# **International Journal on**

# **Advances in Telecommunications**















The International Journal on Advances in Telecommunications is published by IARIA. ISSN: 1942-2601 journals site: http://www.iariajournals.org contact: petre@iaria.org

Responsibility for the contents rests upon the authors and not upon IARIA, nor on IARIA volunteers, staff, or contractors.

IARIA is the owner of the publication and of editorial aspects. IARIA reserves the right to update the content for quality improvements.

Abstracting is permitted with credit to the source. Libraries are permitted to photocopy or print, providing the reference is mentioned and that the resulting material is made available at no cost.

Reference should mention:

International Journal on Advances in Telecommunications, issn 1942-2601 vol. 15, no. 3 & 4, year 2022, http://www.iariajournals.org/telecommunications/

The copyright for each included paper belongs to the authors. Republishing of same material, by authors or persons or organizations, is not allowed. Reprint rights can be granted by IARIA or by the authors, and must include proper reference.

Reference to an article in the journal is as follows:

<Author list>, "<Article title>" International Journal on Advances in Telecommunications, issn 1942-2601 vol. 15, no. 3 & 4, year 2022, <start page>:<end page> , http://www.iariajournals.org/telecommunications/

IARIA journals are made available for free, proving the appropriate references are made when their content is used.

Sponsored by IARIA www.iaria.org

Copyright © 2022 IARIA

# **Editors-in-Chief**

Tulin Atmaca, Institut Mines-Telecom/ Telecom SudParis, France Marko Jäntti, University of Eastern Finland, Finland

# **Editorial Advisory Board**

Ioannis D. Moscholios, University of Peloponnese, Greece Ilija Basicevic, University of Novi Sad, Serbia Kevin Daimi, University of Detroit Mercy, USA György Kálmán, Gjøvik University College, Norway Michael Massoth, University of Applied Sciences - Darmstadt, Germany Mariusz Glabowski, Poznan University of Technology, Poland Dragana Krstic, Faculty of Electronic Engineering, University of Nis, Serbia Wolfgang Leister, Norsk Regnesentral, Norway Bernd E. Wolfinger, University of Hamburg, Germany Przemyslaw Pochec, University of New Brunswick, Canada Timothy Pham, Jet Propulsion Laboratory, California Institute of Technology, USA Kamal Harb, KFUPM, Saudi Arabia Eugen Borcoci, University "Politehnica" of Bucharest (UPB), Romania Richard Li, Huawei Technologies, USA

# **Editorial Board**

Fatma Abdelkefi, High School of Communications of Tunis - SUPCOM, Tunisia Seyed Reza Abdollahi, Brunel University - London, UK Habtamu Abie, Norwegian Computing Center/Norsk Regnesentral-Blindern, Norway Rui L. Aguiar, Universidade de Aveiro, Portugal Javier M. Aguiar Pérez, Universidad de Valladolid, Spain Mahdi Aiash, Middlesex University, UK Akbar Sheikh Akbari, Staffordshire University, UK Ahmed Akl, Arab Academy for Science and Technology (AAST), Egypt Hakiri Akram, LAAS-CNRS, Toulouse University, France Anwer Al-Dulaimi, Brunel University, UK Muhammad Ali Imran, University of Surrey, UK Muayad Al-Janabi, University of Technology, Baghdad, Iraq Jose M. Alcaraz Calero, Hewlett-Packard Research Laboratories, UK / University of Murcia, Spain Erick Amador, Intel Mobile Communications, France Ermeson Andrade, Universidade Federal de Pernambuco (UFPE), Brazil Cristian Anghel, University Politehnica of Bucharest, Romania Regina B. Araujo, Federal University of Sao Carlos - SP, Brazil Pasquale Ardimento, University of Bari, Italy Ezendu Ariwa, London Metropolitan University, UK Miguel Arjona Ramirez, São Paulo University, Brasil Radu Arsinte, Technical University of Cluj-Napoca, Romania Tulin Atmaca, Institut Mines-Telecom/ Telecom SudParis, France Mario Ezequiel Augusto, Santa Catarina State University, Brazil Marco Aurelio Spohn, Federal University of Fronteira Sul (UFFS), Brazil

Philip L. Balcaen, University of British Columbia Okanagan - Kelowna, Canada Marco Baldi, Università Politecnica delle Marche, Italy Ilija Basicevic, University of Novi Sad, Serbia Carlos Becker Westphall, Federal University of Santa Catarina, Brazil Mark Bentum, University of Twente, The Netherlands David Bernstein, Huawei Technologies, Ltd., USA Eugen Borcoci, University "Politehnica" of Bucharest (UPB), Romania Fernando Boronat Seguí, Universidad Politecnica de Valencia, Spain Christos Bouras, University of Patras, Greece Martin Brandl, Danube University Krems, Austria Julien Broisin, IRIT, France Dumitru Burdescu, University of Craiova, Romania Andi Buzo, University "Politehnica" of Bucharest (UPB), Romania Shkelzen Cakaj, Telecom of Kosovo / Prishtina University, Kosovo Enzo Alberto Candreva, DEIS-University of Bologna, Italy Rodrigo Capobianco Guido, São Paulo State University, Brazil Hakima Chaouchi, Telecom SudParis, France Silviu Ciochina, Universitatea Politehnica din Bucuresti, Romania José Coimbra, Universidade do Algarve, Portugal Hugo Coll Ferri, Polytechnic University of Valencia, Spain Noel Crespi, Institut TELECOM SudParis-Evry, France Leonardo Dagui de Oliveira, Escola Politécnica da Universidade de São Paulo, Brazil Kevin Daimi, University of Detroit Mercy, USA Gerard Damm, Alcatel-Lucent, USA Francescantonio Della Rosa, Tampere University of Technology, Finland Chérif Diallo, Consultant Sécurité des Systèmes d'Information, France Klaus Drechsler, Fraunhofer Institute for Computer Graphics Research IGD, Germany Jawad Drissi, Cameron University, USA António Manuel Duarte Nogueira, University of Aveiro / Institute of Telecommunications, Portugal Alban Duverdier, CNES (French Space Agency) Paris, France Nicholas Evans, EURECOM, France Fabrizio Falchi, ISTI - CNR, Italy Mário F. S. Ferreira, University of Aveiro, Portugal Bruno Filipe Margues, Polytechnic Institute of Viseu, Portugal Robert Forster, Edgemount Solutions, USA John-Austen Francisco, Rutgers, the State University of New Jersey, USA Kaori Fujinami, Tokyo University of Agriculture and Technology, Japan Shauneen Furlong, University of Ottawa, Canada / Liverpool John Moores University, UK Emiliano Garcia-Palacios, ECIT Institute at Queens University Belfast - Belfast, UK Ana-Belén García-Hernando, Universidad Politécnica de Madrid, Spain Bezalel Gavish, Southern Methodist University, USA Christos K. Georgiadis, University of Macedonia, Greece Mariusz Glabowski, Poznan University of Technology, Poland Katie Goeman, Hogeschool-Universiteit Brussel, Belgium Hock Guan Goh, Universiti Tunku Abdul Rahman, Malaysia Pedro Gonçalves, ESTGA - Universidade de Aveiro, Portugal Valerie Gouet-Brunet, Conservatoire National des Arts et Métiers (CNAM), Paris Christos Grecos, University of West of Scotland, UK Stefanos Gritzalis, University of the Aegean, Greece William I. Grosky, University of Michigan-Dearborn, USA Vic Grout, Glyndwr University, UK Xiang Gui, Massey University, New Zealand Huagun Guo, Institute for Infocomm Research, A\*STAR, Singapore

Song Guo, University of Aizu, Japan Kamal Harb, KFUPM, Saudi Arabia Ching-Hsien (Robert) Hsu, Chung Hua University, Taiwan Javier Ibanez-Guzman, Renault S.A., France Lamiaa Fattouh Ibrahim, King Abdul Aziz University, Saudi Arabia Theodoros Iliou, University of the Aegean, Greece Mohsen Jahanshahi, Islamic Azad University, Iran Antonio Jara, University of Murcia, Spain Carlos Juiz, Universitat de les Illes Balears, Spain Adrian Kacso, Universität Siegen, Germany György Kálmán, Gjøvik University College, Norway Eleni Kaplani, University of East Anglia-Norwich Research Park, UK Behrouz Khoshnevis, University of Toronto, Canada Ki Hong Kim, ETRI: Electronics and Telecommunications Research Institute, Korea Atsushi Koike, Seikei University, Japan Ousmane Kone, UPPA - University of Bordeaux, France Dragana Krstic, University of Nis, Serbia Archana Kumar, Delhi Institute of Technology & Management, Haryana, India Romain Laborde, University Paul Sabatier (Toulouse III), France Massimiliano Laddomada, Texas A&M University-Texarkana, USA Wen-Hsing Lai, National Kaohsiung First University of Science and Technology, Taiwan Zhihua Lai, Ranplan Wireless Network Design Ltd., UK Jong-Hyouk Lee, INRIA, France Wolfgang Leister, Norsk Regnesentral, Norway Elizabeth I. Leonard, Naval Research Laboratory - Washington DC, USA Richard Li, Huawei Technologies, USA Jia-Chin Lin, National Central University, Taiwan Chi (Harold) Liu, IBM Research - China, China Diogo Lobato Acatauassu Nunes, Federal University of Pará, Brazil Andreas Loeffler, Friedrich-Alexander-University of Erlangen-Nuremberg, Germany Michael D. Logothetis, University of Patras, Greece Renata Lopes Rosa, University of São Paulo, Brazil Hongli Luo, Indiana University Purdue University Fort Wayne, USA Christian Maciocco, Intel Corporation, USA Dario Maggiorini, University of Milano, Italy Maryam Tayefeh Mahmoudi, Research Institute for ICT, Iran Krešimir Malarić, University of Zagreb, Croatia Zoubir Mammeri, IRIT - Paul Sabatier University - Toulouse, France Herwig Mannaert, University of Antwerp, Belgium Michael Massoth, University of Applied Sciences - Darmstadt, Germany Adrian Matei, Orange Romania S.A, part of France Telecom Group, Romania Natarajan Meghanathan, Jackson State University, USA Emmanouel T. Michailidis, University of Piraeus, Greece Ioannis D. Moscholios, University of Peloponnese, Greece Djafar Mynbaev, City University of New York, USA Pubudu N. Pathirana, Deakin University, Australia Christopher Nguyen, Intel Corp., USA Lim Nguyen, University of Nebraska-Lincoln, USA Brian Niehöfer, TU Dortmund University, Germany Serban Georgica Obreja, University Politehnica Bucharest, Romania Peter Orosz, University of Debrecen, Hungary Patrik Österberg, Mid Sweden University, Sweden Harald Øverby, ITEM/NTNU, Norway

Tudor Palade, Technical University of Cluj-Napoca, Romania Constantin Paleologu, University Politehnica of Bucharest, Romania Stelios Papaharalabos, National Observatory of Athens, Greece Gerard Parr, University of Ulster Coleraine, UK Ling Pei, Finnish Geodetic Institute, Finland Jun Peng, University of Texas - Pan American, USA Cathryn Peoples, University of Ulster, UK Dionysia Petraki, National Technical University of Athens, Greece Dennis Pfisterer, University of Luebeck, Germany Timothy Pham, Jet Propulsion Laboratory, California Institute of Technology, USA Roger Pierre Fabris Hoefel, Federal University of Rio Grande do Sul (UFRGS), Brazil Przemyslaw Pochec, University of New Brunswick, Canada Anastasios Politis, Technological & Educational Institute of Serres, Greece Adrian Popescu, Blekinge Institute of Technology, Sweden Neeli R. Prasad, Aalborg University, Denmark Dušan Radović, TES Electronic Solutions, Stuttgart, Germany Victor Ramos, UAM Iztapalapa, Mexico Gianluca Reali, Università degli Studi di Perugia, Italy Eric Renault, Telecom SudParis, France Leon Reznik, Rochester Institute of Technology, USA Joel Rodrigues, Instituto de Telecomunicações / University of Beira Interior, Portugal David Sánchez Rodríguez, University of Las Palmas de Gran Canaria (ULPGC), Spain Panagiotis Sarigiannidis, University of Western Macedonia, Greece Michael Sauer, Corning Incorporated, USA Marialisa Scatà, University of Catania, Italy Zary Segall, Chair Professor, Royal Institute of Technology, Sweden Sergei Semenov, Broadcom, Finland Dimitrios Serpanos, University of Patras and ISI/RC Athena, Greece Adão Silva, University of Aveiro / Institute of Telecommunications, Portugal Pushpendra Bahadur Singh, MindTree Ltd, India Mariusz Skrocki, Orange Labs Poland / Telekomunikacja Polska S.A., Poland Leonel Sousa, INESC-ID/IST, TU-Lisbon, Portugal Cristian Stanciu, University Politehnica of Bucharest, Romania Liana Stanescu, University of Craiova, Romania Cosmin Stoica Spahiu, University of Craiova, Romania Young-Joo Suh, POSTECH (Pohang University of Science and Technology), Korea Hailong Sun, Beihang University, China Jani Suomalainen, VTT Technical Research Centre of Finland, Finland Fatma Tansu, Eastern Mediterranean University, Cyprus Ioan Toma, STI Innsbruck/University Innsbruck, Austria Božo Tomas, HT Mostar, Bosnia and Herzegovina Piotr Tyczka, ITTI Sp. z o.o., Poland John Vardakas, University of Patras, Greece Andreas Veglis, Aristotle University of Thessaloniki, Greece Luís Veiga, Instituto Superior Técnico / INESC-ID Lisboa, Portugal Calin Vladeanu, "Politehnica" University of Bucharest, Romania Benno Volk, ETH Zurich, Switzerland Krzysztof Walczak, Poznan University of Economics, Poland Krzysztof Walkowiak, Wroclaw University of Technology, Poland Yang Wang, Georgia State University, USA Yean-Fu Wen, National Taipei University, Taiwan, R.O.C. Bernd E. Wolfinger, University of Hamburg, Germany Riaan Wolhuter, Universiteit Stellenbosch University, South Africa

Yulei Wu, Chinese Academy of Sciences, China Mudasser F. Wyne, National University, USA Gaoxi Xiao, Nanyang Technological University, Singapore Bashir Yahya, University of Versailles, France Abdulrahman Yarali, Murray State University, USA Mehmet Erkan Yüksel, Istanbul University, Turkey Pooneh Bagheri Zadeh, Staffordshire University, UK Giannis Zaoudis, University of Patras, Greece Liaoyuan Zeng, University of Electronic Science and Technology of China, China Rong Zhao , Detecon International GmbH, Germany Zhiwen Zhu, Communications Research Centre, Canada Martin Zimmermann, University of Applied Sciences Offenburg, Germany Piotr Zwierzykowski, Poznan University of Technology, Poland

# CONTENTS

pages: 23 - 41 Reliability of Erasure-Coded Storage Systems with Latent Errors Ilias Iliadis, IBM Research Europe – Zurich, Switzerland

pages: 42 - 51 Digital Sensing Platform with High Accuracy Time Synchronization Function of Vibration and Camera Sensors Narito Kurata, Tsukuba University of Technology, Japan

pages: 52 - 59 Combined Algorithm for Voronoi Diagram Construction as it Applies to Dynamic Ride Sharing Anton Butenko, University of Oldenburg, Germany Jorge Marx Gómez, University of Oldenburg, Germany

pages: 60 - 69 Base Station Assisted (BSA) Reinforcement Learning for Resource Allocation in Wireless Industrial Environment Idayat O. Sanusi, University of Greenwich, UK Karim M. Nasr, University of Greenwich, UK

pages: 70 - 80 Linking Radio Access Network QoE and QoS with Ensemble Multiple Regression Adrien Schaffner, LivingObjects, France Louise Travé-Massuyès, LAAS-CNRS & ANITI, University of Toulouse, France Simon Pachy, LivingObjects, France Bertrand Le Marec, LivingObjects, France

# Reliability of Erasure-Coded Storage Systems with Latent Errors

Ilias Iliadis IBM Research Europe – Zurich 8803 Rüschlikon, Switzerland email: ili@zurich.ibm.com

Abstract—Large-scale storage systems employ erasure-coding redundancy schemes to protect against device failures. The adverse effect of latent sector errors on the Mean Time to Data Loss (MTTDL) and the Expected Annual Fraction of Data Loss (EAFDL) reliability metrics is evaluated. A theoretical model capturing the effect of latent errors and device failures is developed, and closed-form expressions for the metrics of interest are derived. The MTTDL and EAFDL of erasure-coded systems are obtained analytically for (i) the entire range of bit error rates, (ii) the symmetric, clustered, and declustered data placement schemes, and (iii) arbitrary device failure and rebuild time distributions under network rebuild bandwidth constraints. For realistic values of sector error rates, the results obtained demonstrate that MTTDL degrades whereas, for moderate erasure codes, EAFDL remains practically unaffected. It is demonstrated that, in the range of typical sector error rates and for very powerful erasure codes, EAFDL degrades as well. It is also shown that the declustered data placement scheme offers superior reliability.

Keywords-Storage; Unrecoverable or latent sector errors; Reliability analysis; MTTDL; EAFDL; RAID; MDS codes; stochastic modeling.

# I. INTRODUCTION

Today's large-scale data storage systems and most cloud offerings recover data lost due to device and component failures by deploying efficient erasure coding schemes that provide high data reliability [1]. The replication schemes and the Redundant Arrays of Inexpensive Disks (RAID) schemes, such as RAID-5 and RAID-6, which have been deployed extensively in the past thirty years [2-5] are special cases of erasure codes. Modern storage systems though use advanced, more powerful erasure coding schemes. The effectiveness of these schemes has been evaluated based on the Mean Time to Data Loss (MTTDL) [2-11] and, more recently, the Expected Annual Fraction of Data Loss (EAFDL) reliability metrics [12-16]. The latter metric was introduced, because Amazon S3 [17], Facebook [18], LinkedIn [19] and Yahoo! [20] consider the amount of lost data measured in time.

The reliability level achieved depends not only on the particular choice of the erasure coding scheme, but also on the way data is placed on storage devices. The reliability assessment presented in [4] demonstrated that, for a replication factor of three, a declustered data placement scheme achieves a superior reliability than other placement schemes. The declustered placement scheme ensures that codewords are spread equally across devices. This is the scheme that was originally used by Google [21], Facebook [22], and Microsoft<sup>®</sup> Azure<sup>1</sup> [23], but, to improve data reliability and storage efficiency further, today they use erasure coding schemes that offer higher efficiency [24-26].

The reliability of storage systems is further degraded by the occurrence of unrecoverable sector errors, that is, errors that can be corrected neither by the standard sector-associated error correction code (ECC) nor by the re-read mechanism of hard-disk drives (HDDs). These sector errors are latent, because their existence is only discovered when there is an attempt to access them. Once an unrecoverable or latent sector error is detected, it can usually be corrected by the erasure coding capability. However, if this is not feasible, it is permanently lost, leading to an unrecoverable failure. Consequently, unrecoverable errors do not necessarily lead to unrecoverable failures. Permanent losses of data due to latent errors are quite pronounced in higher-capacity HDDs and storage nodes, because of the higher frequency of their occurrence [27-30]. The risk of permanent loss of data rises in the presence of latent errors.

Previous works have shown that actual latent-error rates degrade MTTDL by orders of magnitude [8][11][28][30]. Does this also apply to the case of the EAFDL metric given that, when a data loss occurs, the amount of sectors lost due to latent errors is much smaller than the amount of data lost due to a device failure? What is the range of error rates that cause EAFDL to deteriorate? This article addresses these critical questions.

Analytical results for the MTTDL and EAFDL metrics in the context of general erasure-coded storage systems, but in the absence of latent errors, were obtained in [12-14]. The first analytical assessment of EAFDL in the presence of latent errors was presented in [15] for the case of RAID-5 systems by presenting a comprehensive theoretical stochastic model that captures all the details of the rebuild process. This model was subsequently extended to a significantly more complex one for the case of RAID-6 systems [16]. Clearly, extending this model further to assess EAFDL in the presence of latent errors for arbitrary erasure coding schemes seems to be a daunting task because of its state explosion. The state of the model developed in this article does not explode, because it takes into account only the most significant details of the rebuild process.

To assess the reliability of erasure-coded systems, we adopt the non-Markovian methodology developed in prior work [12-

<sup>&</sup>lt;sup>1</sup>Microsoft is a trademark of Microsoft Corporation in the United States, other countries, or both.

14] to evaluate MTTDL and EAFDL of storage systems and extending it to assess the effect of latent errors. The validity of this methodology for accurately assessing the reliability of storage systems has been confirmed by simulations in several contexts [4][9][12][31]. It has been demonstrated that theoretical predictions of the reliability of systems comprising highly reliable storage devices are in good agreement with simulation results. Consequently, the emphasis of the present work is on theoretically evaluating the reliability of storage systems with latent errors.

The reliability results obtained by the model developed here are shown to be in agreement with previous specific theoretical and simulation results presented in the literature. We verify its validity by comparing the results obtained for the cases of RAID-5 and RAID-6 systems and showing that they match with those derived by the detailed models in [16]. Furthermore, we demonstrate that the model developed yields theoretical reliability results that match well with the simulation results obtained in [32], which studies the effect of erasure codes deployed in a realistic distributed storage configuration. This establishes a confidence for the model presented, the results obtained, and the conclusions drawn. The model developed is a practical one that takes into account the characteristics of latent errors observed in real systems. It is realistic, because it considers general device failure distributions including realworld ones, such as Weibull and gamma. It can also be used to assess system reliability when scrubbing is employed by applying the methodology described in [8]. This is the first work to study the effect of latent errors on EAFDL for general erasure-coded storage systems.

Note that the storage model considered in this work is relevant and realistic, because it properly captures the characteristics of erasure coding and of the rebuild process associated with the declustered data placement scheme currently used by Google [24], Microsoft<sup>®</sup> Azure [26], Facebook [33], and DELL/EMC [34]. Consequently, the theoretical results derived here are important, because they can be used to assess the reliability of the above schemes and also determine the parameter values that ensure a desired level of reliability. It can also be used to assess system reliability when scrubbing is employed by applying the methodology described in [8].

The key contributions of this article are the following. We consider the reliability of erasure-coded storage systems with latent errors that was derived analytically for the MTTDL and EAFDL reliability metrics for the entire range of sector error rates, and for the symmetric, clustered, and declustered data placement schemes [1]. In this article, we extend our previous work by also presenting results for the additional reliability metric of interest E(H), namely, the conditional expected amount of lost user data, given that data loss has occurred. We subsequently demonstrate that, in the range of typical sector-error rates, unrecoverable failures are frequent, which degrades MTTDL. However, in [1], it was shown that the relative increase of the amount of data loss is negligible, which leaves EAFDL practically unaffected in this range. In the present work, we demonstrate that this result holds for moderate erasure codes, but for very powerful erasure codes, it may not be the case. We have confirmed that reliability results obtained by the model developed here are in agreement with previous specific theoretical and simulation results

TABLE I. NOTATION OF SYSTEM PARAMETERS

Parameter	Definition	
n	number of storage devices	
c	amount of data stored on each device	
l	number of user-data symbols per codeword $(l > 1)$	
m	total number of symbols per codeword $(m > l)$	
(m, l)	MDS-code structure	
8	symbol size	
k	spread factor of the data placement scheme, or	
	group size (number of devices in a group) $(m \le k \le n)$	
b	average reserved rebuild bandwidth per device	
$B_{\rm max}$	upper limitation of the average network rebuild bandwidth	
X	time required to read (or write) an amount $c$ of data at an average	
	rate b from (or to) a device	
$F_X(.)$	cumulative distribution function of $X$	
$F_{\lambda}(.)$	cumulative distribution function of device lifetimes	
$P_{\rm bit}$	probability of an unrecoverable bit error	
seff	storage efficiency of redundancy scheme $(s_{\text{eff}} = l/m)$	
U	amount of user data stored in the system $(U = s_{\text{eff}} n c)$	
$\tilde{r}$	MDS-code distance: minimum number of codeword symbols lost	
	that lead to permanent data loss	
	$(\tilde{r} = m - l + 1 \text{ and } 2 \leq \tilde{r} \leq m)$	
$f_X(.)$	probability density function of $X$ ( $f_X(.) = F'_X(.)$ )	
C	number of symbols stored in a device $(C = c/s)$	
$\mu^{-1}$	mean time to read (or write) an amount $c$ of data at an average rate	
	b from (or to) a device $(\mu^{-1} = E(X) = c/b)$	
$\lambda^{-1}$	mean time to failure of a storage device	
	$(\lambda^{-1} = \int_0^\infty [1 - F_\lambda(t)] dt)$	
$P_s$	probability of an unrecoverable sector (symbol) error	
$P_{\rm DL}$	probability of data loss during rebuild	
$P_{\rm UF}$	probability of data loss due to unrecoverable failures during rebuild	
$P_{ m DF}$	probability of data loss due to a disk failure during rebuild	
Q	amount of lost user data during rebuild	
H	amount of lost user data, given that data loss has occurred during	
	rebuild	
S	number of lost symbols during rebuild	

presented in the literature. We also assess the reliability of realworld erasure coding schemes employed by enterprises. The model developed provides useful insights into the benefits of the erasure coding schemes and yields results for the entire parameter space, which allows a better understanding of the design tradeoffs.

The remainder of the article is organized as follows. Section II reviews prior relevant work and analytical models presented in the literature for assessing the effect of latent errors on the reliability of erasure-coded systems. Section III describes the storage system model and the corresponding parameters considered. Section IV presents the general framework and methodology for deriving the MTTDL and EAFDL metrics analytically for the case of erasure-coded systems and in the presence of latent errors. Closed-form expressions for relevant reliability metrics are derived for the symmetric, clustered, and declustered data placement schemes. Section V presents numerical results demonstrating the adverse effect of unrecoverable or latent errors and the effectiveness of these schemes for improving system reliability. The reliability of real-world erasure coding schemes employed by enterprises to protect their stored data is assessed in Section VI. Finally, we conclude in Section VII.

#### II. RELATED WORK

The adverse effect of latent errors on the MTTDL reliability metric of RAID-5, RAID-6, replication, and erasure-coded systems has been demonstrated in [8][11][27-30]. Analytical reliability expressions for MTTDL that take into account the effect of latent errors have been obtained predominately using Markovian models, which assume that component failure and rebuild times are independent and exponentially distributed [8][11][28][29]. The effect of latent errors on MTTDL of erasure-coded storage systems for the realistic case of non-exponential failure and rebuild time distributions was assessed in [30][35] for a limited range of error rates. In this article, we consider the entire range of sector error rates and assess the effect of latent errors not only on MTTDL, but also on the amount of lost data for the realistic case of non-exponential failure and rebuild time distributions.

Disk scrubbing has been used to mitigate the adverse effect of latent errors on system reliability [8][36][37][38]. The scrubbing process identifies latent errors at an early stage and attempts to correct them before disk failures occur. This in effect reduces the probability of encountering a latent error during the rebuild process. The resulting latent-error probability was derived in [8] as a function of the scrubbing and workload parameters. Subsequently, it was shown that the reliability level achieved when scrubbing is used can be obtained from the reliability level of a system that does not use scrubbing by adjusting the probability of encountering a latent error accordingly. The methodology presented in [8] for deriving the adjusted latent error probability when scrubbing is employed is also applicable for assessing the efficiency of other scrubbing schemes, such as the adaptive scrubbing schemes proposed in [37][38]. Moreover, this methodology can also be applied in conjunction with the reliability results presented in this article to assess the reliability of erasure-coded systems when scrubbing is used.

A simulation analysis of reliability aspects of erasure-coded data centers was presented in [39]. Various configurations were considered and it was shown that erasure codes and redundancy placement affect system reliability. In [32] it was recognized that it is hard to get statistically meaningful experimental reliability results using prototypes, because this would require a large number of machines to run for years. This underscores the usefulness of the analytical reliability results derived in this article.

#### III. STORAGE SYSTEM MODEL

To assess the reliability of erasure-coded storage systems, we adopt the model used in [14] and extend it to cover the case of latent errors. The storage system comprises n storage devices (nodes or disks), where each device stores an amount c of data such that the total storage capacity of the system is n c. This does not account for the spare space used by the rebuild process.

#### A. Redundancy

User data is divided into blocks (or symbols) of a fixed size s (e.g., sector size of 512 bytes) and complemented with parity symbols to form codewords. We consider (m, l) maximum distance separable (MDS) erasure codes, which map l user-data symbols to a set of m (> l) symbols, called a codeword, having the property that any subset containing l of the m symbols can be used to reconstruct (recover) the codeword. The corresponding storage efficiency  $s_{\rm eff}$  and amount U of user data stored in the system is

$$s_{\text{eff}} = l/m$$
 and  $U = s_{\text{eff}} n c = l n c/m$ . (1)



Figure 1. Clustered and declustered placement of codewords of length m = 3 on n = 6 devices. X1, X2, X3 represent a codeword (X = A, B, C, ..., L).

