When hardware events for monitoring are executed frequently, however, performance degradation occurs. We adjusted the numerical values without affecting performance through a heuristic method. To assume the system situation in which applications that require much larger capacity than the capacity of the MCDRAM are running, we perform the NPB programs in the number of multiples.

B. Performance Evaluation

Fig. 4 shows the changes in the roofline model when we run 16 NPB programs in parallel. We have selected several programs due to page limitation. Benchmarks in the loop-line model show different shapes, but the green-points representing HyDM have overall improved flop values than the baseline. On the roofline model, we can see that HyDM works effectively through the rise of the overall flops values.

To evaluate the performance impact of the proposed HyDM, we randomly assigned programs to the cores and calculated the average execution times. Fig. 5 show the averaged execution times for 100 runs. The results are normalized to the baseline. Several benchmarks show large increases in the execution time such as CG.0(15%), CG.1(22%), and FT.1(10%). On average, the reduced execution time is 6.5% with HyDM. This performance improvement of HyDM is due to memory-intensive applications are effectively migrated to MCDRAM.

TABLE I NAS PARALLEL BENCHMARK (NPB) CHARACTERISTICS

Name	Average execution time (sec.)	FP operations (GFlop)	Memory accesses (read/write) (mil.)	Peak memory use(MB)
IS.C	31	23	758 / 338	1,572
EP.C	462	494	2,739 / 1,028	20
CG.C	319	261	31,174 / 810	1,102
MG.C	129	213	7,507 / 2,940	3,536
FT.C	335	837	18,555 / 10,197	7,188
BT.C	920	2,560	45,682 / 17,270	1,676
SP.C	626	1,918	92,526 / 47,727	1,416
LU.C	733	2,504	84,716 / 38,302	760

TABLE II HyDM parameters

Descriptions	Values
The time window period: p	100 (ms)
The number of applications: $length(L)$	16, 32
Size of monitoring windows: $W = \{W_0, W_1,, W_n\}$	5, 10
The threshold ratio of MCDRAM use: t	90 (%)



Fig. 4. Roofline performance models



Fig. 5. Averaged execution time results

V. RELATED WORK

There are some categories of works that are closely related to this paper.

Hybrid memories: Many memory devices have been developed for decades to replace DRAM, which has fast but nonvolatile characteristics [17], [18]. PRAM is easier to integrate than DRAM, but the number of writable times per cell is limited thus memory life is short. STT-RAM has a fast write speed and good write durability, but it is relatively hard to integrate, and therefore, there is less need to replace DRAM in terms of economy. When examining the new memory technology to date, it is difficult to pursue universal memory, and it is judged to be a non-volatile technology that lacks performance rather than DRAM. To use those memories, many researches have been done on hybrid memories in the form of DRAM and other types of memories together. One of the hybrid memory systems uses DRAM as a cache and PRAM as a main memory [19], [20]. They mitigate the durability of PRAM and write delay by filtering the write operations to the main memory using the DRAM cache. In [21], DRAM and PRAM are located at the same level. In order to compensate for the delay in the write operation and the lifetime of the PRAM, a page manager selectively allocates pages among PRAM and DRAM. All of the above techniques are designed to reduce write activities in PRAM, however, this paper addresses the usage of HBM with DRAM.

GPU is the most commonly used hybrid memory to date in HPC [22]. GPU employs GDDR as high-speed memory and relatively slow DRAM as main memory. By storing critical data using prefetch techniques in GDDR, GPU supports fast operation. The GPU operates as an accelerator with respect to DRAM. By comparison, our research explores general processors for HPC environments.

Roofline performance model: The roofline model is used in a number of scientific applications to analyze bottlenecks in the performance of an architecture and to guide software optimizations [9]. Various types of roofline models are proposed in previous works [6], [23]–[25]. In [23], energy version of the roofline model is proposed to show bounds on performance due to energy limitations. This model focuses on identifying the balance between performance and energy in architectural design. In [24], the roofline model is extended to support the cache hierarchy. Recently, the roofline model is extended for specific applications and platforms such as GPUs [6].

Page Migration: A variety of page migration methods using NUMA nodes have been studied [26]–[28]. A basic methodology to efficiently use memories in a NUMA system is to store the data in the same location as the processor that frequently references the data. In [26], the migration of the pages between nodes is performed by using the characteristic that the memory access pattern repeatedly appears in applications. In [29], a sampling-based approach is used in which pages with excessive remote references are migrated to nodes close to the accessing core. The system continuously samples the excess miss counters to produce a list of candidate pages for migration and replication.

We propose a dynamic memory management methodology using the roofline model, the key contribution of our work is the algorithm that efficiently uses different types of memories in HPC systems without any hardware or software modifications. Although our proposed HyDM targets the Intel KNL processor in this paper, adopting the methodology to the systems employing hybrid memories is possible.

VI. CONCLUSION

The hybrid memory is a promising memory technology for future HPC systems. However, effective use of the system is becoming increasingly difficult as the HPC environment is diversifying. In this paper, we proposed a dynamic data migration strategy using the roofline performance model called HyDM. HyDM uses a hardware monitoring tool to observe the status of programs running on the system and perform migration based on the collected data. Also, a feedback mechanism is implemented for the case where the total memory usage used for the programs is larger than the size of the high bandwidth memory. Experiments using a real system, including the Intel KNL processor, we demonstrate that the proposed HyDM improves system performance.

ACKNOWLEDGMENT

This work was supported by Korea Institute of Science and Technology Information (KISTI) grant funded by the Korea government (MSIT) (No.K-18-L12-C07-S01).