Also, the number C of symbols stored in a device is

$$C = c/s . (2)$$

Our notation is summarized in Table I. The derived parameters are listed in the lower part of the table. To minimize the risk of permanent data loss, the m symbols of each codeword are spread and stored on m distinct devices. This way, the system can tolerate any  $\tilde{r} - 1$  device failures, but  $\tilde{r}$  device failures may lead to data loss, with

$$\tilde{r} = m - l + 1$$
,  $1 \le l < m$  and  $2 \le \tilde{r} \le m$ . (3)

Examples of MDS erasure codes are the replication, RAID-5, RAID-6, and Reed–Solomon schemes.

#### B. Symmetric Codeword Placement

In a symmetric placement scheme, the system effectively comprises n/k disjoint groups of k devices, and each codeword is placed entirely in one of these groups. Within each group, all  $\binom{k}{m}$  possible ways of placing m symbols across k devices are used equally to store all the codewords in that group [40]. In particular, we consider the *clustered* and *declustered* placement schemes, as shown in Figure 1, which are special cases of symmetric placement schemes with k being equal to m and n, respectively. In the case of clustered placement, the storage system comprises n/m independent groups, referred to as *clusters*. Each codeword is stored across the devices of a particular cluster. In the case of declustered placement, all  $\binom{n}{m}$  possible ways of placing m symbols across n devices are used equally to store all the codewords in the system.

#### C. Codeword Reconstruction and Rebuild Process

When storage devices fail, codewords lose some of their symbols, which immediately triggers the rebuild process.

1) Exposure Levels: The system is at exposure level u  $(0 \le u \le \tilde{r})$  when there are codewords that have lost u symbols owing to device failures, but there are no codewords that have lost more symbols. These codewords are referred to as the *most-exposed* codewords. Transitions to higher exposure levels are caused by device failures, whereas transitions to lower ones are caused by successful rebuilds. We denote by  $C_u$  the number of most-exposed codewords upon entering exposure level u,  $(u \ge 1)$ . Upon the first device failure it holds that

$$C_1 = C , \qquad (4)$$

where C is determined by (2). In Section IV, we will derive the reliability metrics of interest using the *direct path* 



Figure 2. Rebuild under declustered placement.



Figure 3. Rebuild under clustered placement.

*approximation*, which considers only transitions from lower to higher exposure levels [4][9][12][31][40]. This implies that each exposure level is entered only once.

2) Prioritized Rebuild: When a symmetric or declustered placement scheme is used, as shown in Figure 2, spare space is reserved on each device for temporarily storing the reconstructed codeword symbols before they are transferred to a new replacement device. The rebuild process to restore the data lost by failed devices is assumed to be both *prioritized* and distributed. A prioritized (or intelligent) rebuild process always attempts first to rebuild the *most-exposed* codewords, namely, the codewords that have lost the largest number of symbols [4][9][14][26][32]. At each exposure level u, it attempts to bring the system back to exposure level u-1 by recovering one of the u symbols that each of the  $C_u$  most-exposed codewords has lost by reading l of the remaining symbols. In a distributed rebuild process, the codewords are reconstructed by reading symbols from an appropriate set of surviving devices and storing the recovered symbols in the reserved spare space of these devices. During this process, it is desirable to reconstruct the lost codeword symbols on devices in which another symbol of the same codeword is not already present.

In the case of clustered placement, the codeword symbols are spread across all k (= m) devices in each group (cluster). Therefore, reconstructing the lost symbols on the surviving devices of a group would result in more than one symbol of the same codeword on the same device. To avoid this, the lost symbols are reconstructed directly in spare devices as shown in Figure 3 and described in [14].

TABLE II. NOTATION OF SYSTEM PARAMETERS AT EXPOSURE LEVELS

Parameter	Definition
u	exposure level
$C_u$	number of most-exposed codewords upon entering exposure level u
$\tilde{n}_u$	number of devices at exposure level $u$ whose failure causes an exposure level transition to level $u + 1$
$V_u$	fraction of the most-exposed codewords that have symbols stored on any given device from the $\tilde{n}_u$ devices
$R_u$	rebuild time at exposure level $u$
$\alpha_u$	fraction of the rebuild time $R_u$ still left when another device fails,
	causing the exposure level transition $u \rightarrow u + 1$
$P_{u \rightarrow u+1}$	transition probability from exposure level $u$ to $u + 1$
$b_{n}$	average rate at which recovered data is written at exposure level $u$

3) Rebuild Process: A certain portion of the device bandwidth is reserved for read/write data recovery during the rebuild process, and the remaining bandwidth is used to serve user requests. Let b denote the actual average reserved rebuild bandwidth per device. The lost symbols are rebuilt in parallel using the rebuild bandwidth available on each surviving device. Let us denote by  $b_u$  ( $\leq b$ ) the average rate at which the amount of data corresponding to the number  $C_u$  of symbols to be rebuilt at exposure level u is written to selected device(s). This rate depends on  $B_{\text{max}}$ , the upper limitation of the average network rebuild bandwidth [14]. Also, let  $1/\mu$ ,  $f_X(.)$ , and  $F_X(.)$  denote the mean, the probability density function, and the cumulative distribution function of the time X required to read (or write) an amount c of data from (or to) a device, respectively. The kth moment of X,  $E(X^k)$ , and its mean E(X) are then given by

$$E(X^k) = \int_0^\infty t^k f_X(t) dt$$
, for  $k = 1, 2, \dots$ , (5)

$$\mu^{-1} \triangleq E(X) = c/b . \tag{6}$$

4) Failure and Rebuild Time Distributions: The lifetimes of the *n* devices are assumed to be independent and identically distributed, with a cumulative distribution function  $F_{\lambda}(.)$  and a mean of  $1/\lambda$ . We consider real-world distributions, such as Weibull and gamma, as well as exponential distributions that belong to the large class defined in [31]. Note that, although the model considered here does not account for correlated device failures, their effect can be assessed by enhancing the model according to the approach presented in [8]. This issue, however, is beyond the scope of this article. The results in this article hold for *highly reliable* storage devices, which satisfy the condition [14][31]

$$\mu \int_0^\infty F_\lambda(t) [1 - F_X(t)] dt \ll 1, \quad \text{with} \quad \frac{\lambda}{\mu} \ll 1.$$
 (7)

This condition expresses the fact that the ratio of the mean time  $1/\mu$  to read all contents of a device (which typically is on the order of tens of hours) to the mean time to failure of a device  $1/\lambda$  (which is typically at least on the order of thousands of hours) is very small, and in particular the fact that it is very unlikely that a given device fails during a rebuild period.

When the devices are highly reliable, the MTTDL and EAFDL reliability metrics of erasure-coded storage systems tend to be insensitive to the device failure distribution, that is, they depend only on its mean  $1/\lambda$ , but not on its density  $F_{\lambda}(.)$ . They are, however, sensitive to the distribution  $F_X(.)$  of the device rebuild times [14].

5) Amount of Data to Rebuild and Rebuild Times at Each Exposure Level: We denote by  $\tilde{n}_u$  the number of devices at exposure level u whose failure causes an exposure level transition to level u + 1 and  $V_u$  the fraction of the  $C_u$  most-exposed codewords that have a symbol stored on any given such device. Note that  $\tilde{n}_u$  depends on the codeword placement scheme. The notation used here is summarized in Table II. Let  $R_u$  denote the rebuild time of the most-exposed codewords at exposure level u. At exposure level 1, the amount of data to be recovered is equal to c. Given that this data is recovered at an average rate of  $b_1$  and that the time required to write an amount c of data at an average rate of b is equal to X, it follows that the rebuild time  $R_1$  is given by

$$R_1 = (b/b_1) X . (8)$$

Let  $\alpha_u$  be the fraction of the rebuild time  $R_u$  still left when another device fails, causing the exposure level transition  $u \rightarrow u+1$ . In [41, Lemma 2], it was shown that, for highly reliable devices satisfying condition (7),  $\alpha_u$  is approximately uniformly distributed in (0, 1), that is

$$\alpha_u \sim U(0,1), \quad u = 1, \dots, \tilde{r} - 1.$$
 (9)

We proceed by considering that the rebuild time  $R_{u+1}$  is determined completely by  $R_u$  and  $\alpha_u$  in the same manner as in [13][14][40]. For the rebuild schemes considered, the fraction of the  $C_u$  most-exposed codewords that were not yet considered by the rebuild process upon the next device failure is roughly equal to the fraction  $\alpha_u$  of the rebuild time  $R_u$  still left. Therefore, upon the next device failure, an approximate number  $\alpha_u C_u$  of the  $C_u$  codewords were not yet considered by the rebuild process. Clearly, the fraction  $V_u$  of these codewords that have symbols stored on the newly failed device depends only on the codeword placement scheme. Consequently, the number  $C_{u+1}$  of the most-exposed codewords upon entering exposure level u + 1 is

$$C_{u+1} \approx V_u \, \alpha_u \, C_u \,, \text{ for } u = 1, \dots, \tilde{r} - 1 \,.$$
 (10)

Repeatedly applying (10) and using (4) and the convention that for any sequence  $\delta_i$ ,  $\prod_{i=1}^0 \delta_i \triangleq 1$ , yields

$$C_u \approx C \prod_{i=1}^{u-1} V_i \alpha_i$$
, for  $u = 1, \dots, \tilde{r}$ . (11)

Unconditioning (11) on  $\alpha_1, \ldots, \alpha_{u-1}$  yields

$$E(C_u) = C\left(\prod_{j=1}^{u-1} V_j\right) E\left(\prod_{j=1}^{u-1} \alpha_j \left| \text{ level } u \text{ was entered} \right),$$
  
for  $u = 1, \dots, \tilde{r}$ . (12)

6) Unrecoverable Errors: The reliability of storage systems is affected by the occurrence of unrecoverable or latent errors. Let  $P_{\rm bit}$  denote the unrecoverable bit-error probability. According to the specifications,  $P_{\rm bit}$  is equal to  $10^{-15}$  for SCSI drives and  $10^{-14}$  for SATA drives [8]. Assuming that bit errors occur independently over successive bits, the unrecoverable sector (symbol) error probability  $P_s$  is

$$P_s = 1 - (1 - P_{\rm bit})^s , \qquad (13)$$

with s expressed in bits. Assuming a sector size of 512 bytes, the equivalent unrecoverable sector error probability is  $P_s \approx$ 

 $P_{\rm bit} \times 4096$ , which is  $4.096 \times 10^{-12}$  in the case of SCSI and  $4.096 \times 10^{-11}$  in the case of SATA drives. In practice, however, and also owing to the accumulation of latent errors over time, these probability values are higher. Indeed, empirical field results suggest that the actual values can be orders of magnitude higher, reaching  $P_s \approx 5 \times 10^{-9}$  [42].

#### IV. DERIVATION OF MTTDL AND EAFDL

The MTTDL metric assesses the expected time until some data can no longer be recovered and therefore is lost forever, whereas the EAFDL assesses the fraction of stored data that is expected to be lost by the system annually. The reliability metrics are derived using the methodology presented in [12][13] [14] and extending it to assess the effect of latent errors. This methodology uses the direct path approximation [11], does not involve Markovian analysis [4][9][12][31][40], and holds for general failure time distributions, which can be exponential or non-exponential, such as the Weibull and gamma distributions that satisfy condition (7).

At any point in time, the system can be thought to be in one of two modes: normal or rebuild mode. During normal mode, all devices are operational and all data in the system has the original amount of redundancy. A *first device* failure causes a transition from normal to rebuild mode. A rebuild process attempts to restore the lost data, which eventually leads the system either to a data loss (DL) with probability  $P_{\rm DL}$  or back to the original normal mode by restoring initial redundancy, with probability  $1 - P_{\rm DL}$ . Any symbols encountered with unrecoverable or latent errors are usually corrected by the erasure coding capability. However, it may not be possible to recover multiple unrecoverable errors in a codeword, which therefore leads to data loss.

Let T be a typical interval of a fully operational period, that is, the interval from the time t that the system is brought to its original state until a subsequent first device failure occurs. For a system comprising n devices with a mean time to failure of a device  $1/\lambda$ , the expected duration of T is [12]

$$E(T) = \frac{1}{n\,\lambda}\,,\tag{14}$$

and MTTDL is

$$\text{MTTDL} \approx \frac{E(T)}{P_{\text{DL}}} = \frac{1}{n \,\lambda \, P_{\text{DL}}} \,. \tag{15}$$

The EAFDL is obtained as the ratio of the expected amount E(Q) of lost user data, normalized to the amount U of user data, to the expected duration of T [12, Eq. (9)]:

$$\mathsf{EAFDL} \approx \frac{E(Q)}{E(T) \cdot U} \stackrel{(14)}{=} \frac{n \,\lambda \, E(Q)}{U} \stackrel{(1)}{=} \frac{m \,\lambda \, E(Q)}{l \, c} , \quad (16)$$

with E(T) and  $1/\lambda$  expressed in years.

The expected conditional amount E(H) of lost user data, given that data loss has occurred, is determined by [12, Eq. (8)]:

$$E(H) = \frac{E(Q)}{P_{\rm DL}} . \tag{17}$$

It follows from (15), (16), and (17) that

$$EAFDL = \frac{E(H)}{MTTDL \cdot U} , \qquad (18)$$

with the MTTDL expressed in years.

TABLE III. NOTATION OF RELIABILITY METRICS AT EXPOSURE LEVELS

Parameter	Definition
u	exposure level
$P_u$	probability of entering exposure level $u$
$P_{\mathrm{UF}u}$	probability of data loss due to unrecoverable symbol errors at
_	exposure level u
$P_{\rm UF}$	probability of data loss due to unrecoverable symbol errors
$P_{\rm DF}$	probability of data loss due to $\tilde{r}$ successive device failures
$P_{\rm DL}$	probability of data loss
$q_u$	probability that, at exposure level $u$ , a codeword that has lost $u$
	symbols can be restored
$\hat{q}_u$	probability that, under instantaneous transitions from exposure level
-	1 to exposure level $u$ , all of the $C_u$ most-exposed codewords, which
	have lost $u$ symbols, can be restored

### A. Reliability Analysis

1

At any exposure level u ( $u = 1, \ldots, \tilde{r} - 1$ ), data loss may occur during rebuild owing to one or more unrecoverable failures, which is denoted by the transition  $u \rightarrow UF$ . Moreover, at exposure level  $\tilde{r} - 1$ , data loss occurs owing to a subsequent device failure, which leads to the transition to exposure level  $\tilde{r}$ . Consequently, the direct paths that lead to data loss are the following:

$$UF'_u$$
: the direct path of successive transitions  $1 \to 2 \to \cdots \to u \to UF$ , for  $u = 1, \dots, \tilde{r} - 1$ , and

 $\overrightarrow{DF}$ : the direct path of successive transitions  $1 \rightarrow 2 \rightarrow$  $\cdots \rightarrow \tilde{r} - 1 \rightarrow \tilde{r}$ ,

with corresponding probabilities  $P_{UF_u}$  and  $P_{DF}$ , respectively. The notation for the probabilities of the events that lead to data loss is summarized in Table III.

1) Data Loss: The probability  $P_{\rm UF}$  of data loss owing to unrecoverable failures is

$$P_{\rm UF} \approx \sum_{u=1}^{\tilde{r}-1} P_{\rm UF_u} , \qquad (19)$$

where  $P_{\text{UF}_{u}}$  denotes the probability of data loss associated with the direct path  $UF'_u$ . Also, it holds that

$$P_{\text{UF}_u} = P_u \ P_{u \to \text{UF}} \ , \ \text{ for } \ u = 1, \dots, \tilde{r} - 1 \ ,$$
 (20)

where  $P_u$  is the probability of entering exposure level u, which is derived in Appendix A as follows:

$$P_u \approx (\lambda c)^{u-1} \frac{1}{(u-1)!} \frac{E(X^{u-1})}{[E(X)]^{u-1}} \prod_{i=1}^{u-1} \frac{\tilde{n}_i}{b_i} V_i^{u-1-i} , \quad (21)$$

and  $P_{u \to UF}$  is the probability of encountering an unrecoverable failure during the rebuild process at this exposure level.

In [11], it was shown that  $P_{DL}$  is accurately approximated by the probability of all direct paths to data loss. Therefore,

$$P_{\rm DL} \approx P_{\rm DF} + \sum_{u=1}^{r-1} P_{\rm UF_u} \stackrel{(19)}{\approx} P_{\rm DF} + P_{\rm UF} .$$
 (22)

Approximate expressions for the probabilities of data loss  $P_{\rm UF_{u}}$  and  $P_{\rm DF}$  are subsequently obtained by the following proposition.

*Proposition 1:* For  $u = 1, \ldots, \tilde{r} - 1$ , it holds that

$$P_{\mathrm{UF}_{u}} \approx -(\lambda c)^{u-1} \frac{E(X^{u-1})}{[E(X)]^{u-1}} \left(\prod_{i=1}^{u-1} \frac{\tilde{n}_{i}}{b_{i}} V_{i}^{u-1-i}\right)$$
$$\cdot \log(\hat{q}_{u})^{-(u-1)} \left(\hat{q}_{u} - \sum_{i=0}^{u-1} \frac{\log(\hat{q}_{u})^{i}}{i!}\right), \quad (23)$$
here 
$$\hat{q}_{u} \triangleq q_{u}^{C} \prod_{j=1}^{u-1} V_{j}, \quad (24)$$

where

$$q_u = 1 - \sum_{j=\tilde{r}-u}^{m-u} {m-u \choose j} P_s^j (1-P_s)^{m-u-j} , \qquad (25)$$

$$P_{\rm DF} \approx (\lambda \, c)^{\tilde{r}-1} \, \frac{1}{(\tilde{r}-1)!} \, \frac{E(X^{\tilde{r}-1})}{[E(X)]^{\tilde{r}-1}} \, \prod_{i=1}^{r-1} \frac{\tilde{n}_i}{b_i} \, V_i^{\tilde{r}-1-i}, \quad (26)$$

where  $E(X^{\tilde{r}-1})$  is obtained from (5).

Proof: Equation (23) is obtained in Appendix A. Equation (26) is obtained from the fact that  $P_{\text{DF}} = P_{\tilde{r}}$  and, subsequently, from (21) by setting  $u = \tilde{r}$ .

The MTTDL metric is obtained by substituting (22) into (15) as follows:

$$\text{MTTDL} \approx \frac{1}{n \lambda \left( P_{\text{DF}} + \sum_{u=1}^{\tilde{r}-1} P_{\text{UF}_u} \right)} , \qquad (27)$$

where  $P_{\mathrm{UF}_n}$  and  $P_{\mathrm{DF}}$  are determined by (23) and (26), respectively.

2) Amount of Data Loss: We proceed to derive the amount of data loss during rebuild. Let Q, H, and S be the amount of lost user data, the conditional amount of lost user data, given that data loss has occurred, and the number of lost symbols, respectively. Let also  $Q_{\rm DF}$  and  $Q_{{\rm UF}_u}$  denote the amount of lost user data associated with the direct paths  $D\dot{F}$  and  $UF'_{u}$ , respectively. Similarly, we consider the variables  $H_{\text{DF}}$ ,  $H_{\text{UF}_u}$ ,  $S_{\text{DF}}$ , and  $S_{\text{UF}_u}$ . Then, the amount Q of lost user data is obtained by

$$Q \approx \begin{cases} H_{\rm DF}, & \text{if } \overrightarrow{DF} \\ H_{\rm UF}_u, & \text{if } \overrightarrow{UF}_u, \\ 0, & \text{otherwise}. \end{cases} \text{ for } u = 1, \dots, \widetilde{r} - 1 \qquad (28)$$

Therefore,

$$E(Q) \approx P_{\text{DF}} E(H_{\text{DF}}) + \sum_{u=1}^{\tilde{r}-1} P_{\text{UF}_u} E(H_{\text{UF}_u}) \qquad (29)$$

$$= E(Q_{\rm DF}) + \sum_{u=1}^{r-1} E(Q_{\rm UF_u})$$
(30)

$$\approx E(Q_{\rm DF}) + E(Q_{\rm UF}) ,$$
 (31)

where

$$E(Q_{\rm DF}) = P_{\rm DF} E(H_{\rm DF}) , \qquad (32)$$

$$E(Q_{\mathrm{UF}_u}) = P_{\mathrm{UF}_u} E(H_{\mathrm{UF}_u}), \text{ for } u = 1, \dots, \tilde{r} - 1$$
 (33)

$$E(Q_{\rm UF}) = P_{\rm UF} E(H_{\rm UF}) \approx \sum_{u=1}^{r-1} E(Q_{\rm UF_u}) , \qquad (34)$$

where  $Q_{\rm UF}$  denotes the amount of lost user data due to unrecoverable failures and  $H_{\rm UF}$  the conditional amount of lost

(24)

user data, given that data loss due to unrecoverable failures has occurred.

Note that the expected amount E(Q) of lost user data is equal to the product of the storage efficiency and the expected amount of lost data, where the latter is equal to the product of the expected number of lost symbols E(S) and the symbol size s. Consequently, it follows from (1) that

$$E(Q) = \frac{l}{m} E(S) s \stackrel{(2)}{=} \frac{l}{m} \frac{E(S)}{C} c.$$
 (35)

Similarly,

$$E(Q_{\rm DF}) = \frac{l}{m} E(S_{\rm DF}) s \stackrel{(2)}{=} \frac{l}{m} \frac{E(S_{\rm DF})}{C} c , \qquad (36)$$

$$E(Q_{\rm UF_u}) = \frac{l}{m} E(S_{\rm UF_u}) s \stackrel{(2)}{=} \frac{l}{m} \frac{E(S_{\rm UF_u})}{C} c.$$
(37)

Approximate expressions for the expected amounts  $E(Q_{\text{UF}_u})$  and  $E(Q_{\text{DF}})$  of lost user data are subsequently obtained by the following proposition.

*Proposition 2:* For 
$$u = 1, \ldots, \tilde{r} - 1$$
, it holds that

$$E(Q_{\mathrm{UF}_{u}}) \approx c \frac{l\tilde{r}}{m} (\lambda c)^{u-1} \frac{1}{u!} \frac{E(X^{u-1})}{[E(X)]^{u-1}} \left( \prod_{i=1}^{u-1} \frac{\tilde{n}_{i}}{b_{i}} V_{i}^{u-i} \right)$$
$$\cdot \binom{m-u}{\tilde{r}-u} P_{s}^{\tilde{r}-u}, \quad \text{for } P_{s} \ll \frac{1}{m-\tilde{r}}, \quad (38)$$

$$E(Q_{\rm DF}) \approx c \, \frac{l}{m} \, (\lambda \, c)^{\tilde{r}-1} \, \frac{1}{(\tilde{r}-1)!} \, \frac{E(X^{\tilde{r}-1})}{[E(X)]^{\tilde{r}-1}} \, \prod_{i=1}^{r-1} \frac{\tilde{n}_i}{b_i} \, V_i^{\tilde{r}-i},$$
(39)

where  $E(X^{u-1})$  and  $E(X^{\tilde{r}-1})$  are obtained from (5).

*Proof:* Equations (38) and (39) are obtained in Appendix B. Note that (39) can also be obtained from (38) by setting  $u = \tilde{r}$ , which is the same result as Eq. (25) of [14].

The EAFDL metric is obtained by substituting (30) into (16) as follows:

$$\mathsf{EAFDL} \approx \frac{m\,\lambda\,[E(Q_{\mathsf{DF}}) + \sum_{u=1}^{\tilde{r}-1} E(Q_{\mathsf{UF}_u})]}{l\,c}\,,\qquad(40)$$

where  $E(Q_{\text{UF}_u})$  and  $E(Q_{\text{DF}})$  are determined by (38) and (39), respectively.

The conditional amounts E(H),  $E(H_{DF})$ ,  $E(H_{UF_u})$ , and  $E(H_{UF})$  of lost user data, given that data loss has occurred, are obtained from (17), (32), (33), and (34), respectively.

Remark 1: From (26), (32), and (39), it follows that

$$E(H_{\rm DF}) \approx \left(\frac{l}{m} \prod_{i=1}^{\tilde{r}-1} V_i\right) c$$
 (41)

Note that when entering exposure level  $\tilde{r}$ , for each of the  $C_{\tilde{r}}$  most-exposed codewords there are  $\tilde{r}$  symbols permanently lost. Consequently, the expected number of user-data symbols permanently lost is  $C_{\tilde{r}} (l/m) \tilde{r}$ , which implies that, for a symbol size of s, the expected amount  $E(H_{\rm DF} | C_{\tilde{r}})$  of user data lost is

$$E(H_{\rm DF} | C_{\tilde{r}}) = C_{\tilde{r}} \frac{l}{m} \tilde{r} s .$$
(42)

Unconditioning (42) on  $C_{\tilde{r}}$  yields

$$E(H_{\rm DF}) = E(C_{\tilde{r}}) \frac{l}{m} \tilde{r} s .$$
(43)

Combining (41), (43), and using (2), yields

$$E(C_{\tilde{r}}) \approx \left(\prod_{i=1}^{\tilde{r}-1} V_i\right) \frac{C}{\tilde{r}} .$$
(44)

From (12) and (44), it follows that

$$E\left(\prod_{j=1}^{\tilde{r}-1} \alpha_j \left| \text{level } \tilde{r} \text{ was entered} \right) \approx \frac{1}{\tilde{r}}.$$
 (45)

Remark 2: It turns out that when a data loss has occurred, the variables  $\alpha_1, \ldots, \alpha_{\tilde{r}-1}$  are not distributed identically. More specifically, for a rebuild time  $R_u$ , the uniform distribution of  $\alpha_u$  in the interval (0, 1), given by (9), holds under the assumption that there is a failure during this rebuild period, that is, an exposure level transition  $u \to u+1$ . However, conditioning on the exposure level transitions  $u \to u+1 \to \cdots \to u' \to u'+1$  $(u' > u), \alpha_u$  is no longer uniformly distributed in (0, 1). This is due to the fact that, conditioning on the fact that additional failures occur during the rebuild times  $R_{u+1}, \ldots, R_{u'}$ , it is more likely that the  $R_{u+1}$  period is long rather than short. In this case, only  $\alpha'_u$  is uniformly distributed in (0, 1). Assuming that the system has entered exposure level u, we deduce from (45) that

$$E\left(\prod_{j=1}^{u-1} \alpha_j \middle| \text{ level } u \text{ was entered}\right) \approx \frac{1}{u}, \quad \text{for } u = 2, \dots, \tilde{r}.$$
(46)

Substituting (46) into (12) and using (4) yields

$$E(C_u) \approx \left(\prod_{i=1}^{u-1} V_i\right) \frac{C}{u}$$
, for  $u = 1, \dots, \tilde{r}$ . (47)

*Remark 3:* For small values of  $P_s$ , it holds that  $P_{\rm UF} \rightarrow 0$ and  $E(Q_{\rm UF}) \rightarrow 0$ . Therefore, from (22) and (31), for small values of  $P_s$ , it holds that  $P_{\rm DL} \rightarrow P_{\rm DF}$  and  $E(Q) \rightarrow E(Q_{\rm DF})$ . Consequently, from (17), (32), and (41), it follows that

$$E(H) \approx E(H_{\rm DF}) \approx \left(\frac{l}{m} \prod_{i=1}^{\tilde{r}-1} V_i\right) c$$
, for  $P_s \to 0$ . (48)

When  $P_s$  is extremely small and an unrecoverable failure occurs, this failure most likely occurs when rebuilding a codeword after it has lost  $\tilde{r} - 1$  of its symbols owing to  $\tilde{r} - 1$ device failures and an unrecoverable error is encountered. In this case,  $\tilde{r}$  of its symbols are lost and therefore the expected number of lost user symbols is equal to the product of the storage efficiency l/m and  $\tilde{r}$ , which implies that

$$E(H_{\rm UF}) \approx \frac{l}{m} \tilde{r} s \stackrel{(2)}{=} \frac{1}{C} \frac{l}{m} \tilde{r} c , \quad \text{for } P_s \to 0 .$$
 (49)

*Remark 4:* For large values of  $P_s$ , it holds that

$$P_{\rm DL} \approx P_{\rm UF} \approx P_{\rm UF_1} \approx 1$$
, for  $P_s \to 1$ , (50)

which, by virtue of (15), implies that

$$\text{MTTDL} \approx \frac{1}{n\,\lambda} \,, \quad \text{for } P_s \to 1 \,. \tag{51}$$

It also holds that

$$E(Q) \approx E(H) \approx l c$$
, for  $P_s \to 1$ . (52)

From (18), and using (51) and (52), it follows that

$$\text{EAFDL} \approx m \lambda$$
, for  $P_s \to 1$ . (53)

#### B. Symmetric and Declustered Placement

We consider the case  $m < k \le n$ . The special case k = m corresponding to the clustered placement scheme has to be considered separately for the reasons discussed in Section III-C2. At each exposure level u, for  $u = 1, \dots, \tilde{r} - 1$ , it holds that [13][14]

$$\tilde{n}_u^{\text{sym}} = k - u , \qquad (54)$$

$$b_u^{\text{sym}} = \frac{\min((k-u)\,b, B_{\max})}{l+1} ,$$
 (55)

$$V_u^{\text{sym}} = \frac{m-u}{k-u} \,. \tag{56}$$

The corresponding parameters  $\tilde{n}_u^{\text{declus}}$ ,  $b_u^{\text{declus}}$ , and  $V_u^{\text{declus}}$  for the declustered placement are derived from (54), (55), and (56) by setting k = n. which yields

$$\tilde{n}_u^{\text{declus}} = n - u , \qquad (57)$$

$$b_u^{\text{declus}} = \frac{\min((n-u)\,b, B_{\text{max}})}{l+1} \,, \tag{58}$$

$$V_u^{\text{declus}} = \frac{m-u}{n-u} \,. \tag{59}$$

# C. Clustered Placement

At any exposure level u  $(u = 1, ..., \tilde{r} - 1)$ , it holds that [13][14]

$$\tilde{n}_{u}^{\text{clus}} = m - u , \ b_{u}^{\text{clus}} = \min(b, B_{\text{max}}/l) , \ V_{u}^{\text{clus}} = 1 .$$
 (60)

*Remark 5:* It follows from (60) that a system is not bandwidth-constrained when  $B_{\text{max}} \ge l b$ . Then,  $b_u^{\text{clus}} = \min(b, B_{\text{max}}/l) = b$ . In the case of RAID-5 and RAID-6, it holds that m - l = 1 and m - l = 2 or, equivalently,  $\tilde{r} = 2$  and  $\tilde{r} = 3$ , respectively, such that (22), (23), and (26) yield

$$P_{\rm DL}^{\rm RAID-5} \approx (m-1) \frac{\lambda c}{b} + 1 - (1 - P_s)^{(m-1)C}$$
, (61)

which is the same result as Eq. (85) of [16] (with  $\frac{c}{b} = \frac{1}{\mu}$ ), and

$$P_{\rm DL}^{\rm RAID-6} \approx 1 - q_1^C + \left[1 + \frac{1 - q_2^C}{\log(q_2^C)}\right] (m-1) \frac{\lambda c}{b} + \frac{(m-1)(m-2)}{2} \left(\frac{\lambda c}{b}\right)^2 \frac{E(X^2)}{[E(X)]^2}, \quad (62)$$

where  $q_1, q_2$  are determined by (25). This result is in agreement with Eq. (243) of [16] (with  $\frac{c}{b} = \frac{1}{\mu}$ ). Also, (30), (38), and (39) yield

$$E(Q^{\text{RAID-5}}) \approx \frac{l}{m} (m-1) \left(\frac{\lambda c}{b} + 2P_s\right) c , \qquad (63)$$

TABLE IV. TYPICAL VALUES OF DIFFERENT PARAMETERS

Parameter	Definition	Values
n	number of storage devices	64
c	amount of data stored on each device	20 TB
s	symbol (sector) size	512 B
$\lambda^{-1}$	mean time to failure of a storage device	876,000 h
b	rebuild bandwidth per device	100 MB/s
m	symbols per codeword	16
l	user-data symbols per codeword	13, 14, 15
U	amount of user data stored in the system	1.04 to 1.2 PB
$\mu^{-1}$	time to read an amount $c$ of data at a rate	55.5 h
	b from a storage device	

which is the same result as Eq. (105) of [16] (with  $\frac{c}{b} = \frac{1}{\mu}$ ), and

$$E(Q^{\text{RAID-6}}) \approx \frac{l}{m} \frac{(m-1)(m-2)}{2} \\ \cdot \left[ \left(\frac{\lambda c}{b}\right)^2 \frac{E(X^2)}{[E(X)]^2} + 3\frac{\lambda c}{b}P_s + 3P_s^2 \right] c.$$
(64)

This result is in agreement with Eq. (264) of [16] (with  $\frac{c}{b} = \frac{1}{\mu}$ ).

# V. NUMERICAL RESULTS

Here we assess the reliability of the clustered and declustered schemes for a system comprised of n = 64 devices (disks) and protected by an erasure coding scheme with m =16, which is the codeword length used by Microsoft<sup>®</sup> Azure [26], and l = 13, 14, and 15. Each device stores an amount of c = 20 TB, which is the capacity of the latest generation of Seagate drives, and the symbol size s is equal to a sector size of 512 bytes [43].

Typical parameter values are listed in Table IV. The annualized failure rate (AFR) of HDDs for the year 2021 is in the range of 0.11% to 4.79% [44], which corresponds to a mean time to failure in the range of 180,000 h to 8,000,000 h. The parameter  $\lambda^{-1}$  is chosen to be equal to 876,000 h (100 years) that corresponds to an AFR of 1%, which is the average AFR across all drive models [44]. Considering that 35% of the maximum transfer rate of 285 MB/s [43] is allocated for recovery operations, the reserved rebuild bandwidth *b* is then equal to 100 MB/s, which yields a rebuild time of a device  $\mu^{-1} = c/b = 55.5$  h. Also, it is assumed that the maximum network rebuild bandwidth is sufficiently large ( $B_{\text{max}} \ge n b = 6.4 \text{ GB/s}$ ), that the rebuild time distribution is deterministic, such that  $E(X^k) = [E(X)]^k$ . The obtained results are accurate, because (7) is satisfied, given that  $\lambda/\mu = 6.3 \times 10^{-5} \ll 1$ .

First, we assess the reliability for the declustered placement scheme (k = n = 64). The probability of data loss  $P_{\text{DL}}$  is determined by (22) as a function of  $P_s$  and shown in Figure 4. The probabilities  $P_{\text{UF}u}$  and  $P_{\text{DF}}$  are also shown, as obtained from (23) and (26), respectively. We observe that  $P_{\text{DL}}$  increases monotonically with  $P_s$  and exhibits a number of  $\tilde{r}$  plateaus. In the interval  $[4.096 \times 10^{-12}, 5 \times 10^{-9}]$  of practical importance for  $P_s$ , which is indicated between the two vertical dashed lines, the probability of data loss  $P_{\text{DL}}$  and, by virtue of (15), the MTTDL are degraded by orders of magnitude. The normalized  $\lambda$  MTTDL measure is obtained from (15) and shown in Figure 5. Increasing the number of parities (reducing l) improves reliability by orders of magnitude.

The normalized expected amount E(Q)/c of lost user data relative to the amount of data stored in a device is



Figure 4. Probability of data loss  $P_{\text{DL}}$  vs.  $P_s$  for l = 13, 14, 15; m = 16, n = k = 64 (declustered scheme).



Figure 5. Normalized MTTDL vs.  $P_s$  for l = 13, 14, and 15; m = 16, n = k = 64 (declustered scheme).



Figure 6. Normalized amount of data loss E(Q) vs.  $P_s$  for l = 13, 14, 15; m = 16, n = k = 64 (declustered scheme).

obtained from (30) and shown in Figure 6. The normalized expected amounts  $E(Q_{\text{UF}_u})/c$  and  $E(Q_{\text{DF}})/c$  are also shown as determined by (38) and (39), respectively. The normalized EAFDL/ $\lambda$  measure is obtained from (16) and shown in Figure 7. We observe that E(Q) and EAFDL increase monotonically, but they are practically unaffected in the interval of interest, because they degrade only when  $P_s$  is much larger than the typical sector error probabilities. For the EAFDL metric too, increasing the number of parities (reducing l) results in a reliability improvement by orders of magnitude.

The normalized expected amount E(H)/c of lost user data, given that a data loss has occurred, relative to the amount of data stored in a device is obtained from (17) and shown in Figure 8. The conditional amounts  $E(H_{\rm DF})$  and  $E(H_{\rm UF})$ obtained from (32) and (34), respectively, are also shown. In contrast to the  $P_{\rm DL}$ , EAFDL, and E(Q) metrics that increase monotonically with  $P_s,$  we observe that E(H) does not do so. The reason for that is the following. As shown in Figure 4, for  $P_s \gg 10^{-14}$ , data loss is more likely to be due to sector errors than to device failures. Given that sector errors result in a negligible amount of data loss compared with the substantial data losses caused by device failures, when  $P_s$  increases over the value of  $10^{-14}$ , the conditional amount of lost data decreases. Clearly, this is reversed for high values of  $P_s$ , and the conditional amount of lost data increases.

The expected amount  $E(H_{\rm UF})$  of user data lost due to unrecoverable failures, given that such failures have occurred, is shown in Figure 8 by the dotted green line. For extremely small values of  $P_s$ , and according to Remark 3, the value of  $E(H_{\rm UF})$  corresponds to a single corrupted codeword that loses



Figure 7. Normalized EAFDL vs.  $P_s$  for l = 13, 14, and 15; m = 16, n = k = 64 (declustered scheme).



Figure 8. Normalized E(H) vs.  $P_s$  for l = 13, 14, and 15; m = 16, n = k = 64 (declustered scheme).

 $\tilde{r}$  of its symbols of which  $\tilde{r}-1$  are lost owing to device failures and one is lost owing to an unrecoverable error. Consequently,  $E(H_{\rm UF})$  is independent of  $P_s$ , as indicated by the horizontal part of the dotted green line. Let us now consider Figure 8(a). When  $P_s \gg 10^{-10}$ ,  $E(H_{\rm UF})$  increases, because there are multiple such codewords, each of which loses  $\tilde{r}$  symbols. Subsequently, for  $\tilde{r} \geq 3$  and when  $P_s \gg 10^{-8}$ , unrecoverable failures may also be caused by a single corrupted codeword that loses  $\tilde{r}$  of its symbols,  $\tilde{r} - 2$  of which are lost owing to device failures and two are lost owing to unrecoverable errors. This in turn reduces the amount of lost data in the interval  $(10^{-10}, 10^{-8})$ , as shown in Figure 8(a). Note that this interval corresponds to that of the second plateau, as shown in Figure 4(a). When  $P_s \gg 10^{-6}$ ,  $E(H_{\rm UF})$  increases again owing to the occurrence of multiple such corrupted codewords. Eventually, when  $P_s \gg 10^{-5}$ , unrecoverable failures are encountered during rebuild prior to a second device failure and are caused by corrupted codewords that lose  $\tilde{r}$  of their symbols, one of which is lost owing to the first device failure and  $\tilde{r} - 1$  are lost owing to unrecoverable errors. This in turn increases  $E(H_{\rm UF})$ , which eventually dominates  $E(H_{\rm DF})$ . Similar observations apply in the cases of Figures 8(b) and 8(c).

The reliability metrics corresponding to the clustered placement scheme (k = m = 16) are plotted in Figures 9, 10, 11, 12, and 13. We observe that the reliability achieved by the clustered data placement scheme does not reach the reliability level achieved by the declustered one.

In the cases considered, EAFDL is practically unaffected by the presence of latent errors, as shown in Figures 7 and 12. Note, however, that for larger values of  $\tilde{r}$ , EAFDL may be affected by the presence of latent errors. For example, when m = 24 and l = 12, which yields  $\tilde{r} = 13$ , and for the case of declustered placement scheme, not only MTTDL, but also EAFDL is affected, as shown in Figure 14.

The performance of certain erasure coding schemes was assessed in [32] by obtaining the probability of data loss  $P_{\rm DL}$  using a detailed distributed storage simulator. The  $P_{\rm DL}$  values corresponding to  $P_s = 4.096 \times 10^{-12}$  ( $P_{\rm bit} = 10^{-15}$ ) for two of the configurations considered are indicated by the squares in Figure 15. This figure also shows the probabilities of data loss  $P_{\rm DL}$  that correspond to these two configurations and obtained from (22) as a function of  $P_s$ . We observe that the theoretical results are in agreement with the simulation results, which confirms the validity of the model and the analytical expressions derived.

#### VI. REAL-WORLD ERASURE CODING SCHEMES

Here we assess the reliability of various practical systems that store an amount of U = 1.2 PB user data on devices (disks) whose capacity is c = 20 TB. This amount of user data can therefore be stored on U/c = 60 devices. The system comprises n devices, where n is determined using (1) as follows:

$$n = \frac{U}{c} \frac{m}{l} = 60 \frac{m}{l} .$$
(65)

Subsequently, we consider the following real-world erasure coding schemes:



Figure 9. Probability of data loss  $P_{\text{DL}}$  vs.  $P_s$  for l = 13, 14, 15; n = 64, k = m = 16 (clustered scheme).



Figure 10. Normalized MTTDL vs.  $P_s$  for l = 13, 14, and 15; n = 64, k = m = 16 (clustered scheme).







Figure 12. Normalized EAFDL vs.  $P_s$  for l = 13, 14, and 15; n = 64, k = m = 16 (clustered scheme).



Figure 13. Normalized E(H) vs.  $P_s$  for l = 13, 14, and 15; n = 64, k = m = 16 (clustered scheme).



Figure 14. Reliability metrics vs.  $P_s$  for n = k = 64 (declustered scheme), m = 24, l = 12.



Figure 15.  $P_{\text{DL}}$  vs.  $P_s$  for n = k = 210,  $\lambda^{-1} = 3$  years,  $\mu^{-1} = 34$  hours,  $\lambda/\mu = 0.0013$ , c = 15 TB, and s = 512 B. Reliability schemes: 3-way replication and MDS(14,10).

- 1) the 3-way replication (triplication) scheme that was initially used by Google's GFS, Microsoft<sup>®</sup> Azure, and Facebook. In this case, m = 3, l = 1, with a corresponding storage efficiency of  $s_{\rm eff} = 33\%$ . According to (65), this scheme requires the employment of n = 180 devices.
- 2) the RS(9,6) erasure coding scheme employed by Google's GFS as well as QFS [24, 45], which for m = 9 and l = 6 achieves a storage efficiency of  $s_{\text{eff}} = 66\%$  and requires a number of n = 90 devices.
- 3) the MDS(16,12) erasure coding scheme akin to the LRC(16,12) code used by Microsoft<sup>®</sup> Azure [26], which for m = 16 and l = 12 achieves a storage efficiency of

 $s_{\rm eff} = 75\%$  and requires a number of n = 80 devices.

4) the RS(14,10) erasure coding scheme employed by Facebook [25], which for m = 14 and l = 10 achieves a storage efficiency of  $s_{\rm eff} = 71\%$  and requires a number of n = 84 devices.

We proceed to assess the reliability of the four erasure coding schemes for two data placement configurations: a symmetric one where the system comprises 2 disjoint groups of k devices, such that k = n/2, and a declustered one, such that k = n. As we will see next, a superior reliability is achieved by employing the declustered data placement scheme.

#### A. Symmetric Data Placement

First, we assess the reliability of the 3-way replication (triplication) scheme that requires the employment of n = 180 devices, which in turn implies that each of the two groups comprises k = 90 devices. The reliability measures are obtained for the parameter values listed in Table IV and shown in Figures 16(a) through 20(a). We observe that MTTDL is significantly degraded by the presence of latent errors. In the interval  $[4.096 \times 10^{-12}, 5 \times 10^{-9}]$  of practical importance for  $P_s$ , which is indicated between the two vertical dashed lines, Figure 17(a) shows that MTTDL is degraded by three to six orders of magnitude, whereas Figure 19(a) reveals that EAFDL is practically unaffected in this range.

Second, we consider the MDS(9,6) erasure coding scheme that requires a number of n = 90 devices, which in turn implies that each of the two groups comprises k = 45 devices. The corresponding reliability measures are shown in Figures

16(b) through 20(b). Figure 17(b) shows that, in the region of interest for  $P_s$ , MTTDL is degraded by three to five orders of magnitude, whereas Figure 19(a) reveals that EAFDL is practically unaffected in this range.

Subsequently, we consider the MDS(16,12) erasure coding scheme that requires a number of n = 80 devices, which in turn implies that each of the two groups comprises k = 40devices. The corresponding reliability measures are shown in Figures 16(c) through 20(c). Figure 17(c) shows that, in the interval of practical importance for  $P_s$ , MTTDL is degraded by three to five orders of magnitude, whereas Figure 19(c) reveals that EAFDL is practically unaffected in this range.

Finally, we consider the RS(14,10) erasure coding scheme that requires n = 84 devices, which in turn implies that each of the two groups comprises k = 42 devices. The corresponding reliability measures are shown in Figures 16(d) through 20(d). Figure 17(d) shows that, in the interval of interest, MTTDL is degraded by three to five orders of magnitude, whereas Figure 19(d) reveals that EAFDL is practically unaffected in this range.

From the above, it follows that erasure coding schemes corresponding to higher values of  $\tilde{r}$  offer a higher level of reliability. Thus, the MDS(16,12) and MDS(14,10) erasure coding schemes, for which  $\tilde{r} = 5$ , offer higher levels of reliability compared with the MDS(9,6) and 3-way replication schemes, for which  $\tilde{r} < 5$ . In particular, MDS(14,10) achieves a higher reliability than that of MDS(16,10), albeit at a lower storage efficiency (71% vs. 75%).

#### B. Declustered Data Placement

Here, we assess the reliability achieved by the erasure coding schemes considered when the declustered data placement scheme is used, such that k = n. The reliability results are shown in Figures 21(a) through 25(a).

First, we assess the reliability of the 3-way replication (triplication) scheme. Comparing Figure 22(a) with Figure 17(a), we deduce that MTTDL is roughly the same. However, comparing Figure 24(a) with Figure 19(a), we deduce that EAFDL improves by one order of magnitude.

Regarding, the reliability of the MDS(9,6) coding scheme, comparing Figure 22(b) with Figure 17(b), we deduce that MTTDL improves slightly by one order of magnitude, especially at smaller values of  $P_s$ . However, Figures 24(b) and 19(b) demonstrate that EAFDL improves by two orders of magnitude.

For the MDS(16,12) coding scheme, Figures 22(c) and 17(c) show that MTTDL improves by two orders of magnitude for values around the left vertical dotted line and by one order of magnitude for values around the right vertical dotted line. Also, from Figures 24(c) and 19(c), we observe that EAFDL improves by three orders of magnitude.

Finally, the reliability improvement regarding the MDS(14,10) coding scheme is similar to that of the MDS(16,12) coding scheme. Figures 22(d) and 17(d) show that MTTDL improves by two orders of magnitude for values around the left vertical dotted line and by one order of magnitude for values around the right vertical dotted line.

Also, Figures 24(d) and 19(d) show that EAFDL improves by three orders of magnitude.

#### C. Reliability Improvement

The reliability improvement of the erasure coding schemes considered over the initial 3-way replication is shown in Figures 26 and 27 for the two data placements, respectively. Clearly, in the interval of practical importance for  $P_s$ , the MDS(14,10) erasure coding scheme achieves superior reliability for both symmetric and declustered data placement schemes.

In particular, for the symmetric data placement, Figure 26(a) demonstrates that in the interval of interest, the MDS(9,6) erasure coding scheme improves MTTDL by three orders of magnitude for  $P_s$  values around the left vertical dotted line and by four orders of magnitude for  $P_s$  values around the right vertical dotted line. The MDS(16,12) erasure coding scheme improves MTTDL by six orders of magnitude for  $P_s$ values around the left vertical dotted line and by seven orders of magnitude for  $P_s$  values around the right vertical dotted line. Also, the MDS(14,10) erasure coding scheme improves MTTDL by seven orders of magnitude for  $P_s$  values around the left vertical dotted line and by eight orders of magnitude for  $P_s$  values around the right vertical dotted line. On the other hand, Figure 26(b) demonstrates that in the interval of practical importance for  $P_s$ , the MDS(9,6), MDS(16,12), and MDS(14,10) erasure coding schemes improve EAFDL by two, five, and six orders of magnitude, respectively.

For the declustered data placement, Figure 27(a) demonstrates that in the interval of interest, the MDS(9,6) erasure coding scheme improves MTTDL by four orders of magnitude, the MDS(16,12) erasure coding scheme improves MTTDL by eight orders of magnitude, whereas the MDS(14,10) erasure coding scheme improves MTTDL by nine orders of magnitude. On the other hand, Figure 27(b) demonstrates that in the interval of practical importance for  $P_s$ , the MDS(9,6), MDS(16,12), and MDS(14,10) erasure coding schemes improve EAFDL by three, seven, and eight orders of magnitude, respectively.

Figures 26(c) and 27(c) show the ratios of the E(H) metrics for the MDS(9,6), MDS(16,12), and MDS(14,10) erasure coding schemes to the E(H) metric for the 3-way replication scheme. In the region of interest for  $P_s$ , all three erasure coding schemes result in greater amounts of lost user data, given that data loss has occurred, compared to the conditional amount of lost user data in the case of a 3-way replication. In particular, for the symmetric data placement, the MDS(9,6), MDS(16,12), and MDS(14,10) erasure coding schemes result in about 20, 82, and 33 times greater conditional amounts of lost user data for  $P_s$  values around the left vertical dotted line and in about 58, 550, and 110 times greater conditional amounts of lost user data for  $P_s$  values around the right vertical dotted line, respectively. This is due to the fact that, for small values of  $P_s$  and according to Remark 3, E(H) depends not only on the values of the m and l parameters, but also on the number nof devices in the system, which also varies. Accordingly, for the declustered data placement, Figure 27(c) demonstrates that the MDS(9,6), MDS(16,12), and MDS(14,10) erasure coding schemes result in about 9, 18, and 7 times greater conditional amounts of lost user data, respectively.



Figure 16. Probability of data loss  $P_{DL}$  vs.  $P_s$  for various MDS coding schemes; symmetric data placement with k = n/2.



Figure 17. MTTDL vs.  $P_s$  for various MDS coding schemes; symmetric data placement with k = n/2.











Figure 20. Normalized E(H) vs.  $P_s$  for various MDS coding schemes; symmetric data placement with k = n/2.



Figure 21. Probability of data loss  $P_{DL}$  vs.  $P_s$  for various MDS coding schemes; declustered data placement (k = n).















Figure 25. Normalized E(H) vs.  $P_s$  for various MDS coding schemes; declustered data placement (k = n).



Figure 26. Ratios of the MTTDL, EAFDL, and E(H) metrics for the MDS(9,6), MDS(16,12), and MDS(14,10) schemes to those corresponding to the 3-way replication scheme; symmetric data placement with k = n/2.



Figure 27. Ratios of the MTTDL, EAFDL, and E(H) metrics for the MDS(9,6), MDS(16,12), and MDS(14,10) schemes to those corresponding to the 3-way replication scheme; declustered data placement (k = n).

#### VII. CONCLUSIONS

The effect of latent sector errors on the reliability of erasure-coded data storage systems was investigated. A methodology was developed for deriving the Mean Time to Data Loss (MTTDL) and the Expected Annual Fraction of Data Loss (EAFDL) reliability metrics analytically. Closedform expressions capturing the effect of unrecoverable latent errors were obtained for the symmetric, clustered and declustered data placement schemes. We demonstrated that the declustered placement scheme offers superior reliability in terms of both metrics. We established that, for realistic unrecoverable sector error rates, MTTDL is adversely affected by the presence of latent errors, whereas EAFDL is not. We considered several real-world erasure coding schemes and demonstrated their efficiency. The analytical reliability expressions derived enable the identification of storage-efficient data placement configurations that yield high reliability.

Applying these results to assess the effect of network rebuild bandwidth constraints is a subject of further investigation. The reliability evaluation of erasure-coded systems when device failures, as well as unrecoverable latent errors are correlated is also part of future work.

# APPENDIX A

We consider the direct path  $\overrightarrow{UF_u} = 1 \rightarrow 2 \rightarrow \cdots \rightarrow u \rightarrow UF$  and proceed to evaluate  $P_{UF_u}(R_1, \vec{\alpha}_{u-1})$ , the prob-

ability of entering exposure level u through vector  $\vec{\alpha}_{u-1} \triangleq (\alpha_1, \ldots, \alpha_{u-1})$  and given a rebuild time  $R_1$ , and then encountering an unrecoverable failure during the rebuild process at this exposure level. It follows from (20) that

$$P_{\text{UF}_u}(R_1, \vec{\alpha}_{u-1}) = P_u(R_1, \vec{\alpha}_{u-2}) \cdot P_{u \to \text{UF}}(R_1, \vec{\alpha}_{u-1}).$$
(66)

It follows from Eq.(111) of [13] by setting  $\tilde{r} = u$  that

$$P_u(R_1, \vec{\alpha}_{u-2}) \approx (\lambda b_1 R_1)^{u-1} \prod_{i=1}^{u-1} \frac{\tilde{n}_i}{b_i} (V_i \, \alpha_i)^{u-1-i} \,.$$
(67)

Given that the elements of  $\vec{\alpha}_{u-2}$  are independent random variables approximately distributed according to (9), such that  $E(\alpha_i^k) \approx 1/(k+1)$ , we have

$$E\left(\prod_{i=1}^{u-1} \alpha_i^{u-1-i}\right) = \prod_{i=1}^{u-1} E(\alpha_i^{u-1-i}) \approx \prod_{i=1}^{u-1} \frac{1}{u-i} = \frac{1}{(u-1)!}.$$
(68)

Unconditioning (67) on  $\vec{\alpha}_{u-2}$  using (68) yields

$$P_u(R_1) \approx (\lambda b_1 R_1)^{u-1} \frac{1}{(u-1)!} \prod_{i=1}^{u-1} \frac{\tilde{n}_i}{b_i} V_i^{u-1-i} .$$
(69)

Unconditioning (69) on  $R_1$  and using (6) and (8) yields (21).

We now proceed to calculate  $P_{u \to UF}(R_1, \vec{\alpha}_{u-1})$ . Upon entering exposure level u, the rebuild process attempts to restore the  $C_u$  most-exposed codewords, each of which has m - u remaining symbols. Let us consider such a codeword, and let  $L_u$  be the number of symbols permanently lost and  $I_u$  be the number of symbols in the codeword with unrecoverable errors. Owing to the independence of symbol errors,  $I_u$  follows a binomial distribution with parameter  $P_s$ , the probability that a symbol has a unrecoverable error. Thus, for  $i = 0, \ldots, m-u$ ,

$$P(I_u = j) = \binom{m-u}{j} P_s^j (1-P_s)^{m-u-j},$$
(70)

$$\approx \binom{m-u}{j} P_s^j$$
, for  $P_s \ll \frac{1}{m-u-j}$ , (71)

such that

$$E(I_u) = \sum_{j=1}^{m-u} j P(I_u = j) = (m-u) P_s .$$
 (72)

Clearly, the symbols lost due to the device failures can be corrected by the erasure coding capability only if at least l of the remaining m - u symbols can be read. Thus,  $L_u = 0$  if and only if  $I_u \leq m - u - l$  or, by virtue of (3),  $I_u \leq \tilde{r} - 1 - u$ . Thus, the probability  $q_u$  that a codeword can be restored is

$$q_u = P(L_u = 0) = 1 - P(I_u > \tilde{r} - u)$$
, (73)

which, using (70), yields (25).