REFERENCES

- A. Sodani, "Knights landing (KNL): 2nd Generation Intel® Xeon Phi processor," in *IEEE Hot Chips 27 Symposium*. IEEE, 2015, pp. 1–24.
- [2] R. Espasa, "Larrabee A Many-Core Intel Architecture for Visual Computing." *HiPEAC*, vol. 5952, no. Chapter 2, pp. 2–2, 2010.
- [3] A. K. Singh, M. Shafique, A. Kumarm, and J. Henkel, "Mapping on multi/many-core systems - survey of current and emerging trends." DAC, p. 1, 2013.
- [4] I. B. Peng, R. Gioiosa, G. Kestor, P. Cicotti, E. Laure, and S. Markidis, "Exploring the Performance Benefit of Hybrid Memory System on HPC Environments," in *IPDPSW*. IEEE, 2017, pp. 683–692.
- [5] I. B. Peng, R. Gioiosa, G. Kestor, J. S. Vetter, P. Cicotti, E. Laure, and S. Markidis, "Characterizing the performance benefit of hybrid memory system for HPC applications," *Parallel Computing*, vol. 76, pp. 57–69, 2018.
- [6] A. Lopes, F. Pratas, L. Sousa, and A. Ilic, "Exploring GPU performance, power and energy-efficiency bounds with Cache-aware Roofline Modeling," in *ISPASS*. IEEE, 2017, pp. 259–268.
- [7] O. Mutlu, "Memory scaling: A systems architecture perspective," in *IMW*. IEEE, 2013, pp. 21–25.
- [8] I. Jabbie, "Performance Comparison of Intel Xeon Phi Knights Landing," SIAM Undergraduate Research Online, vol. 10, 2017.
- [9] S. W. Williams, A. Waterman, and D. A. Patterson, "Roofline: An insightful visual performance model for floating-point programs and multicore architectures," EECS Department, University of California, Berkeley, Tech. Rep., 2008.
- [10] Colfax, "Clustering Modes in Knights Landing Processors," 2016.
- [11] J. Jeffers, J. Reinders, and A. Sodani, *Knights Landing architecture*. Morgan Kaufmann, Jan. 2016.
- [12] D. H. Bailey, "NAS Parallel Benchmarks." Encyclopedia of Parallel Computing, no. Chapter 133, pp. 1254–1259, 2011.
- [13] V. M. Weaver, "Linux perf_event Features and Overhead," in *FastPath Workshop*, 2013.
- [14] J. Treibig, G. Hager, and G. Wellein, "LIKWID: A Lightweight Performance-Oriented Tool Suite for x86 Multicore Environments," in *ICPPW*. ACM, 2010, pp. 207–216.
- [15] L. E. Ramos, E. Gorbatory, and R. Bianchini, "Page placement in hybrid memory systems," in *ICS*. ACM, 2011, pp. 85–95.
- [16] D. Shin, S. Park, S. Kim, and K. Park, "Adaptive page grouping for energy efficiency in hybrid PRAM-DRAM main memory," in ACM Research in Applied Computation Symposium. ACM, 2012, pp. 395– 402.

- [17] J. Meena, S. Sze, U. Chand, and T.-Y. Tseng, "Overview of emerging nonvolatile memory technologies," *Nanoscale Research Letters*, vol. 9, no. 1, p. 526, 2014.
- [18] R. F. Freitas and W. W. Wilcke, "Storage-class memory: The next storage system technology," *IBM Journal of Research and Development*, vol. 52, no. 4.5, pp. 439–447, 2008.
- [19] M. K. Qureshi, V. Srinivasan, and J. A. Rivers, "Scalable high performance main memory system using phase-change memory technology," *ACM SIGARCH Computer Architecture News*, vol. 37, no. 3, p. 24, Jun. 2009.
- [20] Y. Ro, M. Sung, Y. Park, and J. H. Ahn, "Selective DRAM cache bypassing for improving bandwidth on DRAM/NVM hybrid main memory systems," *IEICE Electronics Express*, vol. 14, no. 11, pp. 20170437– 20170437, 2017.
- [21] G. Dhiman, R. Z. Ayoub, and T. Rosing, "PDRAM a hybrid PRAM and DRAM main memory system." DAC, p. 664, 2009.
- [22] K. Nakano, "The Hierarchical Memory Machine Model for GPUs," in *IPDPSW*. IEEE, 2013, pp. 591–600.
- [23] J. Choi, D. Bedard, R. J. Fowler, and R. W. Vuduc, "A Roofline Model of Energy." *IPDPS*, pp. 661–672, 2013.
- [24] A. Ilic, F. Pratas, and L. Sousa, "Cache-aware Roofline model -Upgrading the loft." *Computer Architecture Letters*, vol. 13, no. 1, pp. 21–24, 2014.
- [25] D. Doerfler, J. Deslippe, S. Williams, L. Oliker, B. Cook, T. Kurth, M. Lobet, T. M. Malas, J.-L. Vay, and H. Vincenti, "Applying the Roofline Performance Model to the Intel Xeon Phi Knights Landing Processor." *ISC Workshops*, vol. 9945, no. 16, pp. 339–353, 2016.
- [26] W. J. Bolosky, R. P. Fitzgerald, and M. L. Scott, "Simple But Effective Techniques for NUMA Memory Management." SOSP, pp. 19–31, 1989.
- [27] R. P. LaRowe, C. S. Ellis, and M. A. Holliday, "Evaluation of NUMA memory management through modeling and measurements," *IEEE Transactions on Parallel and Distributed Systems*, vol. 3, no. 6, pp. 686–701, 1992.
- [28] Z. Majo and T. R. Gross, "Memory management in NUMA multicore systems," in *ISMM*. New York, New York, USA: ACM Press, 2011, p. 11.
- [29] L. Noordergraaf and R. van der Pas, "Performance experiences on Sun's Wildfire prototype," in ACM/IEEE conference on Supercomputing. New York, New York, USA: ACM Press, 1999.

An Index Table Based Optimally Matched Reversible Image Watermarking Scheme

C. C. Chen, J. M. Yang, and H. Q. Liu

Abstract— A reversible image watermarking scheme recovers the original image after extracting the embedded watermarks. The conventional Class Index Table (CIT) method changes the order of an original block to fit the order of a target block to acquire the watermarked image having a similar visual effect like the target image. In this study, we propose two improvements *RPCIT* and *OMCIT* to improve visual quality of the watermarked image. Experimental results show that the proposed *RPCIT*, using a randomized preprocessing, improves visual quality a little with random information included in the watermarked image. The proposed *OMCIT* outperforms related CIT-based schemes on good visual quality on watermarked image.

Index Terms— Reversible Image Watermarking, Class Index Table (CIT), Random Preprocessing, Optimally Matching.

I. INTRODUCTION

RECENTLY, hiding information into data is an interesting topic, especially on audio, video or image. Lots of Reversible Image Watermarking techniques have been reported in the past two decades and almost all recent RDH methods first generate PEs as the host sequence [4] with then reversibly embedding the message into the host sequence by modifying its histogram with methods like histogram shifting. Moreover, most of these methods can be viewed as a process of semantic lossless compression [5], in which some space is saved for embedding extra data.

In recent years, RDH in encrypted images (RDH-EI) has caught many researchers' attention. In RDH-EI, the original image is encrypted by content provider and the secret data is embedded by data hider with the data hider having no knowledge about the encrypted image. Zhang [6] presented a RDH-EI scheme for by using 3 specific bit value on each pixels of a block. Hong et al. [7] then improved REH-EI by using better estimations and a side match technique. The main idea on EDH-EI is that receiver can obtain different contents (like secret data or original image) with different access rights. If an unauthorized access happens, the content owner encrypts his personal images before store them.