Note that if a codeword is corrupted, then at least one of its l user-data symbols is lost. Owing to the independence of symbol errors, codewords are independently corrupted. Consequently, the conditional probability  $P_{\text{UF}|C_u}$  of encountering an unrecoverable failure during the rebuild process of the  $C_u$  codewords is

$$P_{\text{UF}|C_u} = 1 - q_u^{C_u}$$
, for  $u = 1, \dots, \tilde{r}$ . (74)

Substituting (11) into (74) and using (24) yields

$$P_{u \to \text{UF}}(R_1, \vec{\alpha}_{u-1}) \approx 1 - q_u^C \prod_{j=1}^{u-1} V_j \alpha_j = 1 - \hat{q}_u^{\prod_{j=1}^{u-1} \alpha_j}.$$
 (75)

Substituting (75) into (66) yields

$$P_{\mathrm{UF}_u}(R_1, \vec{\alpha}_{u-1}) \approx P_u(R_1, \vec{\alpha}_{u-2}) \left[ 1 - \hat{q}_u^{\prod_{j=1}^{u-1} \alpha_j} \right].$$
 (76)

Unconditioning (76) on  $\vec{\alpha}_{u-1}$  and using (67) yields

$$P_{\mathrm{UF}_{u}}(R_{1}) \approx P_{u}(R_{1}) - (\lambda b_{1}R_{1})^{u-1} \left(\prod_{i=1}^{u-1} \frac{\tilde{n}_{i}}{b_{i}} V_{i}^{u-1-i}\right)$$
$$\cdot E_{\vec{\alpha}_{u-1}} \left[ \left(\prod_{i=1}^{u-1} \alpha_{i}^{u-1-i}\right) \hat{q}_{u} \prod_{j=1}^{u-1} \alpha_{j} \right]. \quad (77)$$

LEMMA 1: For  $\alpha_i \sim U(0,1)$  and for all  $q \in \mathbb{R}$ , it holds that

$$E\left[\left(\prod_{i=1}^{u-1} \alpha_i^{u-1-i}\right) q \prod_{i=1}^{u-1} \alpha_i\right]$$
  
=  $\frac{1}{(u-1)!} + \log(q)^{-(u-1)} \left(q - \sum_{i=0}^{u-1} \frac{\log(q)^i}{i!}\right)$ . (78)

Proof: It holds that

$$q^{\prod_{i=1}^{u-1}\alpha_i} = e^{\log(q)\prod_{i=1}^{u-1}\alpha_i} = \sum_{j=0}^{\infty} \frac{\log(q)^j (\prod_{i=1}^{u-1}\alpha_i)^j}{j!} ,$$
(79)

which implies that

$$\begin{pmatrix} \prod_{i=1}^{u-1} \alpha_i^{u-1-i} \end{pmatrix} q \prod_{i=1}^{u-1} \alpha_i \\ = \left( \prod_{i=1}^{u-1} \alpha_i^{u-1-i} \right) \left( \sum_{j=0}^{\infty} \frac{\log(q)^j (\prod_{i=1}^{u-1} \alpha_i)^j}{j!} \right) \\ = \sum_{j=0}^{\infty} \frac{\log(q)^j \prod_{i=1}^{u-1} \alpha_i^{u-1-i+j}}{j!} .$$
(80)

Consequently,

$$E\left[\left(\prod_{i=1}^{u-1} \alpha_{i}^{u-1-i}\right) q \prod_{i=1}^{u-1} \alpha_{i}\right]$$

$$= \sum_{j=0}^{\infty} \frac{\log(q)^{j} \prod_{i=1}^{u-1} E(\alpha_{i}^{u-1-i+j})}{j!}$$

$$\approx \sum_{j=0}^{\infty} \frac{\log(q)^{j} \prod_{i=1}^{u-1} \frac{1}{u-i+j}}{j!}$$

$$= \sum_{j=0}^{\infty} \frac{\log(q)^{j}}{(u-1+j)!} = \frac{1}{(u-1)!} + \sum_{j=1}^{\infty} \frac{\log(q)^{j}}{(u-1+j)!} \quad (81)$$

$$= \frac{1}{(u-1)!} + \log(q)^{-(u-1)} \sum_{i=u}^{\infty} \frac{\log(q)^{i}}{i!}$$

$$= \frac{1}{(u-1)!} + \log(q)^{-(u-1)} \left(\sum_{i=0}^{\infty} \frac{\log(q)^{i}}{i!} - \sum_{i=0}^{u-1} \frac{\log(q)^{i}}{i!}\right)$$

$$= \frac{1}{(u-1)!} + \log(q)^{-(u-1)} \left(e^{\log(q)} - \sum_{i=0}^{u-1} \frac{\log(q)^{i}}{i!}\right)$$

From (69) and (78), (77) yields

$$P_{\mathrm{UF}_{u}}(R_{1}) \approx -(\lambda b_{1}R_{1})^{u-1} \left(\prod_{i=1}^{u-1} \frac{\tilde{n}_{i}}{b_{i}} V_{i}^{u-1-i}\right)$$
$$\cdot \log(\hat{q}_{u})^{-(u-1)} \left(\hat{q}_{u} - \sum_{i=0}^{u-1} \frac{\log(\hat{q}_{u})^{i}}{i!}\right). \quad (82)$$

Unconditioning (82) on  $R_1$ , and using (6) and (8), yields (23).

#### APPENDIX B

At exposure level u, when  $I_u \ge m - u - l + 1 = \tilde{r} - u$ , the number  $L_u$  of lost symbols is  $I_u + u$ . Consequently, the expected number  $E(L_u)$  of lost symbols is

$$E(L_u) = \sum_{i=\tilde{r}-u}^{m-u} (i+u) P(I_u = i) , \qquad (83)$$

where  $P(I_u = i)$  is given by (70). Considering approximation (71), it follows that

$$E(L_u) \approx \tilde{r} \begin{pmatrix} m-u\\ \tilde{r}-u \end{pmatrix} P_s^{\tilde{r}-u}, \quad \text{for } P_s \ll \frac{1}{m-\tilde{r}}.$$
 (84)

40

The expected number  $E(S_U|C_u)$  of symbols lost due to unrecoverable failures during the rebuild of the  $C_u$  codewords at exposure level u is equal to  $C_u E(L_u)$ , which yields

$$E(S_{\rm U}|C_u) \stackrel{(84)}{\approx} C_u \,\tilde{r} \begin{pmatrix} m-u\\ \tilde{r}-u \end{pmatrix} P_s^{\tilde{r}-u}, \ P_s \ll \frac{1}{m-\tilde{r}} \,. \tag{85}$$

Substituting (11) into (85) yields

$$E(S_{\rm U}|\vec{\alpha}_{u-1}) \approx C\left(\prod_{j=1}^{u-1} V_j \alpha_j\right) \tilde{r} \binom{m-u}{\tilde{r}-u} P_s^{\tilde{r}-u} .$$
(86)

Subsequently, the expected number  $E(S_{UF_u}|R_1, \vec{\alpha}_{u-1})$  of symbols lost due to unrecoverable failures encountered during rebuild in conjunction with entering exposure level *u* through vector  $\vec{\alpha}_{u-1}$ , and given a rebuild time  $R_1$ , is determined as follows:

$$E(S_{\mathrm{UF}_{u}}|R_{1},\vec{\alpha}_{u-1}) = P_{u}(R_{1},\vec{\alpha}_{u-1}) E(S_{\mathrm{U}}|\vec{\alpha}_{u-1}) .$$
(87)

Substituting (67) and (86) into (87) yields

$$E(S_{\mathrm{UF}_{u}}|R_{1},\vec{\alpha}_{u-1}) \approx (\lambda b_{1}R_{1})^{u-1} \left[\prod_{i=1}^{u-1} \frac{\tilde{n}_{i}}{b_{i}} (V_{i} \alpha_{i})^{u-i}\right]$$
$$\cdot C \tilde{r} \binom{m-u}{\tilde{r}-u} P_{s}^{\tilde{r}-u}, \quad P_{s} \ll \frac{1}{m-\tilde{r}} . \quad (88)$$

From (68), we have that  $E(\prod_{i=1}^{u-1} \alpha_i^{u-i}) = E(\prod_{i=1}^{u} \alpha_i^{u-i}) \approx 1/u!$ . Thus, unconditioning (88) on  $\vec{\alpha}_{u-1}$  yields

$$E(S_{\mathrm{UF}_{u}}|R_{1}) \approx (\lambda b_{1}R_{1})^{u-1} \left(\prod_{i=1}^{u-1} \frac{\tilde{n}_{i}}{b_{i}} V_{i}^{u-i}\right) \frac{1}{u!} C \tilde{r}$$
$$\cdot \binom{m-u}{\tilde{r}-u} P_{s}^{\tilde{r}-u}, \quad P_{s} \ll \frac{1}{m-\tilde{r}} . \quad (89)$$

Unconditioning (89) on  $R_1$ , and using (6) and (8), yields

$$E(S_{\mathrm{UF}_{u}}) \approx (\lambda c)^{u-1} \frac{E(X^{u-1})}{[E(X)]^{u-1}} \left(\prod_{i=1}^{u-1} \frac{\tilde{n}_{i}}{b_{i}} V_{i}^{u-i}\right) \frac{1}{u!} C \tilde{r}$$
$$\cdot \binom{m-u}{\tilde{r}-u} P_{s}^{\tilde{r}-u}, \quad P_{s} \ll \frac{1}{m-\tilde{r}} . \quad (90)$$

Substituting (90) into (37) yields (38).

Remark 6: From (21), (47), (84), and (90), it follows that

$$E(S_{\text{UF}_u}) \approx P_u E(C_u) E(L_u) . \tag{91}$$

Upon entering exposure level u, the expected number  $E(S_U|C_u)$  of symbols lost due to unrecoverable failures during the rebuild of the  $C_u$  codewords is equal to  $C_u E(L_u)$ , as determined by (85). Consequently, upon entering exposure level u, the expected number  $E(S_U)$  of symbols lost due to unrecoverable failures during the rebuild of the most-exposed codewords is  $E(C_u) E(L_u)$ . Therefore, the expected number  $E(S_{UF_u})$  of symbols lost due to unrecoverable failures at exposure level u is obtained by also considering the probability  $P_u$  of entering exposure level u, as determined by (91).

Note that when entering exposure level  $\tilde{r}$ , for each of the  $C_{\tilde{r}}$  most-exposed codewords there are  $\tilde{r}$  symbols permanently

lost. Therefore, the number of data symbols permanently lost is  $C_{\tilde{r}} \tilde{r}$ . Consequently,

$$E(S_{\rm DF}) \approx P_{\rm DF} E(C_{\tilde{r}}) \tilde{r}$$
 (92)

Substituting (92) into (36), and using (44), yields (39).

#### REFERENCES

- I. Iliadis, "Reliability assessment of erasure-coded storage systems with latent errors," in Proceedings of the 14th International Conference on Communication Theory, Reliability, and Quality of Service (CTRQ), Apr. 2021, pp. 15–24.
- [2] D. A. Patterson, G. Gibson, and R. H. Katz, "A case for redundant arrays of inexpensive disks (RAID)," in Proceedings of the ACM International Conference on Management of Data (SIGMOD), Jun. 1988, pp. 109– 116.
- [3] P. M. Chen, E. A. Lee, G. A. Gibson, R. H. Katz, and D. A. Patterson, "RAID: High-performance, reliable secondary storage," ACM Comput. Surv., vol. 26, no. 2, Jun. 1994, pp. 145–185.
- [4] V. Venkatesan, I. Iliadis, C. Fragouli, and R. Urbanke, "Reliability of clustered vs. declustered replica placement in data storage systems," in Proceedings of the 19th Annual IEEE/ACM International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS), Jul. 2011, pp. 307–317.
- [5] I. Iliadis, D. Sotnikov, P. Ta-Shma, and V. Venkatesan, "Reliability of geo-replicated cloud storage systems," in Proceedings of the 2014 IEEE 20th Pacific Rim International Symposium on Dependable Computing (PRDC), Nov. 2014, pp. 169–179.
- [6] M. Malhotra and K. S. Trivedi, "Reliability analysis of redundant arrays of inexpensive disks," J. Parallel Distrib. Comput., vol. 17, no. 1, Jan. 1993, pp. 146–151.
- [7] A. Thomasian and M. Blaum, "Higher reliability redundant disk arrays: Organization, operation, and coding," ACM Trans. Storage, vol. 5, no. 3, Nov. 2009, pp. 1–59.
- [8] I. Iliadis, R. Haas, X.-Y. Hu, and E. Eleftheriou, "Disk scrubbing versus intradisk redundancy for RAID storage systems," ACM Trans. Storage, vol. 7, no. 2, Jul. 2011, pp. 1–42.
- [9] V. Venkatesan, I. Iliadis, and R. Haas, "Reliability of data storage systems under network rebuild bandwidth constraints," in Proceedings of the 20th Annual IEEE International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS), Aug. 2012, pp. 189–197.
- [10] J.-F. Pâris, T. J. E. Schwarz, A. Amer, and D. D. E. Long, "Highly reliable two-dimensional RAID arrays for archival storage," in Proceedings of the 31st IEEE International Performance Computing and Communications Conference (IPCCC), Dec. 2012, pp. 324–331.
- [11] I. Iliadis and V. Venkatesan, "Most probable paths to data loss: An efficient method for reliability evaluation of data storage systems," Int'l J. Adv. Syst. Measur., vol. 8, no. 3&4, Dec. 2015, pp. 178–200.
- [12] —, "Expected annual fraction of data loss as a metric for data storage reliability," in Proceedings of the 22nd Annual IEEE International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS), Sep. 2014, pp. 375–384.
- [13] —, "Reliability evaluation of erasure coded systems," Int'l J. Adv. Telecommun., vol. 10, no. 3&4, Dec. 2017, pp. 118–144.
- [14] I. Iliadis, "Reliability evaluation of erasure coded systems under rebuild bandwidth constraints," Int'l J. Adv. Networks and Services, vol. 11, no. 3&4, Dec. 2018, pp. 113–142.
- [15] —, "Data loss in RAID-5 storage systems with latent errors," in Proceedings of the 12th International Conference on Communication Theory, Reliability, and Quality of Service (CTRQ), Mar. 2019, pp. 1–9.
- [16] —, "Data loss in RAID-5 and RAID-6 storage systems with latent errors," Int'l J. Adv. Software, vol. 12, no. 3&4, Dec. 2019, pp. 259– 287.
- [17] Amazon Web Services, "Amazon Simple Storage Service (Amazon S3)," 2022. [Online]. Available: http://aws.amazon.com/s3/ [retrieved: December 7, 2022]

- [18] D. Borthakur et al., "Apache Hadoop goes realtime at Facebook," in Proceedings of the ACM International Conference on Management of Data (SIGMOD), Jun. 2011, pp. 1071–1080.
- [19] R. J. Chansler, "Data availability and durability with the Hadoop Distributed File System," ;login: The USENIX Association Newsletter, vol. 37, no. 1, Feb. 2012, pp. 16–22.
- [20] K. Shvachko, H. Kuang, S. Radia, and R. Chansler, "The Hadoop Distributed File System," in Proceedings of the 26th IEEE Symposium on Mass Storage Systems and Technologies (MSST), May 2010, pp. 1–10.
- [21] S. Ghemawat, H. Gobioff, and S.-T. Leung, "The Google file system," in Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP), Oct. 2003, pp. 29–43.
- [22] D. Borthakur. HDFS and Erasure Codes (HDFS-RAID), Aug. 2009. [Online]. Available: https://hadoopblog.blogspot.com/2009/08 [retrieved: December 7, 2022]
- [23] B. Calder et al., "Windows Azure Storage: a highly available cloud storage service with strong consistency," in Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP), Oct. 2011, pp. 143–157.
- [24] D. Ford, F. Labelle, F. I. Popovici, M. Stokely, V.-A. Truong, L. Barroso, C. Grimes, and S. Quinlan, "Availability in globally distributed storage systems," in Proceedings of the 9th USENIX Symposium on Operating Systems Design and Implementation (OSDI), Oct. 2010, pp. 61–74.
- [25] S. Muralidhar et al., "f4: Facebook's Warm BLOB Storage System," in Proceedings of the 11th USENIX Symposium on Operating Systems Design and Implementation (OSDI), Oct. 2014, pp. 383–397.
- [26] C. Huang, H. Simitci, Y. Xu, A. Ogus, B. Calder, P. Gopalan, J. Li, and S. Yekhanin, "Erasure coding in Windows Azure Storage," in Proceedings of the USENIX Annual Technical Conference (ATC), Jun. 2012, pp. 15–26.
- [27] E. Pinheiro, W.-D. Weber, and L. A. Barroso, "Failure trends in a large disk drive population," in Proceedings of the 5th USENIX Conference on File and Storage Technologies (FAST), Feb. 2007, pp. 17–28.
- [28] A. Dholakia, E. Eleftheriou, X.-Y. Hu, I. Iliadis, J. Menon, and K. Rao, "A new intra-disk redundancy scheme for high-reliability RAID storage systems in the presence of unrecoverable errors," ACM Trans. Storage, vol. 4, no. 1, May 2008, pp. 1–42.
- [29] I. Iliadis, "Reliability modeling of RAID storage systems with latent errors," in Proceedings of the 17th Annual IEEE/ACM International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS), Sep. 2009, pp. 111–122.
- [30] V. Venkatesan and I. Iliadis, "Effect of latent errors on the reliability of data storage systems," in Proceedings of the 21st Annual IEEE International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS), Aug. 2013, pp. 293–297.
- [31] —, "A general reliability model for data storage systems," in Proceedings of the 9th International Conference on Quantitative Evaluation of Systems (QEST), Sep. 2012, pp. 209–219.
- [32] M. Silberstein, L. Ganesh, Y. Wang, L. Alvisi, and M. Dahlin, "Lazy means smart: Reducing repair bandwidth costs in erasure-coded distributed storage," in Proceedings of the 7th ACM International Systems and Storage Conference (SYSTOR), Jun. 2014, pp. 15:1–15:7.
- [33] K. V. Rashmi, N. B. Shah, D. Gu, H. Kuang, D. Borthakur, and K. Ramchandran, "A "Hitchhiker's" guide to fast and efficient data reconstruction in erasure-coded data centers," in Proceedings of the 2014 ACM conference on SIGCOMM, Aug. 2014, pp. 331–342.
- [34] DELL/EMC Whitepaper, "PowerVault ME4 Series ADAPT Software," Feb. 2019. [Online]. Available: https://www.dellemc.com/ [retrieved: December 7, 2022]
- [35] I. Iliadis and V. Venkatesan, "Rebuttal to 'Beyond MTTDL: A closedform RAID-6 reliability equation'," ACM Trans. Storage, vol. 11, no. 2, Mar. 2015, pp. 1–10.
- [36] T. J. E. Schwarz, Q. Xin, E. L. Miller, D. D. E. Long, A. Hospodor, and S. Ng, "Disk scrubbing in large archival storage systems," in Proceedings of the 12th Annual IEEE/ACM International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS), Oct. 2004, pp. 409–418.

- [37] A. Oprea and A. Juels, "A clean-slate look at disk scrubbing," in Proceedings of the 8th USENIX Conference on File and Storage Technologies (FAST), Feb. 2010, pp. 57–70.
- [38] B. Schroeder, S. Damouras, and P. Gill, "Understanding latent sector errors and how to protect against them," ACM Trans. Storage, vol. 6, no. 3, Sep. 2010, pp. 1–23.
- [39] M. Zhang, S. Han, and P. P. C. Lee, "SimEDC: A simulator for the reliability analysis of erasure-coded data centers," IEEE Trans. Parallel Distrib. Syst., vol. 30, no. 12, 2019, pp. 2836–2848.
- [40] V. Venkatesan and I. Iliadis, "Effect of codeword placement on the reliability of erasure coded data storage systems," in Proceedings of the 10th International Conference on Quantitative Evaluation of Systems (QEST), Sep. 2013, pp. 241–257.
- [41] —, "Effect of codeword placement on the reliability of erasure coded data storage systems," IBM Research Report, RZ 3827, Aug. 2012.
- [42] I. Iliadis and X.-Y. Hu, "Reliability assurance of RAID storage systems for a wide range of latent sector errors," in Proceedings of the 2008 IEEE International Conference on Networking, Architecture, and Storage (NAS), Jun. 2008, pp. 10–19.
- [43] Seagate, exos x20, data sheet. [Online]. Available: https://www.seagate.com/products/enterprise-drives/exos-x/x20/ [retrieved: December 7, 2022]
- [44] Backblaze drive stats for 2021. [Online]. Available: https://www.backblaze.com/blog/backblaze-drive-stats-for-2021/ [retrieved: December 7, 2022]
- [45] M. Ovsiannikov, S. Rus, D. Reeves, P. Sutter, S. Rao, and J. Kelly, "The quantcast file system," in Proceedings of the 39th International Conference on Very Large Data Bases (VLDB), vol. 6, no. 11. VLDB Endowment, Aug. 2013, pp. 1092–1101.

# Digital Sensing Platform with High Accuracy Time Synchronization Function of Vibration and Camera Sensors

Narito Kurata Faculty of Industrial Technology Tsukuba University of Technology Tsukuba City, Ibaraki, Japan e-mail: kurata@home.email.ne.jp

Abstract - The author is conducting research and development of different types of sensor systems for the maintenance of civil infrastructures, such as aging bridges and highways, and buildings. Highly accurate time information is added to measurement data, and a set of sensing data that ensures time synchronization is acquired and used for multimodal analysis of risk. In research so far, as a first step, a sensor device was developed that uses a digital high-precision accelerometer to perform highly-accurate time-synchronized sensing of civil infrastructures and buildings. A vibration table test was performed on the sensor device, and its time synchronization performance was verified. In this paper, a different type of digital sensing platform has been developed that allows a camera sensor to be connected in addition to a digital accelerometer. By adding a unified time stamp synchronized with the absolute time to data from the digital accelerometer and the image from the camera sensor when data is acquired. the vibration and the image can be measured synchronously. First, a Chip-Scale Atomic Clock (CSAC), which is an ultrahigh-precision clock, is mounted on the sensor device, and a mechanism is implemented whereby a time stamp is added to the outputs of the digital accelerometer and the camera sensor with timekeeping precision. A Field-Programmable Gate Array (FPGA) dedicated to adding time stamps is prepared. Tests were performed on the developed sensor device using a shaking table, and the time synchronization performance was checked by comparing the measurement results of multiple sensor devices and a servo type acceleration sensor. Next, the results of a performance verification experiment on a camera sensor mounted on the platform, are described.

Keywords-Time Synchronization; Chip Scale Atomic Clock; Earthquake Observation; Structural Health Monitoring; Micro Electro Mechanical Systems; Camera Sensor

# I. INTRODUCTION

As civil infrastructure, such as bridges and highways, and buildings deteriorate over time, it has become important to automate inspections for maintenance of these structures. In addition, since there are many disasters, such as earthquakes and typhoons in Japan, it is necessary to detect damage to structures immediately after a disaster, and to estimate the damage situation. Data collection and analysis by sensor groups is effective for automating the detection of such abnormalities, but to analyze data sets measured by multiple sensors and evaluate structural safety, time synchronization between sensors is required [1]. The authors applied wireless sensor network technology to develop sensors for seismic observation and structural health monitoring, and demonstrated their performance in skyscrapers [2][3]. In this system, time synchronization was achieved by sending and receiving wireless packets between sensors [2]. However, it is impossible to target multiple buildings, long structures, such as bridges, and wide-area urban spaces with wireless sensor network technology. On the other hand, if sensors installed in various places can autonomously retain accurate time information, this problem can be resolved. The method of using GPS signals is effective outdoors, but it cannot be used inside buildings, underground, under bridges, or in tunnels, etc.

Therefore, the author developed a sensor device that autonomously retains accurate time information using a Chip-Scale Atomic Clock (CSAC) [4]-[6], which is an ultra-highprecision clock [7]-[9]. In order to apply the developed sensor device to earthquake observation, a logic was implemented that detects the occurrence of an earthquake and saves data on earthquake events, and its function was verified in a vibration table experiment [10]. The developed sensor device was also installed in an actual building and on an actual bridge, and used for seismic observation and evaluation of structural health [11]. However, as the developed sensor device had an analog MEMS accelerometer, it was difficult to measure minute vibrations accurately, and there was a risk of noise being mixed with the analog signal. The risk of noise was therefore eliminated by making the accelerometer mounted on the sensor device a digital type [1]. In this paper, a different type of digital sensor platform is further developed wherein a camera sensor can be connected to the sensor device. The details of the digital sensing platform and a mechanism whereby ultra-high-accuracy time information is added to sensor data by the CSAC, are described. The results of experiments performed to verify the time synchronization performance of the camera sensor are also reported.

In this paper, Section II shows the existing time synchronization methods and describes their problems and achievement of the development of digital sensing platform proposed in this research. Section III describes the mechanism for providing ultra-high accurate time information to digital sensor data by the CSAC, and explains the configuration of digital sensing platform and the development of the actual sensor device. Section IV describes the configuration of digital sensing platform and the development of the actual sensor device. Further, Sections V and VI describe the performance verification tests on the time synchronization of developed sensor device and camera sensor device, and it is confirmed that time synchronization among the developed digital sensor devices with camera is achieved.

#### II. STATE OF THE ART

A time synchronization function is indispensable for sensor devices used to monitor structural health and make seismic observations of civil infrastructure, such as bridges and highways, and buildings. This is because time series analysis using phase information cannot be performed unless a data set is obtained in which time synchronization is achieved. Regarding time synchronization of sensing, many studies have already been performed, such as the use of GNSS signals from artificial satellites and the Network Time Protocol (NTP) for time synchronization on the Internet [12]. There are also studies where time synchronization is achieved by utilizing the characteristic of wireless sensor networks that propagation delay is small. For example, time synchronization protocols, such as Reference Broadcast Synchronization (RBS), Timing-sync Protocol for Sensor Networks (TPSN), and Flooding Time Synchronization Protocol (FTSP) are being studied [13]-[17].

However, although time synchronization using wireless technology is convenient, wireless communication may not always be possible. Particularly, if wireless communication is interrupted during an earthquake, it will not be possible to perform sensing where time synchronization is guaranteed. A technology that realizes highly accurate time synchronization in a room includes IEEE1588 Precision Time Protocol (PTP). PTP uses an Ethernet cable in a general Local Area Network (LAN) as a transmission line, and achieves accurate synchronization within one microsecond using time packets. However, it has many problems. For example, it is difficult to achieve stable synchronization accuracy due to packet delay fluctuation and packet loss due to congestion in the LAN. Moreover, since the delay is corrected by packet switching, the PTP devices that can be connected to the master device are limited, and it cannot be deployed in a Wide Area Network (WAN) environment where the delay amount varies drastically.

To obtain a set of sensor data that guarantees long-term, stable time synchronization even when GPS signals are not available, wireless transmission/reception is unstable, and wired network connection is not possible, it would be ideal if different sensors autonomously retain accurate time information. If accurate time information could be added to the data measured by each sensor, a sensor data set that guarantees time synchronization would be obtained. It was therefore decided to develop a sensing system that autonomously retains accurate time information by employing a CSAC [5]-[7], which is a clock with high timekeeping accuracy. The CSAC is a clock that achieves ultra-highaccuracy time measurement to several tens of picoseconds (5 x 10<sup>-11</sup> seconds), while having an ultra-small external shape that can be mounted on a board. Development began in 2001 with the support of the US Defense Advanced Research Projects Agency (DARPA), and consumer products were launched in 2011.

Applications include measures against GPS positioning interference by jamming signals, high-precision positioning by installing on smartphones, etc., and high-level assessment of disaster situations, and further price reduction is expected as it becomes more widespread. CSAC has a small error of about 4 to 8 orders of magnitude less than that of timekeeping by a crystal oscillator, and time synchronization by NTP or GPS signals. If this CSAC is installed in each sensor device and a mechanism is implemented that gives a highly accurate time stamp to sampling of the data to be measured, sensor data sets that guarantee time synchronization can be collected even if GPS signals cannot be used, wireless transmission/reception is unstable, and wired network connection is unavailable. In development so far, the sensor device had a MEMS accelerometer, and the system configuration allowed for any analog sensor to be connected via an external input interface.

However, as the accuracy of the analog MEMS accelerometer is not high, noise might still be mixed with the analog signal, and it was desired to be able to connect a camera sensor, which is a type of digital sensor, it was decided to develop a new platform with full digital sensing. Specifically, a technology was developed to mount a digital accelerometer on the sensor device to enable high-precision acceleration measurement without the risk of noise contamination, connect a camera sensor, and assign accurate time stamps by CSAC to both digital outputs. Data sets obtained by the sensor device described in this paper, where time synchronization is guaranteed, can be used to analyze the structural health of civil infrastructure and buildings, and understand seismic phenomena. In addition, although research has been conducted on the application of camera sensors to structural health monitoring [18][19], time synchronization between sensors has not been rigorously studied.

# III. TIME STAMPING MECHANISM USING CHIP SCALE Atomic Clock

A CSAC, shown in Figure 1, has time accuracy equivalent to that of a rubidium atomic clock and is very accurate compared with crystal resonators [6][8]. The CSAC can achieve ultra-precision time measurement at a level of some ten picoseconds, consumes low power and is small enough to be mounted on a circuit board (Table I). The development of the CSAC started with the support of Defense Advanced Research Projects Agency, and the commercial product was released by an American company in 2011 and is still available for purchase. Recently, ultra-small atomic clock systems, which can be mounted on general communication terminals, such as smart phones, have been proposed, and further downsizing and price reduction are expected. If the sensor device is equipped with a CSAC and a mechanism that adds time stamping for every sample of measured data, the sensor device can create data having high-accuracy time information. Each sensor device autonomously keeps highly accurate time information even if the GPS signals and network communication are unavailable. Therefore, by collecting the measured data by means, such as 3G, Wi-Fi, Ethernet, etc., a data group ensuring time synchronization can be obtained.



Figure 1. Chip Scale Atomic Clock (CSAC).

Rise/fall time: < 10 ns Pulse width: 100 μs	

TABLE I. SPECIFICATIONS OF CSAC

To configure a sensing system composed of multiple sensor devices equipped with a CSAC, one device is set as a master device and other devices as slave devices must be synchronized by defining absolute time information. The main controller of each sensor device is equipped with an input/output connector for 1 Pulse Per Second (PPS) of the CSAC. Using this connector, the master device outputs 1 PPS signal, and each slave device inputs it to synchronize and match the phase of the CSAC in each slave device. The CSAC keeps accurate time, but it does not have absolute time information. Therefore, it must be defined separately. At initial settings, the GPS module installed in the main controller is used. Absolute time information is transmitted from the master device to the slave device by the IEEE 1588 standard. Once all the sensor devices are synchronized at the beginning, they continue keeping highly accurate time information autonomously. It is only necessary to install the sensor device in an arbitrary place and collect data. As mentioned above, any means of data collection, such as Ethernet, Wi-Fi, or 3G, are available as the measured data records accurate time stamping.

# IV. DIGITAL SENSING PLATFORM AND CIRCUIT CONFIGURATION

An ordinary sensor device consists of a CPU, a sensor, a memory and a network interface, and a crystal oscillator is used for the CPU. If a CSAC is mounted on such a sensor device to correct the time information of the CPU and perform measurement, a delay will occur because the timing accuracy of the CSAC is too high. Therefore, in order to directly add time information by the CSAC to the measurement data of the sensor in terms of hardware, a mechanism utilizing a Field-Programmable Gate Array (FPGA), which is a dedicated integrated circuit, was contrived. As a result, the CPU of the sensor device is not overloaded, and it became possible to save measurement data to which time information is added by the FPGA in a memory, and to collect the data via a network. Moreover, since the FPGA is programmable, it can not only handle CSAC time information, but also incorporate logic, for example to detect abnormalities, etc., using the measurement data. In this paper, a mechanism is developed whereby accurate time stamps by the CSAC are added to the outputs of the digital accelerometer and camera sensor.

#### 4.1 System Configuration

As shown in Figures 2 and 3, the sensor device in this research consists of an oscillator board that synchronizes GPS Time (GPST) with the CSAC and supplies a stable reference signal, a sensor board on which a digital accelerometer and external analog sensor input interface are mounted, a signal processing board on which a CPU and FPGA are mounted, and a camera that captures images. A high-precision 10 MHz reference clock and one PPS signal are supplied by the oscillator board, and a time stamp and trigger signal for acquiring data are generated by the FPGA. The sensor board is provided with the digital accelerometer and external analog sensor input interface. Any analog sensor can be connected to the external analog sensor input interface. The sensor devices, which are individually equipped with the sensor board containing the digital accelerometer, can be freely combined with the other sensor devices having displacement sensors or strain sensors via external analog sensor input interface.



Figure 2. System configuration.



Figure 3. Oscillator board configuration.

The digital accelerometer outputs data according to the trigger signal via a Universal Asynchronous Receiver /Transmitter (UART). The data of the sensor connected to the external analog sensor input interface is converted by an A/D converter according to the trigger signal, and output as a 16-bit serial value. The camera sensor can release the shutter according to the trigger signal, and outputs it as an RGB value. The data thus acquired is saved in a connected storage (SSD). The sensor device is operated via a wired LAN, Wi-Fi, or USB.

#### 4.2 Oscillator board

The configuration and external appearance of the oscillator board are shown in Figures 4 and 5. The oscillator board has a function to synchronize the CSAC with 1 PPS output by the GPS module or 1 PPS input from outside. When time synchronization is required between multiple sensor devices, all the sensor devices are designated as "master", or one sensor device is designated as a "master" and the other sensor devices are designated as "slave".



Figure 4. Oscillator board configuration.



Figure 5. External appearance of the oscillator board.

When the sensor device is a master, GPS and the CSAC are synchronized by receiving GPS signals, and inputting 1 PPS obtained from the GPS module to the CSAC. When the sensor device is a slave, the master and slave are synchronized by inputting 1 PPS output by the master. In addition, commands for setting the CSAC synchronization cycle or resetting the phase value, as well as signal selection commands, are executed by the signal processing board through the connector. The 10MHz and 1 PPS output by this board are clock sources for this system.

#### 4.3 Sensor board

The configuration and external appearance of the sensor board are shown in Figures 6 and 7. The sensor board is provided with a digital accelerometer and an external analog sensor input interface. The data obtained by the digital accelerometer can be sampled at the timing of the trigger. For the external analog sensor input interface, three channels are provided assuming that a servo-type analog accelerometer is connected for comparison with the digital accelerometer. Depending on the measurement purpose, any analog sensor, such as a displacement sensor, strain sensor or crack detection sensor, can be connected in addition to a servo-type analog accelerometer.



Figure 6. Sensor board configuration.



Figure 7. Sensor board.

The signal input from the external analog sensor interface is converted by the A/D converter to output as a 16-bit serial value. Note that the signal is split into two paths, and one of them is amplified 64 times. Therefore, even with an A/D converter having a resolution of 16 bits, a resolution equivalent to 22 bits can be obtained. The data obtained by the digital accelerometer is output by the UART, and the data obtained by the sensor connected to the external analog sensor input interface is output by a Serial Peripheral Interface (SPI), both to the signal processing board.

#### 4.4 Signal processing board

The configuration and external appearance of the signal processing board are shown in Figures 8 and 9. The FPGA shown in Table 2 is mounted on the signal processing board. As shown in Table 1, Ubuntu is installed as the OS of the CPU and Cyclone V is installed in the FPGA, enabling the use of 1 GB of SDRAM. Also installed is a USB On-The-Go (USB OTG), and it is possible to be connected to an SSD and a Wi-Fi antenna, which are extension devices. The Samba function can be used to acquire and browse built-in data via a LAN, and Secure Shell (SSH) allows to perform operations, such as settings and starting measurement. In addition, since it has a USB slave function, it can be operated via USB even when a LAN cannot be used. In addition to the above functions, the CPU mainly performs time setting, sorting of acquired data, format conversion, and filing. The FPGA adjusts the internal RTC (real-time clock), and generates the trigger signal based on the clock obtained from the oscillator. It also constitutes a data acquisition block for the various sensors and camera.

# 4.5 Digital accelerometer

The external appearance and specifications of the digital accelerometer mounted on the sensor board are shown in Figure 10 and Table 3, respectively. The on-board digital MEMS has a built-in 3-axis crystal accelerometer with high precision and stable performance manufactured by finely processing a crystalline material with superior precision and stability. As shown in Table 3, high-resolution vibration measurement with low noise (0.5  $\mu$  G/ $\sqrt{}$  Hz) and low power (66 mW), is possible. It enables measurement bandwidths of up to 100 Hz, and an increased data output rate of up to 1000 Sps.



Figure 8. Signal processing board configuration.



Figure 9. Signal processing board.

TABLE II. SPECIFICATIONS OF MAIN FPGA AND DE10-NANO(CPU)

Model	DE10-NANO	
OS	Ubuntu 16.04.6 LTS (GNU/Linux 4.5.0-00185-g3bb556b armv7l)	
CPU	800MHz Dual-core ARM Cortex-A9	
Memory	1GB DDR3 SDRAM	
LAN	1 Gigabit Ethernet PHY with RJ45 connector	
USB	USBOTG×1, USBUART×1	
UART	UART×2	
FPGA	CYCLONE V	
FPGAROM	EPCS64	
Storage	16GB (MicroSD)	



Figure 10. Digital Accelerometer.

Model	EPSON M-A351AS
Range	±5 G
Noise Density	0.5 µG/√Hz (Average)
Resolution	0.06 µG/LSB
Bandwidth	100 Hz (selectable)
Output Range	1000 sps (selectable)
Digital Serial Interface	SPI
Outside Dimensions (mm)	$24 \times 24 \times 18$
Weight	12 grams
Operating Temperature	-20 °C to +85 °C
Power Consumption	3.3 V, 66 mW
Output Mode Selection	Acceleration, Tilt Angle, or Tilt Angle Speed

#### 4.6 Camera sensor

The external appearance and specifications of the camera sensor are shown in Figure 11 and Table 4, respectively.



Figure 11. Camera sensor.

Model	OmniVision OV5642
Active Array Size	2592 x 1944
Power Supply	core: 1.5VDC±5%, analog: 2.6- 3.0 V, I/O: 1.7-3.0 V
Temperature Range	operating: -30 °C to +70 °C stable image: 0 °C to +50 °C
Output Formats (8-bit)	YUV(422/420)/YCbCr422, RGB565/555/444, CCIR656, 8-bit compression data, 8/10-bit raw RGB data
Lenz Size	1/4"
Input Clock Frequency	6-27 MHz
Shutter	rolling shutter
Maximum Image Transfer Rate	5 megapixel (2592x1944): 15 fps 1080p (1920x1080): 30 fps 720p (1280x720):60fps VGA (640x480): 60 fps QVGA (320x240): 120 fps
Scan mode	progressive
Maximum Exposure Interval	1968 x t <sub>row</sub>
Pixel Size	1.4 $\mu$ m x 1.4 $\mu$ m
Image Area	3673.6 μ m x 2738.4 μ m

OmniVision OV5642, which is a CMOS camera module, is used as the camera sensor. It is compact, has low power consumption, supports digital data (YUV422) output, performs very well in poorly lite environments, and can acquire images at the timing of a trigger by inputting an external trigger.

# V. PERFORMANCE CONFIRMATION TESTS ON THE TIME SYNCHRONIZATION FUNCTION OF SENSOR DEVICES WITH DIGITAL ACCELEROMETER

Tests were performed using a shaking table to confirm the performance of the developed digital sensor device. The objective was to confirm the measurement performance and time synchronization performance of the digital sensor device. Three sensor devices and a servo acceleration sensor for comparison were fixed on a shaking table as shown in Figures 12 and 13, the same vibrations were applied in one horizontal direction, and the results were compared. The analog output of the comparative servo acceleration sensor was input to the sensor devices via the external input interface. In the test, a sweep wave of 2 to 20 Hz as shown in Figures 14 and 15 was applied to excite the shaking table as an input wave, and the measurement sampling frequency of the sensor devices was set to 1,000 Hz.



Figure 12. Experimental setup.



Figure 13. Sensor boards on shaking table.



Figure 15. Power spectrum of input swept sine wave



Figure 16. Spectrum ratios of Fourier amplitudes of three sensor modules to servo-type acceleration sensor (X direction).



Figure 17. Spectrum ratios of Fourier amplitude of two slave modules to master module (X direction).



Figure 18. Spectrum ratios of Fourier phase of two slave modules to master module (X direction).

Excitation was applied in the X-direction of the sensor devices, and measurement was performed by the sensor devices and the comparative servo acceleration sensor. Figure 16 shows the results of calculation of the Fourier amplitude spectrum ratios of the acceleration time history measured by the three sensor devices and the comparative servo acceleration sensor. The amplitude of the former three devices relative to the latter reflected the low pass filter characteristics of the digital sensors, so it can be seen that the digital acceleration sensor mounted on the sensor boards have good performance.

Next, Figure 17 shows the results of obtaining the Fourier amplitude spectrum ratios for the measured acceleration data from two sensor devices (slaves) on the shaking table when the other sensor device was used as the master. The amplitude of the former two devices relative to the latter was perfectly flat over the frequency band 2 to 20 Hz, despite the measured results. In addition, Figure 18 shows the results of obtaining the Fourier phase spectrum ratios for the measured acceleration data from two sensor devices (slaves) on the shaking table when the other sensor device was used as the master. There is no phase delay between the sensor devices, and if the time synchronization was maintained it would be expected that it should be about zero over the frequency band 2 to 20 Hz. From the figure, it can be seen that time synchronization has been achieved between the sensor devices.

#### VI. PERFORMANCE VERIFICATION EXPERIMENT ON THE TIME SYNCHRONIZATION FUNCTION OF CAMERA SENSOR

A performance verification experiment was carried out on the time synchronization function of the camera sensor. The experimental system configuration is shown in Figure 19. The trigger signal generated by the signal processing board is transmitted to the camera sensor and the FPGA of the LED control simultaneously. Since the shutter of the camera sensor and the lighting of the LEDs are synchronized, if the image can be acquired when the LEDs light up, it is considered that the image acquired is synchronized with the trigger. The FPGA for LED control shown in Figure 20 is configured to control the lighting of the LEDs at the clock timing of the trigger signal. As shown in Figure 12, 5×5 matrix LEDs light up one by one from upper left to lower right according to the rise of the trigger signal. As shown in Figure 21, the matrix LEDs are photographed with the camera sensor fixed.

Figure 22 shows the result of imaging by the camera sensor. Time elapses from the upper left to the lower right. Since the LEDs in the image light up one by one, images can be acquired according to the signal input to the FPGA for LED control. In the images, the LEDs light up two at a time, and it can be seen that the camera sensor has a certain delay with respect to the trigger. As shown in Figure 23, it was verified that images could be continuously acquired in synchronization with the trigger.



Figure 19. Camera data measurement system.



Figure 20. FPGA for LED control.



Figure 21. Matrix LEDs.



Figure 22. Imaging arrangement.

The performance verification experiment on the time synchronization of developed camera sensor device was carried out as shown above. It is confirmed that time synchronization among the developed digital sensor devices with camera is achieved. The development of this digital sensing platform has enabled time-synchronized measurements between digital accelerometers, camera sensors, and many types of external input analog sensors, as well as multimodal analysis between measurement data.

#### VII. CONCLUSION

In this paper, research and development relating to a digital sensing platform that autonomously retains highly accurate time information by applying a CSAC, was reported. First, a system based on a digital sensor and autonomous time synchronization by a CSAC was described, in which the development of a mechanism and a sensor device that add ultra-high-accuracy time information to sensor data using the CSAC were explained in detail. A function was added to assign the same time stamp as that of the output of the built-in digital accelerometer, to the output of the camera sensor.

Next, the results of tests performed to confirm the time synchronization performance of the sensor device were reported. Three sensor devices were mounted on a shaking table and tests were performed by applying vibrations simultaneously, and by checking the phase properties of the measurement results it was confirmed that time synchronization was achieved for sampling at 1,000 Hz.

The results of an experiment carried out to verify the time synchronization performance of the camera sensor were also reported. In the future, the author plan to apply this new, different type of digital sensing platform to actual structures to acquire acceleration and video data that retain accurate time information. One possible problem is that although the timing accuracy of the CSAC is high, aging will occur in the long term, so it may be necessary to consider how to operate the sensing system according to the purpose and object of measurement. Further, as CSAC are currently expensive, it is hoped that they will be used more extensively in many fields.



Figure 23. Continuous images from the experimental result.

#### ACKNOWLEDGMENT

This research was partially supported by the New Energy and Industrial Technology Development Organization (NEDO) through the Project of Technology for Maintenance, Replacement and Management of Civil Infrastructure, Cross-ministerial Strategic Innovation Promotion Program (SIP). This research was also partially supported by JSPS KAKENHI Grant Number JP19K04963.

#### REFERENCES

- N. Kurata, "High-precision Time Synchronization Digital Sensing Platform Enabling Connection of a Camera Sensor," The Twelfth International Conference on Sensor Device Technologies and Applications (SENSORDEVICES 2021) IARIA, Nov. 2021, pp. 98-104, ISSN: 2308-3514, ISBN: 978-1-61208-918-8
- [2] N. Kurata, M. Suzuki, S. Saruwatari, H. Morikawa, "Actual Application of Ubiquitous Structural Monitoring System using Wireless Sensor Networks," 14th World Conference on

Earthquake Engineering (14WCEE) IAEE, Oct. 2008, Paper ID:11-0037, pp. 1-8.

- [3] N. Kurata, M. Suzuki, S. Saruwatari, H. Morikawa, "Application of Ubiquitous Structural Monitoring System by Wireless Sensor Networks to Actual High-rise Building," Fifth World Conference on Structural Control and Monitoring (5WCSCM) IASCM, July 2010, Paper No. 013, pp. 1-9.
- [4] N. Kurata, "Disaster Big Data Infrastructure using Sensing Technology with a Chip Scale Atomic Clock," World Engineering Conference and Convention (WECC2015) WFEO, Dec. 2015, pp. 1-5.
- [5] N. Kurata, "Basic Study of Autonomous Time Synchronization Sensing Technology Using Chip Scale Atomic Clock," 16th International Conference on Computing in Civil and Building Engineering (ICCCBE2016) ISCCBE, July 2016, pp. 67-74.
- [6] N. Kurata, "An Autonomous Time Synchronization Sensor Device Using a Chip Scale Atomic Clock for Earthquake Observation and Structural Health Monitoring" The Eighth International Conference on Sensor Device Technologies and Applications (SENSORDEVICES 2017) IARIA, Sep. 2017, pp.31-36, ISSN: 2308-3514, ISBN: 978-1-61208-581-4

- [7] S. Knappe, et al., "A Microfabricated Atomic Clock," Applied Physics Letters, vol. 85, Issue 9, pp. 1460-1462, Aug. 2004, doi:10.1063/1.1787942.
- [8] Q. Li and D. Rus, "Global Clock Synchronization in Sensor Networks," IEEE Transactions on Computers, vol. 55, Issue 2, pp. 214-226, Jan. 2006, ISSN: 0018-9340.
- [9] R. Lutwak, et al., "The Chip-Scale Atomic Clock Prototype Evaluation," 39th Annual Precise Time and Time Interval (PTTI) Meeting, Nov. 2007, pp. 269-290.
- [10] N. Kurata, "Improvement and Application of Sensor Device Capable of Autonomously Keeping Accurate Time Information for Buildings and Civil Infrastructures," The Ninth International Conference on Sensor Device Technologies and Applications (SENSORDEVICES 2018) IARIA, Sep. 2018, pp. 114-120, ISSN: 2308-3514, ISBN: 978-1-61208-660-6
- [11] N. Kurata, "A Sensing System with High Accurate Time Synchronization for Earthquake Observation and Structural Health Monitoring of Structures," 17th World Conference on Earthquake Engineering (17WCEE) IAEE, Oct. 2021, Paper No. 9a-0008, pp. 1-9.
- [12] D. Mills, "Internet Time Synchronization: the Network Time Protocol," IEEE Transactions on Communications, vol. 39, Issue 10, Oct. 1991, pp. 1482-1493, doi:10.1109/26.103043.
- [13] M. Maroti, B. Kusy, G. Simon, and A. Ledeczi, "The Flooding Time Synchronization Protocol," Proc. the 2nd International Conference on Embedded Networked Sensor Systems (SenSys '04) ACM, Nov. 2004, pp. 39-49, doi:10.1145/1031495.1031501.
- [14] J. Spilker Jr., P. Axelrad, B. Parkinson and P. Enge, "Global Positioning System: Theory and Applications," Vol. I,

American Institute of Aeronautics and Astronautics (AIAA), 1996, ISBN: 978-1-56347-106-3.

- [15] J. Elson, L. Girod, and D. Estrin, "Fine-Grained Network Time Synchronization using Reference Broadcasts," Proc. 5th Symposium on Operating Systems Design and Implementation (OSDI'02), Dec. 2002, pp. 147-163, doi:10.1145/844128.844143.
- [16] S. Ganeriwal, R. Kumar, and M. B. Srivastava, "Timing-sync Protocol for Sensor Networks," Proc. the 1st International Conference on Embedded Networked Sensor Systems (SenSys '03) ACM, Nov. 2003, pp. 138-149, doi:10.1145/958491.958508.
- [17] K. Romer, "Time Synchronization in Ad Hoc Networks," Proc. the 2nd ACM International Symp. on Mobile Ad Hoc Networking & Computing (MobiHoc'01) ACM, Oct. 2001, pp. 173-182, doi:10.1145/501436.501440.
- [18] A. Alzughaibi, A. Ibrahim, Y. Na, S. El-Tawil, A. Eltawil, "Feasibility of Utilizing Smart-phone Cameras for Seismic Structural Damage Detection," 2020 IEEE International Instrumentation and Measurement Technology Conference (I2MTC) IEEE, May 2020, pp. 173-182, pp. 1-5, ISSN: 2642-2069, ISBN: 978-1-7281-4461-0
- [19] S. Yaryshev, L. Li, M. Marinov, T. Djamiykov, "Development of a Digital Camera-Based Method for Bridge Deformation Measurement," 2020 XXIX International Scientific Conference Electronics (ET) IEEE, Sep. 2020, pp. 1-4, ISBN: 978-1-7281-7427-3

# Combined Algorithm for Voronoi Diagram Construction as it Applies to Dynamic Ride Sharing

Anton Butenko and Jorge Marx Gómez

Department of Computing Science Carl von Ossietzky University of Oldenburg Oldenburg, Germany anton.butenko@uol.de, jorge.marx.gomez@uni-oldenburg.de

Abstract—Standard Voronoi diagram decomposes a plane into cells with a common closest site. This structure is widely used in computational geometry in application to the nearest neighbor problem. Using Euclidean metric is the most straightforward solution, however, in urban environment it may lead to insufficient accuracy that is crucial in such applications as dynamic ride sharing. Deviations in determining the nearest meeting point are especially significant under the presence of obstacles: water reservoirs, railway tracks, highways, industrial zones as well as hilly terrain. Here, we propose a combined approach for city Voronoi diagram construction in general metric space. A transportation network is modelled as weighted graph, so that the route consists of a foot-walking part and shortest path in graph. Presented algorithm constructs continuous Voronoi diagram for a plane using the individual graph Voronoi cells as generator objects. Evaluation for the specific city topography shows that the described algorithm provides more accurate results in comparasion with the standart Voronoi diagram.

Index Terms—Voronoi diagram; dynamic ride sharing.

# I. INTRODUCTION

Ride sharing applications are aimed at connecting drivers and passengers in an optimal way. What this optimal way means depends a lot on the specific mobility solution philosophy and its target audience. Nevertheless, most of them face such optimization problem as the nearest neighbor search: identifying the point from a set of points which is the closest to a given point according to some measure. The mobility application Instaride [3] developed for the spontaneous shared trips is driven by an instant matching algorithm. It connects drivers and passengers in real time based on the user's mobile device positioning (satellite navigation data, triangulation in mobile network) [2]. In order to minimize the driver's efforts and his route detour, the finite set of preselected fixed points is used for passengers' pick-up and drop-off (named meeting points, in general). Preselection of the meeting points is determined by the environmental conditions and is based on criteria such as parking opportunity, presence of pedestrian zones and easily recognizable landmarks. Such an approach leads to the problem of finding the nearest meeting point for users (both drivers and passengers) based on their realtime positions. The paper structure is the following. The Introduction explains the problem's origin. In Section II, we describe the concept of the presented approach and introduce the terms and notation. Sections III and IV describe two parts

of the algorithm: discrete and continuous. In Section V, the algorithm steps are given in detail. Section VI presents the algorithm efficiency evaluation for the specific city topography. Section VII concludes our work.

## II. VORONOI DIAGRAM IN A GENERALIZED METRIC SPACE

One of the most effective ways to solve problems related to the nearest neighbor search is to use the Voronoi diagram. We introduce the following notation here:  $L_{\rho}$  is a metric space with the corresponding function  $\rho: L \times L \to \mathbb{R}^+$  that satisfies metric axioms. Then,  $O_r(x) = \{z : \rho(x, z) < r\}$  is the open metric ball with radius  $r \in \mathbb{R}^+$ ,  $S_r(x) = \overline{O}_r(x) \setminus O_r(x)$ is the metric sphere and  $\Lambda(x, y) = \{z : \rho(x, z) = \rho(y, z)\}$ is the bisector of x and y. It splits  $L_{\rho}$  into the half-spaces  $D(x, y) = \{z : \rho(x, z) < \rho(y, z)\}$  and, lying on the other bisector side  $D(y, x) = \{z : \rho(y, z) = \rho(x, z)\}$ . For a given finite set of seeds  $S = \{s_1, ..., s_k\} \in L_{\rho}$ , the Voronoi cell related to  $s_i$  is expressed as

$$VR(s_i, S) = \bigcap_{i \neq j} D(s_i, s_j) \tag{1}$$

and the Voronoi diagram of S:

$$V(S) = \bigcup_{i \neq j} \overline{VR}(s_i, S) \cap \overline{VR}(s_j, S).$$
(2)

Being the most straightforward solution, a Voronoi diagram based on the Euclidean distance provides tolerable approximation in the urban environment if the points are located quite far apart within the uniform transportation network. In other cases, the results are significantly worse: for short distances, for a sparse roads network, in areas with irregular topography, under the presence of one way roads, or in application to suburbs stretched along the roads forming axon-like structures. Natural obstacles such as rivers, lakes, vegetation zones, ravines and mountains as well as urban (railways, highway, industrial zone, pipelines) play a particularly important role in complicating the route of movement between two points. The use of other metric functions may improve the accuracy; however, another problem arises: in some cases, the bisector dimension may be more than 1 (this is true even for the Manhattan distance  $ho(x,y) = \sum |x_1 - y_1| + |x_2 - y_2|$ ) applicable to the regular rectangular streets network.

In a number of works, the graph represents the streets network. The discrete network Voronoi diagram is then constructed while the metric used is the link between nodes (e.g., Yomono [6]). However, such models do not allow shortcuts, which are often used by pedestrians to shorten the routes. Aichholzer et al. [4] consider a plane with Manhattan distance and isothetic transportation network. There are also several works that use the generalized concept of Voronoi region (*needle*) proposed by Bae and Chwa [7].

The approach presented in this work is aimed at being applicable for the non-orthogonal street structure with curvilinear street segments. At the same time, as the ride sharing is spontaneous, we strive to avoid excessive model complexity; only walking to/from meeting points is assumed for the passenger. In addition, being flexible to the possibility of using the available network bandwidth data, the model should also work with the minimum information of this kind. Thus, we believe, the task of developing an optimal method for constructing a Voronoi diagram for a similar class of problems is to find a balance between complexity, accuracy and flexibility in using available data, as the latter may vary a lot in different regions.

The main idea of the approach presented below is to construct a discrete Voronoi diagram on a graph and then transform the obtained cells into the seeds or generator objects for the continuous Voronoi diagram on the plane. The latter represents the partition of the plane with a transportation network into proximity regions for the set of the given meeting points. The algorithm overview is presented below while individual steps are discussed in detail in sections III-IV.

- 1. Geospatial data preprocessing. The necessary information is: land use, coordinates of the roads.
- 2. Voronoi diagram construction on the graph representing transportation network.
- Cell of the constructed discrete diagram with their coordinates are used as the seeds for continuous Voronoi diagram on a plane.

As a result, continuous combined Voronoi diagram for the meeting points is constructed. This algorithm had been tested for Oldenburg city centre [1] and then applied for pedestrians on the area with radius 7 km around the city centre containing 79 meeting points (Fig. 1).

#### III. VORONOI DIAGRAM ON THE GRAPH

We consider the area of interest as a rectangular domain  $\Omega \subset \mathbb{R}^2$  containing the city transportation network, providing fixed routes. This network is modelled as a weighted graph G(V, E), where  $E = \{e_i\}$  is the set of edges, representing roads and streets and  $V = \{v_k\}$  are the graph vertices, corresponding to the intersections and the deadlocks. Nonnegative edge weights  $w(e_i)$  determine some proximity measure between the vertices connected by the edge  $e_i$ .

Depending on data availability, it can be, e.g., edge length or edge travel time. The latter depends on the segment's capacity, inclination, or traffic. Setting  $\rho_G(v_i, v_j)$  in an ordinary way as the weights sum of the shortest path between  $v_i$  and  $v_j$ , one can consider  $V_{\rho_G}$  as a metric space. Without additional



Fig. 1. Oldenburg with the suburbs (OpenStreetMap [5]).



Fig. 2. Graph representing transportation network with the meeting points.

constraints, it is true for the undirected graph as  $\rho_G$  always satisfies the symmetry axiom. It is not so in the directed graph case, nevertheless inward and outward Voronoi diagrams with corresponding asymmetric metric function can be considered instead. In application to the present nearest neighbor search inward diagram is a general approach that represents both types of movement: driving by car and walking. Taking into account that there are no one-way pedestrian roads and omitting the height difference for simplicity, one can assume that the undirected graph provides good representation for walking on foot in the mild regions.

Hence, we can build a Voronoi diagram on the graph G(V, E) with respect to the meeting points  $S = \{s_1, ..., s_k\} \subset V$  (Fig. 2).

The Voronoi diagram brakes up the set of vertices into the direct sum of the Voronoi cells  $V = V_1 \oplus ... \oplus V_k$ , where  $V_i = VR(s_i, S)$ . Let  $E_i(s_i)$  be a set of edges connecting vertices within  $V_i$ . Then  $E = E_1 \oplus ... \oplus E_k + E_0$ , where  $E_0$  is the set of "border" edges whose vertices belong to the different cells.

The following steps describe a computational algorithm for constructing a Voronoi diagram on a graph. The city transportation network representation as a graph G(V, E) is obtained from the OpenStreetMap (OSM) project geodata [5]. The project provides free editable geographic database of the world. In this work we use Python package Osmnx to download, model and project geospatial OSM data. The rest of the code is also written in Python using such packages as NumPy, Shapely, Matplotlib, Networkx, GeoPandas and others.

At the first step geospatial data of this region is downloaded and projected to Gauss-Krüger projection in which all further computations take place. Thus, current data structure appears as a weighted graph with the certain geometrical coordinates for nodes and edges. Second, locations of the meeting points are added to the set of graph vertices (Fig. 2). Since graph order for the individual town lets allows it, brute-force can be used for the Voronoi diagram construction:  $\forall v \in V$  find the distance on the graph  $\rho_G(v, s_i), i = \overline{1, k}$  using the Dijkstra algorithm. If  $s_j$  satisfies  $\rho_G(v, s_j) = \min \rho_G(v, s_i)$ , then  $v \in VR(s_j)$ . This computation can be easily and effectively parallelized as long as there is no need for data transfer between the threads. Set V is split up into disjoint subsets by the processor cores number. Then nodes of each subset are divided into the groups according to their proximity to a certain seed. Finally, the results are combined together. Finding terms  $V_k$  for direct sum decomposition of V allows to determine corresponding graph edges subsets  $E_k$  belonging to which clearly indicates the nearest seed - meeting point for each  $e \in E \setminus E_0$ .