Inspired by the technique of image transformation proposed by Lee and Tsai [3], their method can transform the original image to a freely-selected target image with the same size, yielding a secret-fragment-visible mosaic image defined in [1]. However, the original image cannot be restored in a lossless way. Therefore, their scheme is not reversible.

Since the reversible data hiding technique studies perfect ways to recover original image after watermark extraction, this paper improves previous scheme to have a RDH-EI with better representation on watermarked image. By using a preprocessing on randomization, the watermarked image has a visually similar to target image and the original image can be perfectly recovery. Further use of Class Index Table (CIT) acquires good results on watermarked images.

The paper is organized as follows. Section 2 briefly reviews important related CIT encryption work [2]. Section 3 introduces our two proposed schemes, including Random Preprocessing CIT (*RPCIT*) and Optimally Matching CIT (*OMCIT*) methods. Section 4 demonstrates our experimental results. Section 5 follows with concluding remarks.

II. REVIEW ON CLASS INDEX TABLE (CIT)

This section briefly reviews the CIT-based reversible data hiding by Zhang et al. [2]. In this scenario, the original image I is the information that we want to embed and the target image T is the cover image that we used. Moreover, the watermarked image I_w is the embedded result. Fig. 1 shows the concept of the CIT-based reversible image watermarking strategy.



Fig. 1. The CIT-based reversible image watermarking strategy.

A CIT (Class Index Table) is selected for transferring the semantic of an original block to the semantic of a target block for acquiring the watermarked block. The usage of the CIT generates a similar visual effect among the target block and the watermarked block. Furthermore, the original block can be perfectly recovery by using the CIT and the watermarked block. Therefore, by partitioning the original image and the target image into original blocks and target blocks, respectively, the

C. C. Chen is with the Department of Computer Science and Information Engineering, Taipei, Taiwan. (corresponding author to provide e-mail: ccchen34@mail.tku.edu.tw).

J. M. Yang was with the Department of Computer Science and Information Engineering, Taipei, Taiwan (e-mail: holt200185@gmail.com).

H. Q. Liu is with the Department of Computer Science and Information Engineering, Taipei, Taiwan (e-mail: a089048@gmail.com).

acquired watermarked image look like the target image.

A CIT for *n* pixels block is *n* bits binary values. A (n, t) CIT includes *t* larger and *n*-*t* smaller pixels, where *t* is user determined parameter. Fig. 2 shows an example of (8, 3) CIT method of an 8 pixels' block with *t*=3. This example books the positions of 3 largest pixels in this block.



Fig. 2. An example of (8, 3) CIT method.

The (n, t) CIT-based embedding process is illustrated as follows.

1. Acquire the *n*-pixels' CITs for the original block and the target block, respectively.

2. Change the pixels' positions in the original block to fit the target's CIT by the following steps.

2.1 In watermarked block, orderly replace the *t* larger pixels in original block to the 1's position in target CIT.

2.2 Replace the remaining positions in watermarked block by the numbers of the remaining *n*-*t* pixels orderly.

3. Acquire the CIT for the new acquired watermarked block.

The (n, t) CIT method transfers the positions of original block to fit the order positions in target block. Therefore, each pair of original block and target block has its transferring data of target block's CIT. Therefore, the embedding process requires the original block and the CIT of the target block. Moreover, the recovery process can perfectly acquire the original block by using the watermarked block and the CIT of original block. Fig. 3 shows an example of (8, 3) CIT-based embedding process, where CITs can be directly acquired from each block.



Fig. 3. An embedding process of (8, 3) CIT.

The embedding process of CIT includes segmentation, matching, and embedding steps. First, the original and the target is segmented to non-overlapped blocks with the same size. The second step is to find a match between original and target blocks. The third step is the embedding process introduced above. Assume that the image size is $M \times N$, the number or $M \times N$

original and target blocks are all n under the (n, t) CIT-

based embedding process. By using the block matching method, the watermarked image can be then acquired. Fig. 4 shows experimental results of the sequential matched CIT method, that matches all blocks by the sequential block number, by the original image is Barbara and target image is Casablanca. Fig. 5 shows the conventional sequential matched CIT method, in which the blocks are matched by the block number, sequentially.

Fig. 4(c) significantly shows that the watermarked image looks like the original image with some noise information. The noise-like property exists because of the pixel distribution difference between the original and the target blocks. From the image structure of two different images, using (n, t) CIT-based embedding process generates same block pixels' order but not visually similar as shown in Fig. 4(c). Our study presents better matching method for improving watermarked image's visual quality.



Fig. 4. Experimental results of the sequential matched CIT method, (a) the original image: Barbara, (b) the target image: Casablanca, and (c) the watermarked image.

0	}←───→	0
1	 ←───→	1
N-2	<>	N-2
N-1	 ←───→	N-1

original block number target block number

Fig. 5. The conventional sequential matched CIT method.

Moreover, Zhang et al. [2] also present a block matching method. However, they are not clearly defined the process of their proposed block matching method. Therefore, in this study, an optimal block matching method is present for acquiring better visual effect.

III. PROPOSED SCHEMES

This section shows our presented embedding and recovery algorithms. The proposed embedding algorithm embeds original image I into target image T to acquire watermarked image I_w . The recovery algorithm recovers the original image I from the watermarked image I_w and corresponding CIT C. Therefore, the correspond CIT C is denoted as the secret key for watermark extraction. We propose two schemes in this study. The first one is *randomized preprocessing CIT* method (*RPCIT*). The second one is *optimally matching CIT* method (*OMCIT*). The proposed *RPCIT* and *OMCIT* schemes are introduced in Sections III.A and III.B, respectively.

A. Randomized Preprocessing CIT method (RPCIT)

This section introduced our proposed RPCIT method. Since the difference on block pixels distribution is a serious problem for acquiring better watermarked quality in the conventional CIT method, the proposed RPCIT method adopts the random preprocessing on all images for acquiring similar distribution on image blocks. Fig. 6 depicts the steps of acquiring the watermarked image W from the original image I and the target image T. Moreover, the CIT of the randomized original image I_R is also generated. In the proposed embedding procedure, the original image I and the target image T are first performing XOR operation with a random image R to acquire the corresponding randomized original image I_R and the randomized target image T_R , respectively. After applying the I_R and T_R to the conventional CIT embedding method, the randomized watermarked image W_R and the CIT table RC of the original image are generated. After applying the XOR operation to the W_R and the random image R, the watermarked image W is then obtained. Fig. 6 shows that the generated data are marked as gray blocks.



Fig. 6. Embedding procedure of the proposed RPCIT method.

Fig. 7 shows the extraction procedure of the proposed *RPCIT* method. In this extraction procedure, two inputs, the watermarked image W and CIT table RC of the original image, are needed. The input W is first performed the XOR operation with the random image R to acquire the randomized watermarked image W_R . By applying the W_R and RC to the conventional CIT recovery method, the randomized original image I_R is then acquired. After performing the XOR operation to the I_R and R, the original image I is obtained.



Fig. 7. Extraction procedure of the proposed RPCIT method.

B. Optimally Matching CIT method (OMCIT)

This section introduces our proposed *OMCIT* method. In the *OMCIT* method, the main novelty is the optimally matching method of blocks between the original image and the target image. Figure 8 shows the structure of the proposed *OMCIT* method, in which two outputs the watermarked image W and CIT table RC are acquired.



Fig. 8. Embedding procedure of the proposed OMCIT method.

The main part of the OMCIT in Figure 8 is the optimally matched CIT embedding procedure and this procedure is depicted in Figure 9. In this proposed procedure, the original image I and the target image T are first partitioned to nonoverlapped 2×4 blocks I_b and T_b , respectively. Then, a loop is created for matching pairs of blocks from lower threshold to higher threshold. A threshold parameter thre is adopted for choosing the matched pair of blocks. There are two selection loops existed in the procedure. The first loop finds all matched blocks through the distance less than the threshold thre. The second loop increases the value of thre by an increment of threshold_step for matching all image blocks, where threshold_step is a user-defined parameter. Moreover, in our experiments, the calculation of distance(i, j) in two image blocks *i* and *j* is determined from the Euclidean distance of block mean and variance.



Fig. 9. Procedure of the proposed optimally matched CIT embedding method.

Fig. 10 shows the recovery process of the proposed *OMCIT* method. The recovery procedure is accomplished like the conventional CIT method. By using the watermarked image W

and the original CIT table CR, the original image can be perfectly recovered.



Fig. 10. The recovery procedure of the proposed OMCIT method.

IV. EXPERIMENTAL RESULTS

This section demonstrates the experimental results of our proposed two schemes. Fig. 11 shows the experimental results of the original image being Barbara as shown in Fig. 11(a) and the target image being Casablanca as show in in Fig. 11(b). Fig. 11(c)-(e) show the watermarked images of the conventional sequential matched CIT, the proposed RPCIT, and the proposed OMCIT methods, respectively. Fig. 11(c) shows that the conventional sequential matched CIT acquires the watermarked image still looks like the original mage rather than the target image. Fig. 11(c) shows that a good matching method is very important. Fig. 11(d) shows that using random preprocessing acquires better quality of the watermarked image. However, noise-like information is also included in this watermarked image. Fig. 11(e) shows the proposed OMCIT method has best visual quality of the watermarked image. Fig. 11 also shows that the CIT table RC is needed in the CIT-based reversible image watermarking method.



Fig. 11. Experimental results, (a) original image: Barbara (b) target image: Casablanca (c) the watermarked image of conventional sequential matched CIT (d) the watermarked image of the proposed *RPCIT* (e) the watermarked image of the proposed *OMCIT*.

V. CONCLUSIONS

This paper presents two improvements *RPCIT* and *OMCIT* of conventional CIT-based reversible image watermarking method. By using CIT, we can transform the same sized original image and target image to acquire the watermarked image. Comparing with the conventional method, our proposed

schemes improve quality of the watermarked image quite a lot. The original image can be perfectly recovered using the watermarked image and the CIT table *RC*. Experimental results show that the watermarked image generated from the proposed *OMNIT* outperforms than other CIT-based schemes. In the future, a block matching is required for acquiring better watermarked image.

REFERENCES

- I. Lai and W. Tsai, "Secret-fragment-visible mosaic image-a new computer art and its application to information hiding," IEEE Trans. Inf. Forensics Security, vol. 6, no. 3, pp. 936–945, Sep. 2011.
- [2] W. M. Zhang, H. Wang, D. Hou, and N. Yu," Reversible Data Hiding in Encrypted Images by Reversible Image Transformation", IEEE Trans. Multimedia, vol. 18, no. 8, Aug. 2016.
- [3] Y. Lee and W. Tsai, "A new secure image transmission technique via secret-fragment-visible mosaic images by nearly reversible color transformation," IEEE Trans. Circuits Syst. Video Technol., vol. 24, no. 4, pp. 695–703, Apr. 2014.
- [4] I. C. Dragoi and D. Coltuc, "Local-prediction-based difference expansion reversible watermarking," IEEE Trans. Image Process., vol. 23, no. 4, pp. 1779–1790, Apr. 2014.
- [5] T. Kalker and F. M. Willems, "Capacity bounds and code constructions for reversible data-hiding," 14th Int. Conf. Digital Signal Processing, pp. 71-76, 2002.
- [6] X. Zhang, "Reversible data hiding in encrypted images", IEEE Signal Process. Lett., vol. 18, pp. 255–258, 2011.
- [7] W. Hong, T. Chen, and H. Wu, "An improved reversible data hiding in encrypted images using side match", IEEE Signal Process. Lett., vol. 19, pp. 199–202, 2012.
- [8] J. C. Chang, Y. H. Chou, C. H. Ni, H. L Wu, "Reversible Data Hiding in Pairwisely Encrypted Images", 2016 Third International Conference on Computing Measurement Control and Sensor Network, pp. 60–63, 2016.

PCA based Performance Analysis with System Profiling Data in Many-core system

Mucheol Kim, Junho Kim, Jongmin Lee Geunchul Park and Chan Yeol Park

Abstract— High performance computers should be performing the high-level computations for various research fields. Then, it is important to enhancing the system performance with statistical approaches and data mining technologies. This paper should analyze predictable correlations from system profiling data. Our approach is based on PCA algorithm, it is determined proper clusters for the performance indices. Experimental analysis suggests the correlation between performance indices, they could suggest decision supports for scale-out architecture in HPC environments.

Index Terms— High Performance Computer, PCA, Data Analytics, System Profile, Many-core System.

I. INTRODUCTION

HIGH Performance Computers(HPC) has an important role in the field of various scientific research fields such as quantum mechanics, weather forecasting, molecular modeling [1]. We expect that High Performance Computer suggests new insights during high-level computations. In addition, parallel processing technology is emerging as a means to avoid performance issues related to processor speed through many-core processing.