#### IV. PLANAR VORONOI DIAGRAM

Constructed according to the previous section, the Voronoi diagram on the graph does not indicate the nearest meeting point for the surface points lying outside the graph edges. As graph V(G, E) can be considered not only as a topological structure but set of geometrical objects: points and lines with the certain coordinates, each subset  $E_m$  corresponds to the lines set  $E'_m$  on  $\mathbb{R}^2$ . It should be noted that in general two seeds may intersect (Fig. 4). Normally it happens under presence of the multi-level roads interchanges, tunnels, crossroads with the prohibition for movement in the certain direction but, in general, may have a connection with the roads congestion and bandwidth. Although such cases represent a small proportion of the total number of cases, the need to process them significantly complicates the procedure for bisector construction (section IV-C).



Fig. 3. Voronoi diagram on the graph.

#### A. Planar metric function

As long as there is no exact information about travel routes outside the transportation lines, it is natural to assume movement along a straight line in the direction of the nearest transportation network segment. However, this simplification does not take into consideration the presence of natural and man-made obstacles: buildings, fences, water reservoirs, ravines, vegetation and industrial zones, farmland as well as private and restricted access areas. In suburbs and rural surroundings such objects can occupy a large area, therefore bypassing them significantly complicates the route. On the one hand, it is possible here to consider a geometrical problem of building an optimal curvilinear route between a given point and a transportation network that does not intersect impassable regions. The length of such a route is then used as a metric function.

On the other hand, the problem of identifying impassable regions from generally available geospatial data can be more difficult than it seems. Although some obstacles can confidently be considered as impassible, for others it is hard to



Fig. 4. Geometrically overlapping cells of Voronoi diagram on the graph.

determine their real degree of obstruction. This applies to a lesser extent to movement by car, but is relevant enough for pedestrians who tend to take shortcuts. For example, taking a shorter route by moving through a vegetation zone may depend on vegetation type, density, soil type, time of year, weather, time of day (due to the illumination factor). Thereby, not only spatial, but short- and long- term time variation of site passability occurs. Even the water reservoirs can freeze in winter and become passable. Additional socio-behavioral aspects play a role in relation to the zones forming artificial obstacles. For them obstruction may depend on such factors: if they are actively used or abandoned; if there security guards and/or CCTV; the kind of fence around the perimeter; legal consequences of a violation. The same applies to the crossing the railways and highways outside the permitted spots.

It should be noted that there is likely a connection between shortcut usage and benefits of route reduction. In contrast, a high local crime rate can drastically reduce pedestrians' willingness to walk outside of the streets. Moreover, tendency to follow formal prohibitions varies in different cultures and regions [9]: while in some cases a prohibition sign is enough, in others even concrete fence is useless. It seems that upto-date information regarding the passibility of shortened or alternative routes should come via some pedestrians' feedback system. Satellite imagery can help with the determining of vegetation properties and recognition of footpaths. Nevertheless, leaving this approach for the future stage of work, we currently use the Euclidean distance as a metric function.

# B. Search for the equidistant points

Considering  $\{E'_1, ..., E'_k\}$  as seeds in (2) and Euclidean metric  $\rho_2$  as  $\rho$ , Voronoi diagram V(E') can be constructed. Obviously,  $\rho_2(M, E'_m)$  is the distance between  $M \in \mathbb{R}^2$  and the nearest to M point of  $E'_m$ .

The first step is to find the metric sphere  $S_r(E'_m)$  for  $E'_m$  with the given radius r. The metric sphere analytically obtained for the straight line segment consists of two couples of straight line segments and circular arcs. As far as even curvilinear roads are represented in  $E'_m$  as polygonal chains,  $S_r(E'_m)$  is expressed as the individual spheres union's perimeter. By the definition, for two seeds and any point  $M \in \Omega : M \in S_r(E'_m) \cap S_r(E'_n) \Rightarrow M \in \Lambda(E'_m, E'_n).$ Therefore, finding sufficient number of such equidistant points as equal radius spheres intersection, allows to determine with some precision the bisector within the domain through further interpolation. Let  $B'_r$  denote the set  $S_r(E'_m) \cap S_r(E'_n)$ . Giving to the radius r variation with some step:  $r_{k+1} = r_k + \Delta r$ (k = 0, 1, ...) we compute all coresponding metric spheres intersections  $B'_r$ . Here  $r_k \in [r_{min}, r_{max}]$ , where  $r_{min} =$  $\frac{1}{2}\rho_2(E'_m,E'_n)$  and  $r_{max}$  is the minimum radius  $r_k$  that satisfies the condition  $B'_{r_k} \not\subset \Omega$ . Choice for the  $\Delta r$  depends on two aspects. First, the set of obtained equidistant points should be adequate for the proper bisector line representation. Second, the excessive precision should be avoided to reduce computational complexity at this algorithm stage. For this reason, the variable radius increment step is chosen:  $\Delta r(k) = \Delta r_k$ :

$$\Delta r_k = \begin{cases} f \cdot \Delta r_{k-1} & \text{if } E'_m \cup E'_n \subset \bigcup_{l=m,n} O_r(E'_l), \\ \Delta r_0 & \text{otherwise} \end{cases}$$
(3)

The radius increment step remains constant unless both seeds are located within the open balls union, hereafter it grows geometrically. Fig. 6 (top) illustrates how it affects computed points distribution. In the computations below  $\Delta r_0 = 6$ meters and f = 1.25. As a result, for each pair  $E'_m, E'_n$  we obtain a set of equidistant points as combination:

$$B'(E'_m, E'_n) = \bigcup_{\forall r_k} B'_{r_k}(E'_m, E'_n) \tag{4}$$

#### C. Bisector construction

A goal of the current step is to construct a continuous bisector from the set of equidistant points  $B'_{mn} = B'(E'_m, E'_n) = \{Q_i\}_{i=0}^{N_B}, Q_i \in \mathbb{R}^2$ . Bisector  $\Lambda = \Lambda_{mn}$  constructed from  $B' = B'_{mn}$  should satisfy the following conditions:

- 1.  $\Lambda$  is a finite set of simple curves without selfintersections.
- 2. A contains maximum number of points from B'.
- 3. Each curve in  $\Lambda$  intersects  $\Omega$  in two points making the domain partition possible.
- 4. A does not intersect with  $E'_m$  and  $E'_n$ .

The problem of such line set construction is to separate points into groups (if necessary) and arrange them in each group in a correct order for interpolation.



Fig. 5. Equidistant points as equal radius metric spheres intersection.

Assuming  $E'_m \bigcap E'_n = \emptyset$ ,  $\Lambda = \Lambda_{mn}$  will consist of one line L that can be obtained with the following procedure. Let denote  $L_{\omega}$  as a tuple of points.

- 1. Select arbitrary the initial point  $Q_0 \in B' : Q_0 \in \Omega$ . Assign  $L_{\omega} = (Q_0); B' = B' \setminus \{Q_0\}.$
- 2. Find  $Q_1 \in B' : Q_1 \in \Omega$  that is the nearest to  $Q_0$  point in B'.

Set  $L_{\omega} = (Q_0, Q_1), B' = B' \setminus \{Q_1\}, n = 2.$ 

3. For  $L_w = (Q_{j_0}, ..., Q_{j_{n-1}})$  find  $Q_n^{\alpha}, Q_n^{\beta} \in B'$  such that: a)  $\rho_2(Q_n^{\alpha}, Q_{j_0}) = \min_{Q \in B'} \rho_2(Q, Q_{j_0}).$ 

If  $\rho_2(Q_n^{\alpha}, Q_{j_0}) < \rho_2(Q_n^{\alpha}, Q_{j_{n-1}})$  then place  $Q_n^{\alpha}$  as the first element in  $L_{\omega}$  and set  $B' = B' \setminus \{Q_n^{\alpha}\}$ .

- b)  $\rho_2(Q_n^{\beta}, Q_{j_{n-1}}) = \min_{Q \in B'} \rho_2(Q, Q_{j_{n-1}}).$ If  $\rho_2(Q_n^{\beta}, Q_{j_{n-1}}) < \rho_2(Q_n^{\beta}, Q_{j_0})$  then place  $Q_n^{\beta}$  as the last element in  $L_{\omega}$  and set  $B' = B' \setminus \{Q_n^{\beta}\}.$
- 4. Assign n = n + 1; repeat step [3.] until B' is not empty
- 5. Compute L as the linear interpolation of  $L_{\omega}$ .

In other words, the process of points arrangement is the sequential increment of the polygonal chain from the two ends. Testing shows that this approach, which is based on simple proximity, provides sufficient accuracy in the vast majority of cases. However, for some closely located seeds with irregular outlines containing combinations of convex and concave elongated segments it may lead to: skipping some of B' points. Also obtained with interpolation line L can: a) contain loops; b) intersect with the seeds. Thus, this resulting L requires examination and, if necessary, must be rebuilt. In exeptional cases in order to enhance the algorithm robustness, a simplified bisector can be constructed. A possible option in such a case is an analytically obtained straight line – bisector



Fig. 6. Computed equidistant points.

of the seeds centroids.

# D. Overlapping seeds processing

In this section case  $E'_m \cap E'_n \neq \emptyset$  will be covered (Fig. 7, 8). The above described procedure for the bisector construction does not work correctly under this condition. In the test performed for Oldenburg this was observed 9 times among 3081 pairs. The approach is the following: the procedure [1] – [5] from the previous section is performed recursively with the additional constraints for  $\Lambda = \{L_j\}_{j=1}^p$ :

$$L_j \bigcap_{j \neq i} L_i = \emptyset$$
 and  $L_j \bigcap E'_l = \emptyset$ . (5)

Here choice l = m or l = n is voluntary. The process of bisector construction for the intersected seeds is presented below. Numbers in square brackets refer to the algorithm steps for the individual line from section IV-C.

- I. Assign a value to l; Set j = 0.
- II. Set  $B' = B'_i$ .
  - 1. Perform step [1.].
  - Find Q<sub>1</sub> as in [2.]. Set B' = B' \ {Q<sub>1</sub>}. If Q<sub>0</sub>Q<sub>1</sub> intersects E'<sub>l</sub> then repeat the current step. Otherwise set L<sub>ω</sub> = (Q<sub>0</sub>, Q<sub>1</sub>).
  - 3. Let  $\xi \in \{\alpha, \beta\}$ . If  $Q_{j_0}Q_n^{\xi}$  does not intersect  $E'_l$  then perform for  $Q_n^{\xi}$  steps [3.a] or [3.b].
  - 4. Perform step [4.].
  - 5. Perform step [5.].
- III. Add L into A. Set  $B'_{j+1} = B'_j \setminus L_{\omega}, \ j = j + 1.$



Fig. 7. Equidistant points for the overlapping seeds: single intersection.



Fig. 8. Equidistant points for the overlapping seeds: multiple intersections.

IV. Repeat steps II, III until  $card(B'_i \cap \Omega) < K^*$ .

Using  $K^* = 0$  is possible although it reduces algorithm robustness. In certain cases, a few equidistant points may be left unused due to the algorithm simplifications. As a result of steps I-III, we obtain lines set  $\Lambda$ .

It is worth to mention that computations on the step of bisesector construction as well as on the step of searching for the equidistant points allows effective parallelization: one thread is allocated for each seeds pair.

# E. Voronoi cells construction

For the certain seed  $E'_m$  each computed bisector  $\Lambda_{mn}$  splits the domain in two parts:  $\Omega = \Omega_n^+ \oplus \Omega_n^-$ :  $E'_m \subset \Omega_n^+, E'_m \not\subset \Omega_n^-$ . If the bisector consists of a single line then it cuts the domain in two polygons. Otherwise multiple bisectors divide the domain into one simply connected region and one multiply connected region consisted of two or more subregions (Fig. 9). According to the Voronoi cell definition:

$$VR(E'_m, E') = \bigcap_{i=\overline{1,k}} \Omega_i^+.$$
 (6)

Making a reverse substitution  $E'_m \to E_m \to s_m$  we get  $VR(s_m, S)$  as the cells of the combined continuous Voronoi diagram based on metric functions  $\rho_2$  and  $\rho_G$ . In practice, due to the limited accuracy for bisector computation  $\Omega = \bigcup VR(s_m, S) \bigcup R$ . While cells overlay is not possible, the voids  $R = \{R_j\}$  between cells may occur. Simple procedure is used here to dispose of this areas of uncertainty: each void is merged with the cell that has the largest common border with it.



Fig. 9. Domain partition. Overlapping seeds case.

#### V. EVALUATION

As mentioned above, test computations were performed for the area with radius 7 km around Oldenburg city centre. The selected region consists of urban, suburban and rural areas. It also includes a variety of natural and man-made obstacles: highways, railways, water reservoirs, industrial and vegetation zones. For a given region with 79 selected meeting points, we construct analytically the standard Voronoi diagram based on the Euclidean distance (diagram I) and the combined one in a way described above (diagram II). After the cells construction through half-spaces intersection, 151 void region remained within the domain. Their total area  $\approx 0.1$  km<sup>2</sup> that is 0.05% of domain area. The area of the two largest voids is approx. 99% of total void area while 131 are smaller than 1 m<sup>2</sup>. After merging the voids with the computed cells the final diagram is obtained (Fig. 11).

For comparison of Voronoi diagrams of types I and II, the following value is used as measure of difference:

$$\Delta S = \frac{1}{S(\Omega)} \sum_{i} S(C_i^1 \setminus C_i^2) = \frac{1}{S(\Omega)} \sum_{i} S(C_i^2 \setminus C_i^1), \quad (7)$$

where  $C_i^1$  and  $C_i^2$  are the cells for the same seed in diagrams I and II correspondingly. For the considered example  $\Delta S \approx 18\%$ . One can expect the bigger difference for higher number of meeting points and, consequently, smaller cells. Also, for 3501 random locations uniformly distributed within the domain, we determine the nearest meeting point in three ways: a) from diagram I; b) from diagram II; c) by computing the routes to all the meeting points with the Openrouteservice engine [10] and detecting the meeting point corresponding to the minimum route length. The results are the following. The nearest meeting points obtained from diagrams I and II are not equal in 18% what is consistent with  $\Delta S$  value. Meeting points are different from the obtained with Openrouteservice for 17% (diagram I) and 10% (diagram II).

#### VI. CONCLUSION

There are several steps to be performed next in the context of this work. First, the potential of using additional bandwidth data must be analyzed. Second, impassable region processing must be implemented in planar Voronoi diagram construction. Third, the presented approach must be tested for the different regions and other methods of travelling, e.g., cars and bicycles. Nevertheless, at this stage, one can conclude that, despite the number of simplifications, the described algorithm provides more accurate results in comparison with a standard Voronoi diagram. At the same time, the processing of complex topography features requires further study since they are probably the main reason for the remaining imprecision. These include multi-level road crossings, tunnels, elongated geometric objects and natural obstacles.



Fig. 10. Classic Voronoi diagram.



Fig. 11. Combined Voronoi diagram.

#### REFERENCES

- A. Butenko and J. Marx Gómez, "Combined algorithm for Voronoi diagram construction in application to dynamic ride sharing," MOBIL-ITY 2022: The Twelfth International Conference on Mobile Services, Resources, and Users, pp. 5-8, Porto, Portugal, 2022.
- [2] M. Eilers et al., "An instant matching algorithm in the context of ridehailing applications, using isochrones and social scoring", In:, (Hrsg.), Informatic 2021.
- [3] Instaride, https://instaride.webflow.io, retrieved: November 2022.
- [4] O. Aichholzer, F. Aurenhammer, and B. Palop, "Quickest Paths, Straight Skeletons, and the City Voronoi Diagram", Discrete and Computational Geometry, 31. pp. 17-35, 2004,doi:10.1007/s00454-003-2947-0. Gesellschaft für Informatik, Bonn, pp. 103-114, 2021, doi:10.18420/informatik2021-007.
- [5] https://www.openstreetmap.org, retrieved: November 2022.
- [6] Yomono H. Yomono, "The Voronoi diagram on a network", Technical report, Nippon Systems Co, Tokyo, 1991.

- [7] S. W. Bae and K.-Y. Chwa, "Voronoi Diagrams with a Transportation Network on the Euclidean Plane," Int. J. Comput. Geometry and Appl., vol. 16, pp. 101-112, 2004, doi:10.1007/978-3-540-30551-4\_11.
- [8] https://github.com/gboeing/osmnx, retrieved: November 2022.
- [9] G. Bierbrauer, "Reactions to violation of normative standards: a crosscultural analysis of shame and guilt", International Journal of Psychology, vol. 27, pp. 181-193, 1992, doi: 10.1080/00207599208246874.
- [10] https://openrouteservice.org, retrieved: November 2022.

# Base Station Assisted (BSA) Reinforcement Learning for Resource Allocation in Wireless Industrial Environments

Idayat O. Sanusi and Karim M. Nasr

Faculty of Engineering and Science, University of Greenwich, Kent, ME4 4TB, UK {i.o.sanusi, k.m.nasr}@gre.ac.uk

Abstract— Device-to-Device (D2D) enabled cellular networks are a promising solution for Ultra-Reliable Low-Latency Communication (URLLC) systems. Integrating D2D into future wireless industrial networks and next-generation manufacturing can support the creation of massive machinetype wireless connections. In this paper, we present a Base Station Assisted (BSA) reinforcement learning approach for resource allocation in a D2D-enabled cellular network targeting smart manufacturing and Industry 4.0 applications. A distributed local Q-table is used for the D2D agents to prevent global information gathering and a stateless Qlearning approach is adopted to reduce the complexity of learning and the dimension of the O-table. The O-tables of the D2D agents are then uploaded to the Base Station (BS) for the resource allocation to be implemented centrally. Simulation case studies show that the presented semi distributed BSA technique results in reduced signalling overheads and a good Quality of Service (QoS) across the network compared to other conventional schemes.

Keywords—Fifth Generation (5G) and beyond networks; Radio Resource Management (RRM); Distributed Algorithms; Device-to-Device Communication (D2D); Reinforcement Learning.

# I. INTRODUCTION

The increasing growth in the number of wireless smart devices and applications necessitates novel and efficient Radio Resource Management (RRM) schemes to address the different challenges faced. Device to Device (D2D) communication is considered one of the key technologies for 5G and beyond networks aiming to provide improvements in performance metrics such as throughput, spectrum, and energy efficiency especially for new verticals such as smart manufacturing and Ultra Reliable Low Latency Communication (URLLC) use cases, e.g., in wireless industrial applications [1]. Machine learning and artificial intelligence techniques are some of the main techniques currently gaining increased interest to realise the expectations of future generation wireless systems [2]-[3].

Spectrum access where a cellular and D2D users share the same resources can potentially result in improved spectrum efficiency. However, if shared resource allocation is not properly coordinated, mutual interference between cellular and D2D links may degrade the Quality of Service/Quality of Experience (QoS/QoE) of end-users.

Future wireless networks are characterised by a high density of devices and dynamic environments with rapidly changing Channel State Information (CSI). Centralised and distributed schemes are two RRM approaches used to allocate resources to users. In a centralised scheme, the global acquisition of CSI by a centralised controller (e.g., a Base Station (BS)) often incurs high signaling overheads and computation complexity which, tend to increase with the number of users, therefore making it impractical to deploy. Furthermore, RRM problems are often formulated as optimisation problems where the QoS requirements are modelled as the constraints. These optimisation problems are often complex and difficult to solve directly. A distributed approach does not need a central entity. Resource allocation is executed by users, therefore reducing the amount of information exchange, computations, and processing by the base station, and resulting in improved QoS across the network.

Game theory and machine learning are important techniques that can be used to realise a distributed RRM scheme. Matching theory, which, has been used to solve assignment or pairing problems between two distinct sets of players with diverse QoS objectives [4], may get complex in a multiuser scenario with rapidly changing channel conditions using full CSI, as in [5].

Reinforcement Learning (RL) has been explored to address RRM problems in dynamic environments [6]-[7]. RL is a machine learning approach, well-suited to support decision making in 5G-and-beyond networks with uncertainties, for example, in distributed resource allocation with unknown or partial information of network conditions. Q-learning is a reinforcement learning technique that uses a look-up table, known as Q-table, to determine an optimal strategy to adopt, by storing the values used to compute the maximum expected future rewards for actions taken at each state. A large number of agents, states and actions can lead to a high-dimension Q-table which, may result in slow convergence and limit the practical applications due to the high memory requirements [8]. These challenges can be addressed by using Deep Reinforcement Learning (DRL) which, uses deep neural networks to approximate the tables [9]. However, DRL is associated with high complexity and large learning data [10]-[11].

RL has been widely investigated to study intelligent power level and spectrum channel allocations for D2Denabled cellular networks in a multi-agent environment. The work in [12] formulated the resource allocation problem as a stochastic non-cooperative game among D2D users. However, the QoS requirements of cellular users sharing the same frequency bands with D2D users were not considered in the reward model. In [13], a multi-agent actor-critic framework was proposed which, involves cooperation between users and sharing of all historical information (states, actions, and policies) in a centralised training scheme to ensure stability. This will consequently increase the amount of signalling overheads and information exchange. In [12]-[15], the reward function captured the QoS metric of cellular users in a centralised Q-learning approach, which, also leads to increased signalling overheads.

In this paper, we present a semi-distributed reinforcement learning scheme for spectrum resource sharing of D2D Users (DUEs) and Cellular Users (CUEs) targeting smart manufacturing environments and URLLC networks. This semi distributed approach relies on two phases. First, a decentralised training of agents is implemented. This is followed by Q-tables being uploaded to the base station for final resource allocation. The reward function is modeled in such a way that there is no information exchange related to other agents' actions or rewards. To address the problem of the 'curse of dimensionality' associated with Q-learning, a stateless Q-learning approach is adopted to reduce the dimension of the Q-table, nonetheless capturing the QoS demands of the D2D users. The main contributions of this work are summarised below:

- A hybrid RRM scheme with distributed D2D training and a centralised channel allocation is presented with an advantage of reduced signalling overheads compared with conventional centralised approaches. This hybrid RRM scheme relies on stateless reinforcement learning algorithm is presented, where there is no state transition, to ensure a reduced dimension of the stateaction mapping, nevertheless capturing the key performance metrics of the DUEs. With this technique, there is a decrease in complexity and signalling overheads making scheme adaptable to high-density networks.
- In previous works [16]-[18], the QoS of cellular users is captured by integrating it in the state space or reward function of the D2D users. Rather than the BS exchange the QoS estimation of the CUE with the DUE at each time slot, a Q-table for the CUEs is maintained and updated.
- Numerical simulations are used verify the performance of the presented algorithm in comparison to other approaches in terms of achieved throughput, signalling overheads and complexity.

The paper is organised as follows: The system model and problem formulation are presented in Section II. In Section III, a stateless reinforcement learning algorithm for base station-assisted resource allocation is presented. Section IV presents simulation case studies and results. The main conclusions and directions for future work are summarised in Section V.

# II. SYSTEM MODEL AND PROBLEM FORMULATION

We consider D2D and cellular users coexisting within a cellular network for uplink spectrum-sharing as illustrated in Fig. 1. There are *N* Cellular Users (CUEs) represented by a set  $C = \{c_1, ..., c_i, ..., c_N\}$  and *M* D2D Users (DUEs) denoted by a set  $D = \{d_1, ..., d_j, ..., d_M\}$  deployed randomly within the coverage of the base station in a single cell system.



Figure 1. An illustration of a D2D enabled cellular network

The DUEs can autonomously select a Resource Block (RB) denoted by a set  $K = \{k_1, ..., k_i, ..., k_N\}$ , from a pool of radio resources [18]19, which, can overlap with that of the CUEs for the benefit of reuse gain. The cellular users have strict performance requirements in the form of minimum Signal-to-Interference Plus-Noise-Ratio (SINR) values to guarantee their throughput. The D2D links also have minimum SINR thresholds to guarantee their throughput demands, in addition, to the reliability and delay requirements.

We assume that each CUE has been pre-allocated a resource block. The transmit power of the CUEs and DUEs are denoted by  $P_{c_i}$  and  $P_{d_j}$  respectively.  $g_{c,B}$ ,  $g_{d_T,B}$ ,  $g_{d_T,d_R}$  and  $g_{c,d_R}$  are the channel gains of cellular communication from the CUE  $c_i$  to the BS, the interference link from the DUE transmitter  $d_T$  to the BS, the D2D communication link

from the DUE transmitter  $d_T$  to the receiver  $d_R$  and the interference link form the CUE transmitter to the DUE receiver  $d_R$ , respectively. The channel gain comprises small-scale fading which, is assumed to be exponentially distributed with a unit mean and large-scale fading which, includes pathloss and shadowing with log-normal distribution.

The instantaneous received SINR at the BS from *i*th CUE and *j*th DUE over *i*th sub-channel at time slot t is given as [1]:

$$\Gamma_{c_i}(t) = \frac{P_{c_i}g_{c,B}(t)}{\sigma^2 + \sum_{d_j \in D} \lambda_j^i(t) P_{d_j}g_{d_T,B}(t)}$$
(1)

$$\Gamma_{d_j}(t) = \frac{P_{d_j}g_{d_T,d_R}(t)}{\sigma^2 + \sum_{c_i \in C} \lambda_j^i(t) P_{c_i}g_{c,d_R}(t)}$$
(2)

 $\lambda_j^i \in \{0,1\}$  denotes the binary resource reuse indicator,  $\lambda_j^i = 1$  implying that the *j*th DUE selects *i*th CUE subchannel at time slot *t* and  $\lambda_j^i(t) = 0$  otherwise. We assume that each DUE can access only one CUE sub-channel i.e.,  $\sum_N \lambda_j^i \leq 1$  and each CUE sub-channel is accessed by only one DUE i.e.,  $\sum_M \lambda_j^i \leq 1$ . The data rates of the *i*th CUE and *j*th DUE is at time slot *t* given by:

$$\mathbf{T}_{c_i}(t) = W_i \log_2 \left( 1 + \Gamma_{c_i}(t) \right), \tag{3}$$

$$T_{d_j}(t) = W_i \log_2 \left( 1 + \Gamma_{d_j}(t) \right), \tag{4}$$

where  $W_i$  is the bandwidth of each resource block. The variance of the Additive White Gaussian Noise (AWGN) is denoted by  $\sigma^2$ .

The channel gains for the different links (q, r) be expressed as follows:

$$\begin{cases} g_{c,B} = G_1 \gamma_{c,B} \chi_{c,B} L_{c,B}^{-\alpha_1} \triangleq \zeta_{c,B} L_{c,B}^{-\alpha_1} \\ g_{d_T,B} = G_2 \gamma_{d_T,B} \chi_{d_T,B} L_{d_T,B}^{-\alpha_2} \triangleq \zeta_{d_T,B} L_{d_T,B}^{-\alpha_2} \\ g_{d_T,d_R} = G_3 \gamma_{d_T,d_R} \chi_{d_T,d_R} L_{d_T,d_R}^{-\alpha_3} \triangleq \zeta_{d_T,d_R} L_{d_T,d_R}^{-\alpha_3} \\ g_{c,d_R} = G_4 \gamma_{c,d_R} \chi_{c,d_R} L_{c,d_R}^{-\alpha_4} \triangleq \zeta_{c,d_R} L_{c,d_R}^{-\alpha_4} \end{cases}$$
(5)

where  $G_r$  is the pathloss constant,  $\gamma_{q,r}$  is the small-scale fading gain due to multipath propagation and assumed to have an exponential distribution with unit mean. The largescale fading comprises pathloss with exponent  $\alpha_r$  and shadowing which, has a slow fading gain  $\chi_{q,r}$  with a lognormal distribution.  $L_{q,r}$  is the distance from terminal q to terminal r [20].

The channel gain  $g_{d_T,d_R}$  and  $g_{c,d_R}$  can be estimated at the DUE receiver,  $d_R$  and made available at its transmitter,  $d_T$  instantaneously [19]. Similarly,  $g_{c,B}$  and  $g_{d_T,B}$  can be obtained at BS through local information since uplink transmission is considered.

The reliability of the DUE  $d_j \in D$ ,  $\xi_{d_j}(t)$ , is defined as the probability of packet delay exceeding a predefined delay bound,  $l_{d_j,\max}$ , on channel *i* at slot *t* is less than a threshold [21]. The objective of the system is to maximise the total throughput,  $T_R$ , of paired CUEs and DUEs while satisfying the QoS demands. The optimisation problem and constraints are described in (6).

$$\mathbf{Max}_{\lambda_j^i} \ T_R = W_i(\lambda_j^i(\sum_{c_i \in C} \log_2(1 + \Gamma_{c_i}) + \sum_{d_j \in D} \log_2(1 + \Gamma_{d_j})))$$
(6)

subject to

$$\lambda_j^i \Gamma_{c_i} - \Gamma_{c_i,\min} \ge 0 \qquad \qquad \forall c_i \in C \tag{6a}$$

$$\Pr\left(l_{d_j} > l_{d_j,\max}\right) < 1 - \xi^*_{d_j} \quad \forall \ d_j \in D \tag{6b}$$

$$\sum_{c_i \in C} \lambda_j^i \leq 1 \qquad \forall a_j \in D \qquad (6c)$$
$$\sum_{d_i \in D} \lambda_i^i \leq 1 \qquad \forall c_i \in C \qquad (6d)$$

The minimum SINR,  $\Gamma_{c_i,\min}$ , to guarantee the throughput requirement of the CUEs is defined in constraint (6a). Constraint (6b) takes into account reliability and delay, where  $l_{d_j}$  is the packet delay constraint for packet transmission of DUE  $d_j$ . The expression captures the fact that the end-to-end delay should be less than  $l_{d_j,\max}$  with a probability of at least  $1 - \xi_{d_j}^*$ . Constraints (6c) and (6d) are channel association criteria. The reliability of the DUE links in (6c) is evaluated using an empirical estimation of number of packets transmitted similar to [21], from  $d_T$  to  $d_R$  whose delay is within the budget  $l_{d_j,\max}$  over the total number of packets sent to  $d_R$  at time slot t i.e.,

$$\xi_{d_j}(t) = 1 - \Pr\left(l_{d_j} > l_{d_j, \max}\right) \approx 1 - \frac{L_{d_j}(t)}{B_{d_j}(t)} \cong \frac{L_{d_j}(t)}{B_{d_j}(t)}, \quad (7)$$

where  $L_{d_j}(t)$  is the number of packets for which,  $l_{d_j} > l_{d_j,\max}$  and  $L'_{d_j}(t)$  is the number of packets transmitted with  $l_{d_j} \leq l_{d_j,\max}$  (or number of packets delivered within the delay bound).  $B_{d_j}(t)$  is total packet transmitted by DUE  $d_j$  at time slot t. Reliability can also be measured in terms of the outage probability, which, is the probability that the measured SINR is lower than a minimum is less than a predefined threshold. The expression of the outage probability of *j*th DUE conditioned on the selected *i*th channel at time slot t is given below [22].

$$p_{R}(t) = \Pr\left(\Gamma_{d_{j}} \leq \Gamma_{d_{j},\min}\right)$$
$$= 1 - \frac{P_{d_{j}}g_{d_{T},d_{R}}\exp\left(-\frac{\Gamma_{d_{j},\min}\sigma^{2}}{P_{d_{j}}g_{d_{T},d_{R}}}\right)}{P_{d_{j}}g_{d_{T},d_{R}} + \Gamma_{d_{j},\min}P_{c_{i}}g_{c_{n},d_{R}}} \leq p_{R_{0}}, \qquad (8)$$

where  $p_R(t)$  is the measured outage probability of DUE  $d_j$ at time slot t and  $p_{R_0}$  is the maximum tolerable outage probability of  $d_j$ .

The reliability of the DUE in terms of outage probability is expressed as [21]:

$$\xi_{d_i}(t) = 1 - p_R(t). \tag{9}$$

Transmission delay is given as the ratio of packet size transmitted within delay bound to the transmission rate [23]. From (7), (8) and (9) the transmission delay of *j*th DUE using the *i*th RB is formulated as:

$$l_{d_j}(t) = \frac{L'_{d_j}(t)}{W_i \log_2(1 + \Gamma_{d_i})}.$$
 (10)

At each time slot t, the resource allocation system implements two functions, namely:

i) determining the SINR,  $\Gamma_{c_i}$  for the *i*th CUE and the SINR  $\Gamma_{d_j}$  that the *j*th DUE to ensure that the minimum SINR and target reliability  $\xi_{d_j}^*$  thresholds are achieved and

ii) allocating RBs to *j*th DUE so that  $T_R$  is maximised.

The resource allocation optimisation problem for D2D communication in a cellular network is NP hard and a direct solution is not feasible. We present a base station-assisted resource allocation scheme which, adopts a semi-distributive RRM approach.

# III. STATELESS REINFORCEMENT LEARNING FOR BASE STATION-ASSISTED RESOURCE ALLOCATION

The goal of the agents is to maximise throughput in a D2D-enabled cellular network. At each time slot t, a DUE observes a state  $s^t$  and takes an action  $a^t$  from the action space (i.e., select an RB  $k_i$ ), according to a policy  $\pi$ . Q-learning enables an agent to determine the optimal strategy that maximises its long term expected cumulative reward. The Q-value is updated as follows [23]:

$$Q^{t+1} = \begin{cases} Q^{t}(s^{t}, a^{t}) + \sigma \left[ r^{t} + \eta \max_{a} Q^{t}(s^{t+1}, a^{t+1}) - Q^{t}(s^{t}, a^{t}) \right] \\ & \text{if } s = s^{t}, \ a = a^{t} \\ Q^{t}(s^{t}, a^{t}), & \text{otherwise} \end{cases}$$
(11)

where  $\sigma \in [0,1]$  is the learning rate. With  $\sigma = 0$ , the Q-values are never updated, hence no learning has taken place; setting  $\sigma$  to a high value such as means that learning can occur quickly and  $0 \le \eta \le 1$  is the discount factor used to balance immediate and future reward [24].

The state space, action space and rewards function in the learning environment are defined as follows:

1) State Space: The state observed by DUE  $d_j \in D$ ,  $S_{d_j}^i(t)$ , using resource block RB  $k_i$  at time slot t is defined by three variables, resulting in eight possible states as defined in Table I.

$$S_{d_j}^{i}(t) = \left\{ S_{\Gamma_{d_j}}^{i}, S_{\xi_{d_j}}^{i}, S_{l_{d_j}}^{i} \right\},$$
(12)

where  $S \in S_{d_j}^i = \{0,1\}$ .  $S_{\Gamma_{d_j}}^i(t)$  indicates the interference level and is defined as:

$$S_{\Gamma_{d_j}}^i(t) = \begin{cases} 1 & \Gamma_{d_j}(t) \ge \Gamma_{d_j,\min}, \\ 0 & \text{otherwise}, \end{cases}$$
(12a)

 $S_{\xi d_j}^i(t)$  indicates the level of reliability and is defined as:

$$S^{i}_{\xi_{d_{j}}}(t) = \begin{cases} 1 & \xi_{d_{j}}(t) \ge \xi^{*}_{d_{j}} \\ 0 & \text{otherwise} \end{cases}, \quad (12b)$$

 $S_{ld_j}^i(t)$  indicates the packet transmission time and is defined as:

$$S_{l_{d_j}}^i(t) = \begin{cases} 1 & l_{d_j}(t) \le l_{d_j,\max}, \\ 0 & \text{otherwise} \end{cases},$$
(12c)

TABLE I. State Space for DUEs

$S^i_{\Gamma_{d_j}}$	$S^i_{\xi_{d_j}}$	$S^i_{l_{d_j}}$	$S_{d_j}^i$
0	0	0	0
0	0	1	0
0	1	0	0
0	1	1	0
1	0	0	0
1	0	1	0
1	1	0	0
1	1	1	1

The state-action dimension is reduced by adopting a stateless learning approach. For the considered scenario, an action  $a_i \in A$  taken by an agent will result in the end of an episode i.e., states 0 and 1 are terminal states, where  $S_{d_j}^i(t) = 1$  is the goal state of the DUEs. Therefore, the learning environment can be modelled entirely using a stateless Q-learning i.e., action-reward only since the state transition is not required. An agent can choose its action based solely on its Q-value. The updated Q-value of the chosen action is based on the current Q-value and the

immediate reward from selecting that action. The update function in (11) is re-formulated as follows:

$$Q^{t+1}(a^{t}) = \begin{cases} Q^{t}(a^{t}) + \sigma[r(a^{t}) - Q^{t}(a^{t})], & \text{if } a = a^{t} \\ Q^{t}(a^{t}), & \text{otherwise} \end{cases}$$
(13)

where  $r(a^t)$  is the immediate reward of selecting a.

In contrast to the standard Q-value update function in (11), it can be seen in (13) that not only the state-action formation (s, a) is not necessary, but also the information of the next state  $s^{t+1}$  is not required because the actions lead to a terminal state. Therefore, the Q-table is defined in terms of the actions only and updated using the immediate reward. This results in  $1 \times |N|$  dimension Q-table for *j*th DUE. This method reduces the learning complexity and the Q-table dimension.

The traditional cellular users in the network need to be protected from the interference caused by the DUEs for their minimum SINR to be satisfied. This may be achieved by integrating the SINR of the CUE,  $\Gamma_{c_i}$  in the state space or by reward function modelling. This way, the DUEs can obtain the information from the BS at time slot t as in [17]-[18], [25]; hence, the DUEs get a reward if the CUE SINR  $\Gamma_{c_i} \ge \Gamma_{c_i,\min}$ , and a penalty otherwise. Rather than the BS exchange the measured CUE SINR,  $\Gamma_{c_i}$ , with the DUEs for every action  $a^t$  taken at each time slot, we adopt a scheme in which, the BS keeps a look-up table of the *i*th CUE based on the actions on the DUEs. Therefore, the Q-table for the *i*th CUE is  $1 \times |M|$  considering a stateless Q-learning structure.

2) Action Space: The action space of DUE  $d_j \in D$  is a set of all actions denoted by  $A = \{a_{1,i}^t, \dots, a_{i,i}^t, \dots, a_{N}^t\}$ , where  $a_i^t$  is the action taken by  $d_j \in D$  at time slot t and defined as the selection of an RB  $k_i$ .

3) Action-Selection Strategy: There are methods to select an action based on the current evaluation of the Q-value at every time slot t using a policy denoted by  $p_{d_j}^t$ . These methods are used to balance the exploration and exploitation of actions taken by the agents [26]. Epsilon greedy ( $\varepsilon$ -greedy) is one of the methods of choosing an optimal Q-value and described as follows:

$$p_{d_j}^t = \begin{cases} \operatorname{argmax} Q(a) & \operatorname{probability} 1 - \varepsilon \text{ (exploitation)} \\ \operatorname{Random action} & \operatorname{probability} \varepsilon \text{ (exploration)} \\ (14) \end{cases}$$

where  $\varepsilon$  is the exploration rate with  $0 \le \varepsilon \le 1$ . The exploration rate is the probability that the agents will explore the environment rather than exploit it.  $\varepsilon \to 1$  results in greater exploration whereas  $\varepsilon \to 0$  means greater exploitation.

4) Reward Function: The reward function is modelled such that it relies only on local observations and can be

implemented in a distributive manner. The rewards of the *j*th DUE and *i*th CUE for taking an action  $a_i^t$  is expressed in terms of the achievable throughput using the Shannon capacity formula. Thus, the reward is directly related to the objective function of the optimisation problem.

Equation (15) shows that *j*th DUE only gets a reward when all the state variables are 1 (i.e., the minimum QoS demands are met), while *i*th CUE gets a reward if its minimum SINR is satisfied at each time slot for the action taken by *j*th DUE. From the reward function defined above, learning can be implemented independently in a decentralised manner such that each agent maintains a local Q-table. There is no information exchange relating to other agents' actions or rewards and no cooperation is needed between the agents, which, results in reduced signalling overheads and reduced complexity compared with a centralised Q-learning approach.

$$r_{d_j}(a^t) = \begin{cases} T_{d_j}^k(t) & S_{d_j}^i(t) = 1\\ 0, & S_{d_j}^i(t) = 0 \end{cases},$$
 (15a)

$$r_{c_i}(a^t) = \begin{cases} T_{c_i}^k(t) & \Gamma_{c_i} \ge \Gamma_{c_i,\min} \\ 0, & \text{otherwise} \end{cases}$$
(15b)

The Q-value of the *j*th DUE for selecting *i*th RB at time slot *t* is updated as follows:

$$Q_{d_{j}}^{i}(a^{t}) = \begin{cases} Q_{d_{j}}^{i}(a^{t}) + \sigma \left[ r_{d_{j}}(a^{t}) - Q_{d_{j}}^{i}(a^{t}) \right], \text{ if } a = a^{t} \\ Q_{d_{j}}^{i}(a^{t}), & \text{otherwise} \end{cases}$$
(16a)

Similarly, the Q-value of the *i*th CUE for action taken by the *j*th DUE is updated as follows:

$$Q_{c_{i}}^{j}(a^{t}) = \begin{cases} Q_{c_{i}}^{j}(a^{t}) + \sigma [r_{c_{i}}(a^{t}) - Q_{c_{i}}^{j}(a^{t})], \text{ if } a = a^{t} \\ Q_{c_{i}}^{j}(a^{t}), & \text{otherwise} \end{cases}$$
(16b)

From (16), it can be seen that after the training, the Q-table of the *j*th DUE,  $Q_{d_j}(a)$ , will return  $Q_{d_j}^i(a) = 0$  for its action on *i*th RB that do not meet its QoS requirements. Similarly, the Q-table of the *i*th CUE,  $Q_{c_i}(a)$ , will return  $Q_{c_i}^j(a) = 0$  for the action of *j*th DUE on *i*th RB that do not meet its QoS requirements.

The BSA algorithm summarised in Algorithm I, aims to optimise the achieved system throughput. After the training phase, each DUE loads its Q-value table,  $Q_{d_j}(a)$ , to the BS for centralised matching. The BS will then allocate cellular resource blocks to D2D links in such a way that spectrum sharing is optimised, network throughput is maximised and there is no need for information exchange between the UEs to find a suitable candidate.

Algorithm I: The BSA Reinforcement Learning Algorithm

1: Initialise the action-value function for the DUEs

$$Q_{d_j}(a) = 0 | Q_{d_j}(a) \equiv Q_{d_j}^i(a^t), i = 1, 2, \dots, N \quad \forall d_j \in D$$

2: Initialise the action-value function for the BS for the actions of the *j*th DUE on the *i*th RB

 $\begin{bmatrix} Q_{c_i}(a) = 0 | Q_{c_i}(a) \equiv Q_{c_i}^j(a^t), \ j = 1, 2, \dots, M \end{bmatrix} \quad \forall \ c_i \in C$ for  $d_i \in D$   $1 \le j \le M$  do 3: 4: while not converge do 5: generate a random number  $x \in \{0,1\}$ 6: if  $x < \varepsilon$  then 7: Select action  $a_i^t$  randomly 8: else Select action  $a_i^t = \operatorname{argmax}_{a \in A} Q_{d_i}(a^t)$ 9: 10: end

11: Evaluate  $\xi_{d_j}$ ,  $\Gamma_{d_j}$  and  $l_{d_j}$  of  $d_j \in D$  for the action  $a^t$ 

12: Measure the SINR,  $\xi_{c_i}$ , of CUE  $c_i \in C$  for the action  $a^t$  taken by  $d_i \in D$ 

13: Observe immediate reward of  $d_i \in D$  and  $c_i \in C$ ,

14: Update action-value for action of  $d_j \in D$  on the ith RB  $Q_{d_i}^i(a) = Q_{d_i}^i(a) + \sigma \left[ r_{d_i}(a^t) + Q_{d_i}^i(a) \right]$ 

15: Update action-value for  $c_i \in C$  for action  $a^t$  of *j*th DUE  $Q_{c_i}^j(a) = Q_{c_i}^j(a) + \sigma [r_{c_i}(a^t) + Q_{c_i}^j(a)]$ 

16: end while

18: Load  $Q_{d_i}(a)$  to the BS  $\forall d_j \in D$ 

19: for  $d_j \in D$   $1 \le j \le M$  do

20: Obtain 
$$Q(a) = \left\{ Q_{d_j}^i(a), Q_{c_i}^j(a) \right\}$$
  $i = 1, 2, ..., N$   
21:  $\overline{Q}(a) \subseteq Q(a) \mid \left\{ Q^i(a), Q^j(a) \right\} \in \mathbb{R}^+$  where  $\mathbb{R}^+$ 

21:  $Q(a) \subseteq Q(a) | \{Q_{d_j}^i(a), Q_{c_i}^j(a)\} \in \mathbb{R}^+, \text{ where } \mathbb{R}^+$ positive real number

22:  $Q_{\text{TOT}} = Q_{d_j}^i(a) + Q_{c_i}^j(a) \quad \forall q \in \overline{Q}(a)$ 23: end for

24: Set up a list for unmatched DUE  $D_u = \{d_j : \forall d_j \in D_u\}$ 25: while  $D_u \neq \emptyset$  do

26: Rank  $D_u$  in increasing order of  $|0 \overline{Q}(a)|$ 

27: Start DUE 
$$d_j \in D_u$$
:  $Q(a) \neq \emptyset$  with the least  $|Q(a)|$ 

28:  $c_i^* = \max_{r_i \in R} Q_{\text{TOT}}$ 29:  $D_u = D_u - d_j$ 

- $\begin{array}{ll} 29: & D_u = D_u d_j \\ 30: & \overline{Q}(a) = \overline{Q}(a) \backslash c_i^* & \forall d_{j'} \in D_u | j' \neq j \end{array}$
- 31: end while

#### IV. SIMULATION CASE STUDY AND PERFORMANCE EVALUATION

The performance of the BSA approach described in Section III, is verified by considering a single-cell network in an industrial scenario. The simulation set-up and channel models are as described in [1] and summarised in Tables II and III. The network dynamics is captured by generating the channel fading effects randomly. The throughput is the main metric used to evaluate the performances of the algorithms. The performance of BSA is compared with centralised optimisation and the game theoretic Deferred Acceptance (DA) techniques [1][20].

#### A. Throughput Performance

The throughput performance of matched DUEs as a function of the number of DUEs in the system M, is shown in Fig. 2. It can be concluded that the sum throughput of the DUEs increases with the number of cellular users M for all the considered algorithms. As expected, the number of admitted DUEs increases with the introduction of new DUEs to the system, but unchanged if a valid cellular resource-sharing partner cannot be found because the minimum QoS requirements are not satisfied. The performances of centralised and BSA approaches are comparable, while the DA method shows the least performance. The BSA algorithm outperforms the DA algorithm by up to 9.69% increase in the DUE throughput performance. However, it is semi-distributive as the final resource allocation is implemented by the BS whereas the DA approach is decentralised (the channel selection is usercentric, and no BS intervention is necessary to achieve autonomy). Players can make their resource allocations choices to maximise their individual throughput and ultimately achieve system stability. The performance of the sum throughput of the matched UEs (that is valid pairings between CUEs and DUEs) with respect to the number of cellular users M is shown in Fig. 3. The sum throughput increases with M. The BSA approach indicates better performance at  $M \leq 35$  with up to 12.05% increase in sum throughput compared to the centralised approach, while the centralised approach performed better at M > 35 with up to 9.39% increase in throughput. The DA algorithm again shows the least performance compared to the BSA technique.

The effects of the outage probability of the DUE,  $p_{R_0}$ , and delay threshold of the DUEs,  $l_{d_j,\text{max}}$  on the sum rate of the matched UEs for all algorithms are shown in Fig. 4 and Fig. 5, respectively. The sum throughput of the matched UEs increases with  $p_{R_0}$  and  $l_{d_j,\text{max}}$ . This is because higher  $p_{R_0}$ causes the interference from the CUEs to be more tolerable by the DUEs, therefore making potential CUE-DUE pairing possible. Similarly, higher  $l_{d_j,\text{max}}$  increases the sum throughput at fixed outage probability and payload since the delay requirement is less stringent. More DUEs are able to

. = . . .

satisfy the delay constraint and the number of admitted DUEs is increased.

# TABLE II. MAIN SIMULATION PARAMETERS [1][20][27]

Parameter	Value
Carrier frequency, $f_c$	2GHz
System bandwidth	10MHz
Number of resource blocks (RB), K	50
RB bandwidth	180 kHz
Maximum CUE transmit power, $P_{c_i,max}$	23dBm
Maximum DUE transmit power, $P_{d_j, max}$	13dBm
D2D distance, $L_{d_T, d_R}$	$10m \le L_{d_T,d_R} \le 20m$
CUE SINR Threshold, $\Gamma_{c_i,\min}$	7 dB
DUE SINR Threshold, $\Gamma_{d_j,\min}$	3 dB
Noise power density	−174 dBm/Hz
Number of CUEs, N	50
Number of DUEs, M	50
Reliability for DUE, $p_{R_0}$	10 <sup>-5</sup>
Exploration rate, $\varepsilon$	0.7
Learning rate, $\sigma$	0.9
DUE Maximum Delay, $l_{d_j,\max}$	50ms
DUE Message Size, $B_{d_i}$	15kB

TABLE III. CHANNEL MODEL FOR LINKS [28]-[30]

Parameter	In-factory	UE-UE link	BS-UE link
	DUE link		
Pathloss	$36.8 \log_{10}(d[m$	$40 \log_{10}(d[m])$	$37.6 \log_{10}(d[m])$
model	+ 35.8	+ 28	+ 15.3
Shadowing	4dB	6dB	8dB
Fast fading	Rayleigh	Rayleigh	Rayleigh



Figure 2. Sum-rate of matched DUEs with varying number of DUEs, M in the System, for N = 50



Figure 3. Sum Throughput of matched UEs as a function of the number of DUEs M, in the system, for N = 50



Figure 4. Effect of the DUE outage ratio  $p_{R_0}$ , on the sum throughput of matched CUE-DUE pair for N = M = 50,  $l_{d_i,\max} = 50$ ms



Figure 5. Effect of the delay bound,  $l_{d_j,\text{max}}$  on the sum throughput of matched CUE-DUE pair for N = M = 50,  $p_{R_0} = 10^{-5}$ 

#### B. Signalling Overheads and Complexity Analysis

We now evaluate and compare the signalling overheads and computation complexity of the investigated algorithms. Signalling overheads are evaluated in terms of the level of involvement of the BS and User Equipment (UE), i.e., BS-UE communication. The signalling overhead evaluated is an aggregation of contributions of channel information acquisition and information exchange by the BS-UE links. The number of iterations T depends on the network dynamics. A summary of the signalling overhead estimation is presented in Table IV. The different approaches are also evaluated in terms of their computation complexity. The run time for the algorithm also depends on the number of iterations and on the number of users. It can be concluded that the centralised algorithm has the highest complexity, while the DA scheme has the least complexity, with a 10.38% reduction in processing time compared with the centralised approach for the studied scenario.

An overall comparison for the studied techniques based on throughput, signalling and complexity metrics is shown in Fig. 6 for different numbers of users. It can be seen that the centralised approach has the best throughput performance, however it has higher signalling overheads and computation complexity in comparison to the other approaches. DA has the lowest throughput performance and complexity, while BSA achieves the lowest signalling overheads. BSA achieves a 49.81% reduction in signalling overheads and 0.94% reduction in complexity with less than 9% lower throughput performance compared to the centralised approach which, is a good tradeoff of throughput and signalling overheads.



Figure 6a. Use-case 1: M = 30, N = 50





Figure 6. Overall performance comparison with the centralised approach as a reference

#### TABLE IV. SIGNALLING OVERHEAD ESTIMATION

Estimation of the Signalling Overhead by the BS				
Centralised	M(1+4T)			
DA	2M(N+T)			
BS-A	2M(1+T)			

In summary, the results indicate that for throughput maximisation in a low-density network in which, selforganisation is not important, the centralised scheme is the best to adopt at the cost of signalling overheads. The DA is a promising technique to achieve good throughput performance at lower signalling overheads and complexity if device autonomy and network stability are essential. On the other hand, BSA is a semi-distributive approach which, offers a good trade-off of throughput, complexity and signalling overheads trade-off compared to DA and centralised optimisation schemes.

Regardless of the limitations of the investigated and developed RRM techniques presented in this paper, the results from adopting these methodologies, suggest the possibility of developing a conceptual qualitative evaluation framework to assist in the selection of an appropriate scheme to achieve specific priorities for the target industrial scenarios, as presented in Table V.

TABLE V. QUALITATIVE COMPARISON OF THE DIFFERENT METHODOLOGIES

Scheme	BSA	Centralised	DA
		Optimisation	
RRM	Semi	Centralised	Distributed
Approach	distributed		
RRM	Reinforcement	Mathematical	Matching
Technique	learning	optimisation	theory
Throughput	Average	Best	Worst
Complexity	Average	Worst	Best
Signalling	Best	Worst	Average
Overheads			

# V. CONCLUSIONS

We presented a semi-distributed BSA scheme for RRM of a D2D enabled cellular network targeting wireless industrial applications. The BSA scheme is an RL based approach which relies on distributed training of the D2D agents. Subsequently, the look-up tables for the D2D agents are loaded to the BS for centralised channel allocation.

The performance of the BSA scheme was compared with centralised optimisation and the game theoretic DA approaches in terms of throughput, signalling overheads and computation complexity. It is concluded that BSA offers a good trade-off of throughput, complexity and signalling overheads compared to DA and the centralised optimisation schemes. However, the BSA scheme is semi distributed. The future work aims at exploring optimised fully distributed techniques with the aim of facilitating an increased DUE autonomy through the combination of game theory and machine learning techniques.

# REFERENCES

- I. O. Sanusi and K. M. Nasr, "A Machine Learning Approach for Resource Allocation in Wireless Industrial Environments," in Proc. of the Eighteenth Advanced International Conference on Telecommunications (AICT), pp. 18-23, Jun. 2022.
- [2] J. Kaur, M. Arif Khan, M. Iftikhar, M, Imran and Q. E. Ul Haq, "Machine learning techniques for 5G and beyond networks," IEEE Access, Vol. 9, pp. 23742-23488, Jan. 2021.
- [3] F. Tariq, M. R. Khandaker, K. K. Wong, M. A. Imran, M. Bennis and M. Debbah. "A speculative study on 6G," IEEE Wireless Communications, Vol. 27, no. 4, pp. 118-125, Aug. 2020.
- [4] Y. Gu, W. Saad, M. Bennis, M. Debbah and Z. Han, "Matching theory for future wireless networks: fundamentals and applications," IEEE Communication Magazine, Vol. 53, no. 5, pp. 52-59, May 2015.
- [5] B. Tian, L. Wang, Y. Ai and A. Fei, "Reinforcement learning based matching for computation offloading in D2D communications," in Proc. of 2019 IEEE/CIC International Conference on Communications in China (ICCC), pp. 984-988, Aug. 2019.
- [6] D. L. Van and C. K. Tham, "A deep reinforcement learning based offloading scheme in ad-hoc mobile clouds," in Proc. of IEEE INFOCOM 2018-IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS), pp. 760-765, Apr. 2018.
- [7] H. Ye and Y. L. Geoffrey, "Deep reinforcement learning for resource allocation in V2V communications," in Proc. of IEEE International Conference on Communications (ICC), pp. 1-6, May 2018.
- [8] B. Fernandez-Gauna, I. Etxeberria-Agiriano and M. Graña, "Learning multirobot hose transportation and deployment by distributed round-robin Q-learning," PloS one, Vol. 10, no. 7, Jul. 2015.
- [9] K.K. Nguyen, T.Q. Duong, N.A. Vien, N.A. Le-Khac and M.N. Nguyen, "Non-cooperative energy efficient power allocation game in D2D communication: A multiagent deep reinforcement learning approach," IEEE Access, Vol. 7, pp. 100480-100490, Jul. 2019.
- [10] S. De Bast, R. Torrea-Duran, A. Chiumento, S. Pollin and H. Gacanin, "Deep reinforcement learning for dynamic network slicing in IEEE 802.11 networks," in Proc. IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS), pp. 264-269, Apr. 2019.
- [11] H. Wu, X. Li and Y. Deng, "Deep learning-driven wireless communication for edge-cloud computing: opportunities and challenges," Springer, Journal of Cloud Computing, Vol. 9, no. 21, Dec. 2020.
- [12] A. Asheralieva and Y. Miyanaga, "An autonomous learning-based algorithm for joint channel and power level selection by D2D pairs in cellular networks," IEEE transactions on communications, Vol. 64, no. 9, pp. 3996-4012, Jul. 2016.
- [13] Z. Li and C. Guo, "Multi-agent deep reinforcement learning based spectrum allocation for D2D underlay communications," IEEE Transactions on Vehicular Technology, Vol. 69, no. 2, pp. 1828-1840, Feb. 2020

- [14] S. Nie, Z. Fan, M. Zhao, X. Gu and L. Zhang, "Qlearning based power control algorithm for D2D communication," in Proc. of IEEE 27th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC), pp. 1-6, Sep. 2016.
- [15] K. Zia, N. Javed, M. N. Sial, S. Ahmed, A. A. Pirzada and F. Pervez, "A distributed multi-agent RL-based autonomous spectrum allocation scheme in D2D enabled multi-tier HetNets," IEEE Access, Vol. 15, no. 7, pp. 6733-6745, Jan. 2019.
- [16] Y. F. Huang, T. H. Tan, Y. L. Li and S.C. Huang, "Performance of resource allocation for D2D communications in Q-Learning based heterogeneous networks," in Proc. of 2018 IEEE International Conference on Consumer Electronics-Taiwan (ICCE-TW), pp. 1-5, May 2019.
- [17] I. Budhiraja, N. Kumar and S. Tyagi, "Deep-Reinforcement-Learning-Based Proportional Fair Scheduling Control Scheme for Underlay D2D Communication," IEEE Internet of Things Journal, Vol. 8, no. 5, pp. 3143-3156, Mar. 2021.
- [18] S. Nie, Z. Fan, M. Zhao, X. Gu and L. Zhang, "Qlearning based power control algorithm for D2D communication," in Proc. of IEEE 27th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC), pp. 1-5, Sep. 2016.
- [19] L. Liang, H. Ye and G.Y. Li, "Spectrum sharing in vehicular networks based on multi-agent reinforcement learning," IEEE Journal on Selected Areas in Communications, Vol. 37, no. 10, pp. 2282- 2292, Aug. 2019.
- [20] I. O. Sanusi, K. M. Nasr and K. Moessner, "Radio resource management approaches for reliable Device-to-Device (D2D) communication in wireless industrial applications," IEEE Transactions of Cognitive Communication and Networking, Vol. 7, no. 3, pp. 905-916, Oct. 2021.
- [21] A. T. Kasgari and W. Saad, "Model-free ultra-reliable low delay communication (URLLC): A deep reinforcement learning framework," in Proc. of IEEE International Conference on Communications (ICC), pp. 1-6, May 2019.
- [22] H. Wang and X. Chu, "Distance-constrained resourcesharing criteria for device-to-device communications underlaying cellular networks," Electronics letters, Vol. 48, no. 9, pp. 528-530, Apr. 2012.
- [23] H. Yang, X. Xie and M. Kadoch, "Intelligent resource management based on reinforcement learning for ultrareliable and low-delay IoV communication networks," IEEE Transactions on Vehicular Technology, Vol. 68, no. 5, pp. 4157-4169, Jan. 2019.
- [24] F. E. Souhir, A. Belghith and F. Zarai, "A reinforcement learning-based radio resource management algorithm for D2D-based V2V communication," in Proc. of 15th International Wireless Communications & Mobile Computing Conference (IWCMC), pp. 1367-1372, Jun. 24, 2019.
- [25] Y. Wei, Y. Qu, M. Zhao, L. Zhang and F.R. Yu, "Resource allocation and power control policy for Device-to-Device communication using multi-Agent

reinforcement learning," Computers, Materials & Continua, Vol. 63, no. 3, pp. 1515-1532, May 2020.