The KNL processor has been able to boot independently without a host CPU in comparison with previous KNights Corner (KNC). It consists of 32 ~ 36 tiles connected in 2D-mesh form, and presents a theoretical maximum performance over 3TFLOPS [2]. The Knights Landing processor has two defferent memories which are DDR4 and MCDRAM. MCDRAM is a high-performance memory with four times the bandwidth difference theoretically compared to DDR4. They could improve both of memory capacity and bandwidth with their combinations [3].

Manuscript received September 30, 2018. (Write the date on which you submitted your paper for review.)

Mucheol Kim is with School of Software, Chung-Ang University, Seoul, Korea. (corresponding author to provide phone: +82-2-820-5327; fax: +82-2-820-5327; e-mail: Mucheol.kim@gmail.com).

Junho Kim., is with dept. Computer Science & Engineering, Chung-Ang University, Seoul, Korea (e-mail: kjhcau@gmail.com).

Jongmin Lee is with the Dept. of Computer & Software Engineering, Wonkwang University, Iksan, Korea (e-mail: square55@wku.ac.kr)

Geunchul Park is with National Institute of Supercomputing and Networking, KISTI, Daejeon 34141, Korea(e-mail : gcpark@kisti.re.kr)

Chan Yeol Park is with National Institute of Supercomputing and Networking, KISTI, Daejeon 34141, Korea(e-mail:<u>chan@kisti.re.kr</u>)

In order to optimize parallel computing performance in HPC environment, it is essential to understand and analyze the architecture of many-core processor. In order to effectively utilize the many-core processor, it is important to parallelize the tasks to divide tasks into numerous small tasks so that they can be effectively allocated over multiple cores. Appropriate scheduling policies are considered to enhancing the system performance. In addition, there are many researches for scale-out architectures, then they are optimizing the performance in many-core system. Therefore, it is necessary to analyze the system profiling data with the data mining technology.

Previous researches have applied various statistical methodologies and data mining techniques to analyze system profiling data [4][5]. However, as the complexity of HPC system increases, it becomes more important to derive the correlation between performance indicators. Then it is necessary to analyze the relationship between the performance indicators. Furthermore it is important to utilize their application for performance enhancement. In this paper, we focus on analyzing predictable correlations and sequencing from system profiling data and mining repeated sequence of events. This is a very important factor in determining the resource allocation policy of HPC according to the characteristics of upcoming tasks that can be performed.

This paper is organized as follows. Section 2 describes how data mining techniques are utilized in our performance analysis with benchmark datasets. In Section 3, we describe the experiment's results and followed analytic discuss. Section 4 concludes with future directions.

II. PERFORMANCE ANALYSIS WITH SYSTEM PROFILING DATA

A new architecture suitable for scale-out is required for distributed processing using many-core processor as a host processor. There is a need for optimizing methodology for improving resource utilization and energy efficiency. In this section, we present an analytic methodology with data mining technology for analyzing system performance in HPC environments.



Figure 1. Process for Analytic Approach for System profiling Data

A. Preprocessing with A-priori Algorithm

In this paper, we perform clustering and extracting the correlation between properties in system profiling data. Then, principal components are extracted using PCA algorithm and the correlation between the properties is analyzed. In this process, the preprocessing with A-priori algorithm [7] was performed in order to exclude needless properties in the correlation deduction process. Then, we changed the zero and nonzero variables according to whether the value was increased or decreased in comparison with the previous row. As a result, when all the values are equal in the column, we should remove the column. Because it could cause the confusion to recognize the association with other properties.

B. PCA based Unsupervised Learning

PCA(Principal Component Analysis)[8] is widely used as an algorithm for identifying clusters of data that has similar characteristics. Meanwhile, some variables make a cluster different from another. That is, it simplify that data by decomposition of variables into each certain cluster named principal component. Then PCA could be performed by eigenvalue decomposition of a data correlation matrix with normalization of the initial data. The results of a PCA also displays the relationship between variables with data variance which is the magnitude of eigenvalues.

In this paper, elbow method was used to find the number of principal components. The Elbow method helps to infer the appropriate number of clusters through consistency and validation of data. This is a method to investigate the ratio of variance according to the change in the number of clusters. The point at which the degree of variance decreases is called the Elbow Point. In addition, it is possible to estimate Elbow points with cumulative proportions. As a result, performance measures corresponding to the number of principal components obtained through inference can be identified through data correlation. In this case, the data correlation could be can be deducted through the eigenvectors and eigenvalues of the individual measure factors.

III. PCA BASED PERFORMANCE ANALYSIS WITH SYSTEM

A. Experimental Environments

We profiled and analyzed the system in HPC environments. In this experiments, we profile the data from Intel Xeon Phi CPU 7250 system which has 68Cores and supports 272Threads. Kernal Version is Linux 3.10.0-514.10.2.e17.x86_64. We collected the system log data with resource monitoring tools ("collect" and "dstat"). "collectl" provides each CPU, memory fragmentation, disk, luster file system, network, fan, power, and temperature information. "dstat" provides CPU usage, disk, network, process, memory, file system, and virtual memory information.

On the other hand, we collected the system profile data with NAS Parallel Benchmarks (NPB) [9]. They are a small set of programs for evaluating the performance of parallel supercomputers. The benchmarks are derived from computational fluid dynamics applications and five kernels (IS, EP, CG, MG, FT) and three pseudo-applications (Block Tri-diagonal solver, Scalar Penta-diagonal solver, Lower-Upper Gauss-Seidel solver).





Figure 2. The number of clusters with elbow method in EP

In this section, we analyzed the relationship between performance indices with PCA algorithm. At first, we determined the number of clusters with elbow method. It means the performance indicators for performing the computational fluid dynamic application in each kernel.

(Figure 2) displayed the number of clusters through from the embarrassingly parallel kernel the elbow method. Then (a) and (b) show the performance indicators obtained through resource monitoring collected by 'collectl' and 'dstate', respectively.

The number of clusters are determined by the point where the proportion of variance decreases, that is, the elbow point. (Figure 2a) showed that number of clusters are determined to 2. On the other hand, monitored resources by 'dstate' in (Figure 2b) selects 3 or 4 clusters. If two or more elbow point candidates appeared, the cumulative contribution rate should be considered.



Figure 3. The number of clusters with elbow method in MG

Figure 3 showed the number of clusters for resources monitored in the MG (Multi-Grid on a sequence of meshes, long- and short-distance communication, memory intensive) kernel. In this case (Figure 3a), two clusters are selected, and in case of (Figure 3b), it is appropriate to select 3-4 clusters, respectively.

The 1st Principal Component (PC1) of the MG is composed of v1 (sys: Time spent in "pure" system time), v6 (KBIn received / sec.) And v7 (PktIn: Received packets / sec., There is a high correlation between v0 (cpu: The CPU number) and v2 (inter: the interrupt summary stats). These can be represented by factors related to execution time and CPU allocation, respectively.

(Figure 4b) shows that the performance indicators of the 2^{st} cluster are v0 (usr: cpu usage by a user processes), v8 (used: amount of used memory), v9 (cach: amount of cached memory) number of running processes), and v17 (int: number of interrupts). It can be inferred that PC1 is the clustered main component of cpu, memory, cache memory usage, number of running processes and number of interrupts by usr process. In the PC2 element, v3 (disk_read: total number of read operations on disks) and v15 (io_read: read I / O requests). These means that the number of disk accesses according to I / O requests is highly related.

(a) Collectl - MG

	PC1	PC2	PC3	PC4
vO	-1.762466e-02	-0.555300839	0.015743116	0.08936169
V1	9.961689e-05	-0.414226427	-0.225255779	-0.35059117
v2	-1.801438e-02	-0.555178168	0.016611802	0.08772109
V3	1.515590e-02	0.443959215	-0.102629502	-0.24246786
v4	-6.250706e-04	-0.056290602	0.585203199	-0.39466743
v5	1.400118e-03	-0.033600088	0.643546656	-0.29048258
V6	9.814762e-02	0.014588301	0.035482888	0.04298882
ν7	9.945072e-01	-0.026060141	0.005532339	0.01341448
V8	4.619016e-03	0.009942001	-0.316103798	-0.65623943
v9	2.082858e-02	-0.097192964	-0.283508175	-0.36026932

(b) dstate - MG

Rota	ation (n x k)	= (19 x 19):				
	PC1	PC2	PC 3	PC4	PC 5	
V0	0.340152202	0.019048120	0.011614280	-0.054767793	0.032020317	
v1	0.143994621	0.099925625	0.186456853	-0.371304175	-0.209459144	e
v2	-0.340105600	-0.024854584	-0.020978998	0.078802369	-0.016228121	1
v3	-0.018171005	0.137361473	-0.626467433	-0.222620305	-0.095614933	
v4	0.023434863	-0.028198802	0.022031638	-0.067806572	0.293009688	
v5	0.281490049	-0.236939331	-0.109917817	-0.100423471	-0.042694458	
v6	0.119793855	-0.505996556	-0.130123179	-0.060871084	-0.025657790	4
v7	-0.068312312	-0.777716529	-0.139314553	-0.006499454	0.002612224	
v8	0.339451964	0.020512289	0.009104954	-0.050665213	0.060465366	
v9	0.327317705	0.015240025	0.042222185	-0.036461363	0.093138028	
v10	-0.339404168	-0.022153559	-0.009704869	0.050744791	-0.060572033	
v11	0.029056754	-0.161547972	0.232978742	0.002502699	-0.369812286	
v12	0.013479101	0.030579519	0.108687003	-0.116323918	-0.656723264	•
v13	0.339120699	0.030374877	0.023643871	-0.079148354	0.052555519	
v14	-0.174687585	-0.050247949	0.175922776	-0.622280208	0.189539418	
v15	-0.018171005	0.137361473	-0.626467433	-0.222620305	-0.095614933	1
v16	0.008299134	0.003681332	-0.001922213	0.058845427	0.456395895	
v17	0.338163494	0.024709420	0.026968695	-0.092641106	0.021017583	
v18	-0.223381169	-0.035196073	0.174749591	-0.556617395	0.129377827	

Figure 4. Intra-cluster Correlation between performance measurement indices in MG

In the PC1 of the EP kernel, v7 (PktIn: Received packets / sec) is the most relevant, while PC2 element has a high correlation of v0 (cpu: The CPU number) and v2 (inter: the interrupt summary stats). PC1 is related to packet reception because it is highly associated to v7 (PktIn: Received packets / sec.) And PC2 is related to v0 (cpu: The CPU number) and v2 (inter: the interrupt summary stats).

In the PC1 of EP kernel, v0 (usr: cpu usage by a user processes), v8 (used: amount of used memory), v13 (run: number of running processes) as well as v17 (int: number of interrupt). That is, PC1 is a component related to the amount of CPU used and memory used by the user, the amount of cache memory used, the total number of bytes received on the network, the number of processes being executed, and the number of interrupts.
(a) Collectl - EP

	PC1	PC2	PC3	PC4	PC 5
VO	0.013338535	0.588570148	-0.32363833	0.04385292	-0.19005135
V1	0.001825800	0.009880875	-0.48204305	-0.25004912	0.47776974
v2	0.014008944	0.581348570	-0.34794218	0.02519467	-0.16723928
V3	0.001955741	-0.203212057	-0.39893204	-0.31091424	0.44082544
v4	-0.001102609	0.140832496	0.27517466	-0.62694176	-0.14497469
v5	0.004150302	0.189652001	0.25654581	-0.62703406	-0.03826693
V6	0.073950048	-0.018993828	-0.02658176	-0.02474585	0.01058389
v7	0.996676607	-0.003159324	0.02424902	0.00940449	0.02114825
v8	0.011534157	-0.357274360	-0.38261862	-0.18161025	-0.44822224
v9	0.025185874	-0.300680662	-0.30570205	-0.13542214	-0.53793684
		(b) (dstate - EP		
Pot at	ion (n x k) -	(10 × 10).			

		(
	PC1	PC2	PC3	PC4	PC 5	
VO	0.3535180488	-0.068370302	0.016271750	0.0329034964	-0.0044875581	
V1	0.0648474530	-0.107887593	-0.262654738	0.1327830721	0.0268766392	
v2	-0.3532747473	0.070849505	-0.011466857	-0.0420959851	0.0025731747	÷
v3	-0.0132003841	0.091914338	0.566153743	0.3068098376	0.2270036422	
v4	-0.0006258241	0.003545678	-0.009441324	-0.0500845798	0.0686855874	
v5	0.2123360239	-0.340537999	0.084262436	0.0262063756	0.0481248784	÷
v6	-0.0437880145	-0.446421786	0.072458955	0.0022808139	0.0318096169	1
v7	-0.2647408221	-0.791171604	0.089624905	-0.0210155362	0.0463737606	•
v8	0.3544472107	-0.058590521	0.012625957	0.0398302595	-0.0047356551	
v9	0.3229929699	-0.027236725	-0.061000607	-0.0004651386	-0.0216479416	,
v10	-0.3534438459	0.063607694	-0.010934749	-0.0383351290	0.0070818608	
v11	0.0351364139	0.038206492	-0.235127569	-0.1991847787	0.6033735126	•
v12	0.0196030190	0.011890373	-0.099419984	-0.1436249107	0.6909476304	•
v13	0.3537286524	-0.066099681	0.005185501	0.0427143835	-0.0062035144	
v14	-0.0081817486	-0.034154823	-0.316119866	0.6338355596	0.0652116192	
v15	-0.0132003841	0.091914338	0.566153743	0.3068098376	0.2270036422	,
v16	0.0031781898	-0.015347898	-0.031331386	-0.0198430851	0.1841090174	
v17	0.3530585610	-0.070330786	0.004341450	0.0520526911	0.0004589051	
v18	-0.1490843620	-0.002562233	-0.314571711	0.5635362382	0.0761358824	

Figure 5. Intra-cluster Correlation between performance measurement indices in EP

In the member properties of PC2 in EP kernel, there are v2 (idl: the number of idle processes), v3 (read: total number of read operations on disks.), v10(free: amount of free memory), v15(read : read I/O requests). Furthermore, they are associated to system performance indices which are associated to v5 (one_m: load average over the last 1 minute), v6 (five_m: load average over the last 1 minute), v7 (fifteen_m: load average over the last 1 minute) on the negative opposite side.





Figure 6. Correlation Analysis between performance measurement indices in EP

(Figure 6) could identify correlations through all performance measurement indices in the EP kernel. The correlation between the v0 and v2 is 1 in (Figure 6(a)), then it means that the CPU number and the interrupt summary stats. The correlation between v8 and v9 is 0.69, and v6 and v7 are 0.86. This shows that there is some relationship in case of output and input process. In case of (Figure 6(b)), v0, v2, v8, v13 have strong correlations. Then it means that cpu usage by a user processes is related to the number of idle process, amount of used memory and running processes.

As a result, we should deduct the correlations associated with each other in clusters. In other words, it is possible to recognize the association between performance issues during the execution of applications with different characteristics. Furthermore, we suggest policies for the decision support with machine learning approach, especially PCA algorithm.

IV. CONCLUSION

With the development of HPC technology, effective use of many-core processors is emerging as an important research issue. Therefore, it is important to perform tasks in parallel processing using multiple cores. In this paper, NAS parallel benchmark application is performed and system profiling data is collected. Clustering is performed on the collected system profiling data with the PCA algorithm. Furthermore, we analyzed the relationships in the respective clusters. In addition, indicators associated to each performance index are derived. It is important to identify clustered performance indicators based on machine learning in many-core systems. As a result, our approach could support decision to present composite indicators through correlation analysis. In the future, we expect to be able to support decision of performance optimization model in designing new scale-out architecture in high performance computers.

ACKNOWLEDGMENT

This work was supported by Korea Institute of Science and Technology Information(KISTI) grant funded by the Korea government(MSIT) (No.K-18-L12-C07-S01).

REFERENCES

- Oliker, L., Biswas, R., Borrill, J., Canning, A., Carter, J., Djomehri, M. J., ... & Skinner, D. (2004, June). A performance evaluation of the Cray X1 for scientific applications. In International Conference on High Performance Computing for Computational Science (pp. 51-65). Springer, Berlin, Heidelberg.
- [2] Sodani, Avinash. "Knights landing (knl): 2nd generation intel® xeon phi processor." Hot Chips 27 Symposium (HCS), 2015 IEEE. IEEE, 2015.
- [3] Park, G., Rho, S., Kim, J. S., & Nam, D. Towards optimal scheduling policy for heterogeneous memory architecture in many-core system. Cluster Computing, 1-13
- [4] Barnes, T., Cook, B., Deslippe, J., Doerfler, D., Friesen, B., He, Y., ... & Oliker, L. (2016, November). Evaluating and optimizing the NERSC workload on knights landing. International Workshop on Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems (PMBS), (pp. 43-53).
- [5] Park, B. H., Hukerikar, S., Adamson, R., & Engelmann, C. (2017, September). Big Data Meets HPC Log Analytics: Scalable Approach to

:

Understanding Systems at Extreme Scale. 2017 IEEE International Conference on Cluster Computing (CLUSTER), (pp. 758-765).

- [6] Fu, X., Ren, R., Zhan, J., Zhou, W., Jia, Z., & Lu, G. (2012, October). LogMaster: mining event correlations in logs of large-scale cluster systems. 2012 IEEE 31st Symposium on Reliable Distributed Systems (SRDS), (pp. 71-80).
- [7] Agrawal, Rakesh, and Ramakrishnan Srikant. "Fast algorithms for mining association rules." Proc. 20th int. conf. very large data bases, VLDB. Vol. 1215. 1994.
- [8] Xing, Fei, Haihang You, and Charngda Lu. "HPC benchmark assessment with statistical analysis." Procedia Computer Science 29 (2014): 210-219.
- [9] NAS Parallel Benchmarks https://www.nas.nasa.gov/publications/npb.html

Blended Learning Design for Computer Programming Courses

Jieming Ma

Abstract—Most of virtual learning environments and online judge systems are content centric and they do not provide an active and personalized learning. This paper presents a blended learning pedagogical approach which combines instructor-led instructions, printed instructions, virtual learning environment, and online judge system to achieve maximum teaching objectives and learning outcomes. The developed blended learning system provides a virtual learning environment providing all the resources and applications as services anywhere at anytime. In addition, the system allows teachers to aggregate data by multiple resources to formulate hypotheses about strategies to raise student achievement.

Index Terms—Blended learning, online judge, virtual learning system.

I. INTRODUCTION

Programming is one of the major skills that computer science and engineering students need. In recent years, it also has been emphasized in other majors such as physics, biology, economics, etc. Practical experience is crucial for students to be proficient and to develop computational thinking skills, which are quite different from the ways of thinking where exam writing is emphasized [1].

Instructor-led classroom learning (ILCL) allows for students and instructors to interact and discuss the traditional training materials, but there are a few downsides from the learner's perspective. The ILCL typically takes a long time and the designated class time makes it difficult for certain students [2]. Moreover, the lack of practice in the classroom may reduce the learning efficiency of programming courses.

The recent information technology enabled online learning to address the difficulties posed by the limitations of the ILCL. When students engaged in online learning, they can study anywhere at any time. However, researchers have shown that students attending online courses don't have enough opportunities to interact with peers and classroom, and eventually they drop out at substantially higher rates than their counterparts in on-campus [3].

Blended learning supplements classroom teaching with Internet support to stimulate discussion and extend the range of learning materials [4]. The provision of differentiated instruction provides an opportunity to give personalized instruction to every student, which caters their needs. The online materials can be accessed easily. In addition, blended learning gives students a chance to communicate with teachers using modern communication technology, such as videoconferencing and teaching forum.

Since a certain amount of effective programming exercises is of importance in learning programming, an efficient learning platform is required to support self-regulation [5]. With the aim of increasing the flexibility of learning time and place, and providing individual learning opportunities for students, a blended learning system is developed for students by combing the virtual learning environment and online judge system. The system will assist instructors and motivate students to practice more and to go beyond the theoretical aspects they learn in class and, as a result, help them sharpen their algorithmically and programming skills.

II. BLENDED LEARNING DESIGN

Blended learning is regarded as the most effective teaching way for students [6]. The proposed pedagogical approach provides the blend of various teaching methods in preference to the learning styles of the students. Four main modules are involved in the approach, including instructor-led instructions, printed instructions, virtual learning environment, and online judge system. Fig. 1 shows the designed blended learning architecture.



Fig. 1. The Blended Learning Architecture.

This work was supported in part by the XJTLU Teaching Development Fund under Grant TDF 17/18-R16-111.

J. Ma is with the Department of Computer Science and Software Engineering, Xi'an Jiaotong-Liverpool University, Suzhou, 215123, P.R. China (e-mail: Jieming.ma@xjtlu.edu.cn).

A. Virtual Learning environment with online judges

New technology has enabled more frequent and varied online assessment measures. An online judge system automatically evaluates programs submitted by the students based on pre-set conditions and giving feedback to them in a fast manner [7,8]. By combining ongoing information on students' activities and needs, a more personalized virtual learning environment is developed with online judges. A question set is maintained by instructors. The system can compile and execute the code submitted by students and test them with pre-defined data. Error messages will appear when syntax or time/memory limit issues are raised. The students are given an initial account in the online judge system and their activity data are automatically collected into servers. This process makes it possible to trace the trail of individual students and thus to better understand students' behaviors. Thus, the proposed online judge system offers important benefits: immediate feedback, objectivity and consistency in the evaluation.

Besides online judges, the virtual learning environment provides an online learning platform allowing students access to many useful facilities and information anywhere at any time. This includes access to quiz, coursework, past exam papers, chat, forum, podcast, timetables, questionnaire, personal files, etc. These online tools enable students to participate in education remotely and to communicate online without traditional face-to-face meetings. Fig. 2 shows the screen shot of the developed Xi'an Jiaotong-Liverpool University (XJTLU) blended learning system.

-	Online Classroom	Online Judge	Online Contests	Authors
	Module	Guild	News	Register
	Coursework	Problem set	Current Contest	Update Your Info
	Quiz	Submit Solutions	Past Contests	Rank list
	Forum	Judge Status	Rules	
		Search Problems		
	Welcome to t	the XJTLU Blend System	ed Learning	
	Welcome to t Xi'an Jiaotong-Liverpoo (BLS) is an online platfo and teaching at XJTUJ, up-to-date with your m assignments, and much	the XJTLU Blend System I University (XJTLU) Blend in designed to support a Access XJTLU BLS anytim iodules, view important re more.	ed Learning ded Learning System and enhance learning e, anywhere and stay esources, submit	
	Welcome to t Xi'an Jiaotong-Liverpoo (BLS) is an online platfor and teaching at XITUJ. up-to-date with your m assignments, and much The Online Judge provi online versions of many	the XJTLU Blend System Il University (XITLU) Blend md designed to support a Access XITLU BLS anytim doulies, view important re more. des you an opportunity tr contests held at the XITI	ed Learning System ind enhance learning e. anywhere and stay sources, submit o participate in LU on a regular basis.	

Fig. 2. Screen shot of XJTLU Blended Learning System.

B. Instructor-led classroom learning (IICL) with online technologies

Instructor-led classroom learning (IICL) is considered a powerful teaching method, whose procedure is systematic and grounded on sound teaching and learning principles. In a programming course, students need plenty of time to practice their programming skills through in-class exercises, quizzes, assignments and labs. The proponents of online technologies as an alternative to the IICL argue that such an approach provides significant advantages for programming learning in terms of flexibility.

The XJTLU blended learning system provides the following features in various learning stages:

• Online classroom is a web-based environment that allows teachers and students to communicate, interact, collaborate, explain ideas. Teachers can auto-generate the coursework and quiz online. The submitted solutions will be marked by the online judge system.

• Online judge provides a platform where students are able to learn, practice and sharpen programming skills. Students can search and browse problems, submit solutions, and check the results in the system.

• Online contests allow students to participate in online versions of many contests held at the XJTLU. The past contest problems can be found in the online system.

The inquiry showed that the XJTLU blended learning system significantly increases the flexibility of learning time and place, and thus support more self-regulated learning in programming courses.

III. CONCLUSION

In this paper, a blended learning system has been developed to promote active and independent learning. The integration of instructor-led classroom learning and online judges allows for students to practice the contents learned in class and to explore new practical problems. The online learning features, such as quiz, coursework and online teaching materials, have been used for flexible and individual learning.

Our future work is to focus on the educational data analysis model for supporting process-oriented learning. It would be interesting to investigate identify the strengths and weaknesses of an entire class as well as individual students by using student achievement and performance data.

REFERENCES

- J. Petit, O. Giménez, and S. Roura. "Jutge.org: an educational programming judge.," in *Proc. 43rd ACM technical symposium on Computer Science Education (SIGCSE)*, 2012, pp. 445-450.
- [2] P.C.N. Subramani, *Effective Teaching and Learning*, Ashok Yakkaldevi, 2016, 224-242.
- [3] Y. Levy, "Comparing dropouts and persistence in e-learning courses," Computers & education, 2007, vol. 48, no. 2, pp. 185-204.
- [4] F. Mohammad. "Blended Learning and the Virtual Learning Environment of Nottingham Trent University", in *Proc. IEEE 2nd Int. Conf. on Developments in E-systems Engineering*, 2010, pp.295-299.
- [5] A. Robins, J. Rountree and N. Rountree, "Learning and Teaching Programming: A Review and Discussion," *Computer Science Education*, 2010, vol. 13, no. 2, pp. 137-172.
- [6] S. T. Selvi and P. Perumal, "Blended learning for programming in cloud based e-Learning system," in *Proc. IEEE Int. Conf. on Recent Trends in Information Technology*, 2012, pp. 197-201.
- [7] K. Andy, A. Lim, and B. Cheang, "Online judge." Computers & Education, 2001, vol. 36, no.4, pp. 299-315.
- [8] K. Adrian, M. Małafiejski, and T. Noiński, "Application of an online judge and contester system in academic tuition," in *Proc. the Advances in International Conference on Web Based Learning*, 2007, pp. 343-354.