- [26] J. Kim, J. Park, J. Noh and S. Cho, "Autonomous power allocation based on distributed deep learning for deviceto-device communication underlaying cellular network," IEEE Access, Vol. 8, pp. 107853-107864, Jun. 2020.
- [27] G. Brown, "Ultra-Reliable Low-Latency 5G for Industrial Automation", Qualcomm white paper, 2018
- [28] WINNER II Channel Models, Standard IST-4-027756 WINNER II D1.1.2 V1.2, Sep. 2007.
- [29] H. Xing and S. Hakola, "The investigation of power control schemes for a device-to-device communication integrated into ofdma cellular system," in Proc. of IEEE Personal Indoor and Mobile Radio Communication (PIMRC), pp. 1775-1780, Sep. 2010.
- [30] Evolved Universal Terrestrial radio Access (E-UTRA), "Further Advancements for E-UTRA Physical Layer Aspects (Release 9)," 3GPP TR 36.814, Tech. Rep., 2010.

# Linking Radio Access Network QoE and QoS with Ensemble Multiple Regression

Adrien Schaffner LivingObjects Toulouse, France adrien.schaffner@livingobjects.com Louise Travé-Massuyès LAAS-CNRS & ANITI, University of Toulouse Toulouse, France louise.trave@cnrs.fr Simon Pachy LivingObjects Toulouse, France simon.pachy@livingobjects.com

Bertrand Le Marec LivingObjects Toulouse, France bertrand.lemarec@livingobjects.com

Abstract—The evaluation of user satisfaction is an essential performance indicator for network operators. It can be impacted by several causes, including causes linked to the network. However, linking the subjective comments of a customer with an objective behavior of the network is an issue. Experience shows that an indicator taken from customer complaints gives a good trend on the level of network quality perceived by customers, but it is difficult to transpose into concrete actions because it is often unrelated to the key performance indicators on which engineers base their action plans. The objective of this work is to learn a model that links the complaint rate, expressed by the Customer Satisfaction Rate indicator, with a set of key performance indicators so that performance engineers better understand customer expectations and act foremost on the indicators that give the most dissatisfaction. To this end, this paper takes advantage of ensemble learning applied to multiple regression, focusing the ensemble strategy on variable selection. The model hence makes it possible to link Quality of Experience and Quality of Service, which is demonstrated by nice interpretable results obtained from applying the method to data from a French telecom case study.

Index Terms—Ensemble learning; Regression models; Data analysis; Knowledge extraction; Radio access networks; QoS/QoE relationship; Quality via QoE and customer reports.

#### I. INTRODUCTION

In the space of a few years, the telecom market has undergone numerous technological and regulatory transformations that have engendered a price war from which operators are now trying to get out. They try to better differentiate themselves by moving towards a better customer experience and better support. The evaluation criteria most often adopted to establish a comparison of mobile networks are field measurement campaigns or user satisfaction surveys. User satisfaction surveys are expressed by the number of complaints received, the presence or absence of unfair terms in contracts, the commercial network and telephone assistance, connection time as well as call drop rate and their management noted by a supervisory authority, such as ARCEP (Regulatory Authority for Electronic Communications and Posts) in France or FCC (Federal Communications Commission) in United States.

The Customer Satisfaction Rate (CSR) is a good performance indicator that helps operators to effectively manage and control their business and decision making. The CSR provides the number of complaints relative to the number of customers for a given area. However, predicting customer behavior, their level of satisfaction (or dissatisfaction) has always been a challenge for operators. It is therefore important to link the CSR to a set of Key Performance Indicators (KPI) that can easily be interpreted by performance engineers to act on the relevant causes of dissatisfaction.

This paper, whose beginnings can be found in [1], presents how to learn a model that links the CSR to a set of KPIs from data while selecting a set of explanatory KPIs from an oversized, but yet relevant, set. Compared to [1], the problem is cast into an ensemble learning framework. Adopting an original point of view, model prediction and variable selection are optimized in an interlinked way by an ensemble multiple regression process. This process considers a set of base models whose results are then combined. Unlike standard approaches, ensemble integration is focused on combining the variable selection results issued from the base models rather than directly the predictions. The final regression model captures the relationship between Quality of Experience (QoE) and Quality of Service (QoS).

The contents of the paper are organized as follows. Section II analyzes related work and positions the method of this paper with respect to the state of the art. Section III formulates the problem as a regression problem and provides the identified issues. Section IV presents two regression methods, Ordinary Least Squares (OLS) and Least Absolute Shrinkage and Selection Operator (LASSO), that are later used in the three base methods for ensemble generation in Section VI. Section V describes the application that aims at explaining the customer complaint indicator CSR that has been driving the design of the method. It also presents the data that has been used and the KPIs that have been considered as candidate explanatory variables. Section VII explains the steps of the ensemble integration method. The results of applying the ensemble integration method to the CSR problem are then interpreted in Section VIII. Finally, Section IX concludes the paper.

70

# II. RELATED WORK

Much research investigated about customer complaint behavior since long [2] [3]. The idea of using complaint data to solve problems in design, marketing, installation, distribution and after sale use and maintenance, is quite natural. Understanding of customer complaint and market behavior has also been investigated so as to provide a framework for interpreting the data and extrapolating it to an entire customer base [4]. Especially in the mobile telecom industry, studies on customer complaint behaviour are numerous and continue today, significantly accentuated by the emphasis on machine learning techniques.

Given the increased competitiveness in this field, many studies have focused on a problem related to customer complaints, which is customer churn. Due to the direct effect on the revenues of the companies, especially in the telecom field, companies are seeking to develop means to predict potential customer to churn. Over the years, many machine learning algorithms have been used to produce churn prediction models and building feature's engineering and selection methods [5] [6] [7]. In the churn problem, not only complaint data but Henley segmentation, call details, line information, bill and payment information, account information, demographic profiles, service orders, etc. are potentially important. In this huge set of features, [8] identifies a subset of relevant features and applies several prediction techniques including Logistic Regressions, Multilayer Perceptron Neural Networks, Support Vector Machines and the Evolutionary Data Mining Algorithm in customer churn as predictors, based on the subset of features. [9] uses classification like the Random Forest algorithm, as well as, clustering techniques to identify the churn customers and provide the factors behind the churning of customers by categorizing the churn customers in groups.

In this paper, the focus is put on using solely complaint data to solve problems in maintenance. To do so, this work aims at linking the complaint rate with a set of technical KPIs that point at the cause of the complaints and suggest reconfiguration or repair actions on the network. This problem is much less explored in the literature than that of the churn. Literature can be exemplified by [10] that achieves correlation analysis and prediction between mobile phone users complaints and telecom equipment failures in three steps involving hierarchical clustering, pattern mining, and decision trees. On the other hand, [11] uses four machine learning algorithms, Artificial Neural Network (ANN), Support Vector Machine (SVM), K-Nearest Neighbors (KNN) and Decision Tree (DT) experimented on a database of 10,000 Korean mobile market subscribers and the variables of gender, age, device manufacturer, service quality, and complaint status. It found that ANN's prediction performance outperformed other algorithms. This last work takes into account much more data than those fixed by the objective of this paper. In addition, the first focuses on equipment failure while we want to handle the KPIs that are the data used on a daily basis by network monitoring operators. Last but not least, the algorithms used

in [11] are certainly good for prediction, but they are limited in their ability to explain predictions. The relation between the prediction and the inputs of the model remains implicit. On the contrary, the objective of this work is to clearly explain this link so that it provides useful information. This is why, the approach has been based on simple regression models while the complexity of the problem is tackled with ensemble learning.

Ensemble learning is an active research topic in different communities, including pattern recognition, machine learning, statistics and neural networks [12]. Ensemble learning [13] relies on combining several learning algorithms to obtain better predictive performance, in particular in terms of robustness and accuracy [14]. Most works on ensemble learning focus on classification problems, however this approach can as well be interesting for regression problems. It is for this latter purpose that we are concerned with it.

In this paper, we adopt the general definition of ensemble learning proposed in [15]:

*Definition 1 (Ensemble learning):* Ensemble learning is a process that uses a set of models, each of them obtained by applying a learning process to a given problem. This set of models (ensemble) is integrated in some way to obtain the final prediction.

The direct approach to ensemble learning is managed in two steps: ensemble generation that generates a set of models and ensemble integration that implements a strategy for combining the prediction results of the base models [16]. This paper adopts an original point of view in considering two tasks at once: prediction and variable selection. Unlike standard approaches, ensemble integration is focused on combining the variable selection results issued from the base models rather than directly the predictions.

#### **III. PROBLEM FORMULATION**

In our approach, the problem of explaining the level of customer satisfaction (or dissatisfaction), i.e., the QoE, is formulated as the one of obtaining a model linking the CSR to a set of KPIs that can be interpreted by performance engineers in terms of operational actions, i.e., improving QoS. To obtain this model, we rely on multiple linear regression theory and cope with the complexity of the problem through ensemble learning.

Multiple linear regression [17] is a classic family of learning algorithms that postulates that a variable is expressed as the weighted sum of other variables. Multiple linear regression defines the conditions and the model according to which a quantitative variable y is explained by several other quantitative variables  $x_j, j = 1, ..., p$ . y is considered *dependent* or *endogenous* and the variables  $x_j, j = 1, ..., p$  are said to be *explanatory* or *predictor* variables. Multiple linear regression assumes that the variation of each explanatory variable has an influence, with not necessarily equal proportions, on the behavior of the dependent variable. The function that relates the dependent variable to the explanatory variables is linear. Summarizing, multiple linear regression is a learning method that postulates that a variable y (in our problem y=CSR) is expressed as the weighted sum of other variables. In our problem, we want to learn the relationship between some KPIs and the CSR, so that performance engineers better identify the causes of customer dissatisfaction and act first and foremost on the indicators that most influence. Formally, for a number p of explanatory KPIs named  $x_j, j = 1, \ldots, p$ , which are instanciated in Section V, and the dependent variable y = CSR, the goal is to learn weights  $\beta_0, \beta_1, \ldots, \beta_p$  such as:

$$y = \beta_0 + \beta_1 x_1 + \dots \beta_p x_p \tag{1}$$

For this, we have a dataset gathering n observed samples, n > p + 1, each of dimension (p + 1) and identified by the index i:

$$(x_1^i, x_2^i, \dots, x_n^i, y^i), \ i = 1, \dots, n.$$
 (2)

Observed samples are used to estimate the parameters  $\beta_k, k = 0, \dots, p$ , that are assumed to be constant. Each sample is assumed to satisfy relation (1) with an error  $\epsilon_i$ :

$$y^{i} = \beta_{0} + \beta_{1}x_{1}^{i} + \dots \beta_{p}x_{p}^{i} + \epsilon_{i}, \ i = 1, \dots, n.$$
 (3)

Under some statistical assumptions on the error terms  $\epsilon_i$ , in particular independence and identical distribution, the vector of parameters  $\beta = (\beta_1, \dots, \beta_p)^T$  and the nuisance parameter  $\sigma^2$  defining the variance of the error  $\epsilon = (\epsilon_1, \dots, \epsilon_n)^T$ , i.e.,  $var(\epsilon) = \sigma^2 I$ , can be estimated by classical methods like least squares minimization [18] or, assuming that the error terms follow a centered normal distribution, likelihood maximization [19].

The model obtained after estimation of the parameters can be evaluated by the coefficient of determination  $R^2$ .

$$R^{2} = \frac{SSR}{SST} = \frac{\sum_{i}^{n} (\hat{y}^{i} - \bar{y})^{2}}{\sum_{i}^{n} (y^{i} - \bar{y})^{2}}$$
(4)

where  $\hat{y}^i$  is the prediction for the i-th sample,  $\bar{y}$  is the mean, SSR is the sum of squares due to regression, i.e., the variability from the mean  $\bar{y}$  that the regression manages to explain, and SST is the sum of squares total, i.e., the variability of the observed variables around the mean.

 $R^2$  represents the proportion of variance for the dependent variable that is explained by explanatory variables in the regression model. The closer the value of  $R^2$  is to 1, the better the regression. However, in practice, the threshold value for  $R^2$  for considering a good regression is highly dependent on the problem.

In our problem, the goal of the ensemble integrated regression model is to extract knowledge, i.e., to determine the KPIs that influence the CSR and to use the coefficients of the regression to quantify their influence on the CSR.

In practice, the issues to be faced are the following :

• Business experience tells us that each of the explanatory KPIs can only worsen the condition of the telecom

network and therefore should increase the CSR (e.g., an increase in the call drop rate, in the expert's mind, naturally increases the CSR). It is hence important to take care of the signs of the coefficients obtained by the regression.

• The number of candidate KPIs for explanation is high and can lead to irrelevant models.

The last issue defines one of the main objectives of this work. Indeed, there are two important elements in a model to highlight the relationship between explanatory KPIs and the dependent variable CSR:

- 1) Which are the relevant explanatory KPIs ?
- 2) How strong is their influence ?

These two elements will come as the result of the ensemble regression method that we propose in Sections VI and VII.

# IV. TWO CLASSICAL LINEAR REGRESSION METHODS

This section presents the principles of two classical multiple regression methods that are used to obtain base models as presented in Section VI. These are then leveraged in the proposed ensemble integration method presented in Section VII.

#### A. Ordinary Least Squares

When trained with data, the Ordinary Least Squares (OLS) method [20] selects parameter values  $\beta_j$ ,  $j = 1, \ldots, p$  of the linear expression (1) by the principle of *least squares*. It minimizes the sum of the squares of the differences between the observed dependent variable value in the observed data  $y^i$ ,  $i = 1, \ldots, n$ , and the value predicted by the linear function of the explanatory variables  $\hat{y}^i$ ,  $i = 1, \ldots, n$ . The optimization criterion, or loss function, is thus given by:

$$\mathcal{L} = \min_{\beta_0, \beta_1, \dots, \beta_p} \frac{1}{2} \sum_{i=1}^n (y^i - \hat{y}^i)^2$$
  
= 
$$\min_{\beta_0, \beta_1, \dots, \beta_n} \frac{1}{2} \sum_{i=1}^n (y^i - \beta_0 - \sum_{j=1}^p \beta_j x_j^i)^2$$
(5)

In geometrical terms, this can be seen as the sum of the squared distances, parallel to the axis of the dependent variable, between each data point in the set and the corresponding point on the regression surface. The smaller the differences, the better the model fits the data.

The OLS estimator is consistent, i.e., has convergence to the real parameters values as the training data is increased, when the regressors are exogenous. It is optimal in the class of linear unbiased estimators when the errors are homoscedastic, i.e., they have the same variance, and are serially uncorrelated. Under these conditions, the OLS method provides minimumvariance mean-unbiased estimation when the errors have finite variances. Under the additional assumption that the errors are normally distributed, OLS is the maximum likelihood estimator.

In this work, the function ols of the Python module statsmodels has been used to implement OLS.



Figure 1. Extract of the training data for four KPIs (in red) over one year. Units of ordinates are pourcentage for top graphs and erlangs for bottom graphs; unit of abscissa is time for all graphs.



Figure 2. Training data for the voice CSR over one year. Unit of the ordinate is a rate between 0 and 1; unit of the abscissa is time.

#### B. Least Absolute Shrinkage and Selection Operator

Least Absolute Shrinkage and Selection Operator (LASSO) is a regression method that performs both variable selection and regularization in order to enhance the prediction accuracy and interpretability of the resulting model [21]. In other words, the LASSO method handles the complexity of the model with L1 regularization [22], so that the variables not having a contribution to the model are automatically removed from the regression. This means that it adds the "absolute value of magnitude" of coefficients as penalty term to the loss function as shown in Equation 6. LASSO shrinks the less important explanatory variables altogether. This method works well for explanatory variable selection, particularly in case of a huge number of explanatory variables.

$$\mathcal{L} = \min_{\beta_0, \beta_1, \dots, \beta_p} \frac{1}{2} \sum_{i=1}^n (y^i - \hat{y}^i)^2 + \lambda \sum_{j=1}^p |\beta_j|$$
  
= 
$$\min_{\beta_0, \beta_1, \dots, \beta_n} \frac{1}{2} \sum_{i=1}^n (y^i - \beta_0 - \sum_{j=1}^p \beta_j x_j^i)^2 + \lambda \sum_{j=1}^p |\beta_j|$$
  
(6)

If  $\lambda$  is set to zero, then LASSO gets back OLS whereas a very large value increases zero coefficients hence it under-fits.

In this work, the fonction lassocv of the Python module statsmodels has been used to implement LASSO.

#### V. DATA AND PRE-PROCESSING

The goal is to predict the CSR and the influencing factors on a global scale, and not on each specific site, so that performance engineers retrieve aggregated information useful for decision making. The project was hence conducted using data at the level of French *departments* (France has 93 departments that define as many territorial communities) by setting as many regression problems as French departments.

As for the explanatory variables used, the advice of telecom experts led to a mixture of KPIs for both 2G, 3G, and 4G for six classes: traffic (like downlink data traffic), availability (like signaling failure rate), drop rates, accessibility, performance (like data\_failure rate), and mobility (like handover\_drop\_rate). In total, 50 KPIs were in the list of explanatory variables, to divide between Data and Voice. Data and Voice are indeed considered to be truly independent from a customer perspective. However, the technical KPIs used by experts to explain voice and data performance have an important common basement. Among the 35 KPIs of the voice list and the 30 KPIs of the data list, 15 KPIs were common to the two lists.



Figure 3. Steps of the fusion regression method

The available data for each department covered a full year. While both daily and weekly values were considered, it was eventually decided to stick with daily ones, to retain a bigger dataset in the training and avoid losing information by averaging over 7 days. An extract of the training data corresponding to four voice KPIs for a specific French department, 2G availability, 3G signaling failure rate, 3G voice traffic, and 4G voice traffic is shown in Figure 1 on the preceding page. The graph of the corresponding CSR is given in Figure 2 on the previous page.

In a context where the number of explanatory variables is high, it is quite often the case that several variables provide the same information or that some variables remain almost constant, or also that some variables have been poorly sensored. To remedy these common problems, classic data pre-processing solutions were applied in a first step, which consisted in:

- Removing strongly correlated variables, more precisely those with correlation coefficient higher or equal to 0.8;
- Removing variables of low variance through the dataset, more precisely those whose relative standard deviation was lower or equal to 10% of the highest;
- Removing variables with more than 10% missing values. Interpolation was used to fill the gaps for the remaining variables.

In addition, all variables were scaled so that they could be ranked according to the magnitude of their corresponding weights in the regressions.

#### VI. ENSEMBLE GENERATION

Despite the pre-processing carried out and the elimination of a subset of the KPIs proposed by experts in the field, the number of KPIs remains high, which suggests that still several of them have no direct impact on the CSR. The idea to tackle this problem is to apply an ensemble learning method leveraging the following three base regression approaches, all including a variable selection mechanisms:

- Multicollinearity analysis with OLS (M-COL),
- Backward Stepwise Regression with OLS (B-STEP),
- Structure learning with LASSO (LASSO).

Learning three base regression models with the three methods above constitutes  $Step \ 1$  of our ensemble regression method.

Each of the base methods has its own way to tackle the problem of selecting the most relevant explanatory variables, as explained in Sections VI-A, VI-B, and VI-C. To obtain the benefits of the three methods and smooth out the inconsistencies, the three methods are then integrated as explained in Section VII and illustrated in Figure 3. The originality of the proposed ensemble regression integration is that it integrates variable selection instead of directly integrating predictions. This ensemble strategy follows the analysis of [23] whose results suggest the need to examine models using multiple variable selection methods, because when they do not agree, they each may expose different aspects of the complicated theoretical relationships among predictors.

Methods M-COL and B-STEP rely on the classical Ordinary Least Squares method (OLS) presented in Section IV-A whereas LASSO, *Least Absolute Shrinkage and Selection Operator*, uses the method of the same name in its original version of linear regression as presented in Section IV-B.

#### A. Multicollinearity analysis with OLS

The M-COL method builds on OLS adding an additional preprocessing step that selects a subset of features based on multicollinearity analysis.

In a regression, multicollinearity is a problem that arises when some explanatory variables in the model measure the same phenomenon. Strong multicollinearity is problematic because it can increase the variance of the regression coefficients and make them unstable and difficult to interpret. Strongly correlated predictor coefficients will vary considerably from sample to sample. They may even present the wrong sign. Multicollinearity does not affect the goodness of the fit or the quality of the forecast. However, the individual coefficients associated with each explanatory variable cannot be interpreted reliably whereas this interpretation is exactly what we are looking for in this work.

Multicollinearity and correlation should not be confused. If collinear variables are de facto strongly correlated with each other, two correlated variables are not necessarily collinear. There is collinearity when two or more variables measure the "same thing".

Classically, in case of quantitative explanatory variables, multicollinearity can be assessed by the *variance inflation factor* (VIF) [24]. The VIF for an explanatory variable is equal to the ratio of the overall model variance to the variance of a model that includes only that single explanatory variable. This ratio is calculated for each explanatory variable. The VIF estimates how much the variance of a coefficient is "increased" due to a linear relationship with other predictors. Thus, a VIF of 1.7 tells us that the variance of this particular coefficient is 70% greater than the variance that should be observed if this factor was absolutely not correlated with the other predictors. The ideal case is obviously when all VIFs are equal to 1, indicating that there is no multicollinearity.

In the case study, multicollinearity analysis was performed considering the 35 and 30 KPIs indicated by the experts in the Voice and Data lists respectively. The VIF threshold was chosen to be 5, beyond which the corresponding KPI was eliminated. Figure 4 shows the results obtained on a specific cell.

#### B. Backward stepwise regression with OLS

After training a regression model, a *p-value* for each KPI can be obtained: it tests the null hypothesis that the coefficient is equal to zero, in other words, whatever its value, the KPI brings no information whatsoever to the model. A low p-value (typically 0.05 or less) indicates that one can reject the null hypothesis: a predictor that has a low p-value is probably a meaningful addition to the model as it changes the model prediction. Conversely, a larger p-value implies that changes in the predictor do not bring changes in the response.

Backward stepwise selection (or backward elimination) is a variable selection method that begins with a model that contains all variables under consideration (called the Full Model), then removes the least significant variable one after the other based on the p-value until a given stopping condition is satisfied. In our case, the stopping condition states that all remaining variables have a p-value smaller than some prespecified threshold.

Summarizing, the algorithm is as follows:

- train a model with all KPIs,
- remove the KPI with the highest p-value if it is not lower than a theshold,
- otherwise, stop.

Stepwise regression methods are known to have some drawbacks like instability in the variable selection and biased re-



Figure 4. KPI selection and relevancy on a scale from 0 to 1 for a specific cell: grey KPIs are those discarded by pre-processing and multicollinearity analysis, green KPIs are those of minor impact on the CSR, magenta KPIs are those that are preponderant according to the obtained regression model. KPI names have been deliberately blurred.

gression coefficients [25]. However, they may provide efficient means to examine multiple models for further investigation.

Note that the problem of biased regression coefficients can be fixed by running a model with the selected variables on a different data set.

#### C. Structure learning with LASSO

The LASSO method is well known in the literature and has already proved itself in numerous regressions. Here is a quick reminder of the presentation of Section IV-B : in the standard regression like OLS, coefficients are obtained through minimization of the residual squared sum. The LASSO method is similar but adds a penalization term to reduce the number of KPIs kept during the regression. The penalization takes the form of an L1 norm of the coefficients that reduces the available domain of values, allowing some coefficients to be precisely zero, thus letting one remove the matching KPIs.

An advantage of LASSO is that it can be used in highdimensional problems where the number of observed samples is much smaller than the number of explanatory variables, a case where more classical methods, like OLS, do not work. However, in this very case, if the true vector  $\beta$  is not hollow enough (too many variables of interest), the lasso will not be able to find all these variables of interest. Another limitation is in case of strong correlations, in particular if variables are highly correlated with each other and are important for the prediction, the lasso will favor one of them over the others. Another case, where correlations are a problem, is when the



Figure 5. Steps 1-2 of the ensemble integration method for the Voice performance problem: count of the number of times an explanatory KPI is ranked 1, 2, or 3 in the base models from M-COL (blue), B-STEP (orange), and LASSO (grey).



Figure 6. Step 3 of the ensemble integration method for the Voice performance problem: sum of the counts of the number of times an explanatory KPI is ranked 1, 2, or 3 by M-COL, B-STEP, LASSO. KPIs framed in red count above the threshold.

variables of interest are correlated with other variables. In this case, the consistency of the variable selection by LASSO is no longer guaranteed.

# VII. ENSEMBLE INTEGRATION

The principle of ensemble learning is to integrate several learning algorithms to obtain better performance. In this work, ensemble integration is not directly performed on the base model predictions, but on variable selection, for which three base algorithms have been proposed in Section VI. The integrated variable selection is used to learn the final models, hence resulting in indirect regression prediction, as illustrated in Figure 3 on page 5.

Each of the base methods has its own way to tackle the problem of selecting the most relevant explanatory variables, as explained in Sections VI-A, VI-B, and VI-C. Each also comes with a set of advantages and drawbacks.

Ensemble integration aims at obtaining the benefits of the three base algorithms and smooth out their drawbacks, in particular the fact that the base algorithms do not always select the best possible combination of variables.

In the regression model given by (1), explanatory variables  $x_j, j = 1, ..., p$ , can be ranked according to the magnitude of their corresponding weight  $\beta_1, ..., \beta_p$ . The idea developed in this work uses this ranking and includes four steps for the

whole ensemble regression method and three steps for the ensemble integration phase:

- Ensemble generation (as presented in Section VI)
  - Step 1 For every regression problem (corresponding to a French department), learn three base regression models with the three selected methods involving explanatory variable selection, namely M-COL, B-STEP, and LASSO;
- Ensemble integration
  - Step 2 For M-COL, B-STEP, and LASSO, count the number of times a given explanatory variable (KPI) has rank 1, 2, or 3 over the corresponding base regression models;
  - Step 3 Sum up the counts over the three sets of base models and select the explanatory variables whose count exceeds a threshold T;
  - Step 4 For every regression problem, learn (on different training data) two integrated regression models with OLS and LASSO considering only the explanatory variables selected at the previous step and deduce the final models and the most impacting variables.

The steps of the ensemble regression method are illustrated in Figure 3 on page 5. The output of the method takes the form of two sets of models called MODELS I and MODELS II,



Figure 7. Examples of final models: on training data (top), on test data with larger time scale (bottom).

from which knowledge about most influencing explanatory variables can be extracted as explained in Section VIII.

The ensemble integration method is exemplified with the CSR prediction problems set at the level of French departments.

*Steps 1-2* are illustrated in Figure 5 on the preceding page that gives the results for the Voice performance problem. For each explanatory KPI, the blue, orange, and grey bars provide the number of times the KPI is ranked 1, 2 or 3 in the base models obtained by the M-COL, B-STEP, and LASSO method, respectively. Let us note a good convergence of the count referring to B-STEP and LASSO.

*Step 3* is illustrated in Figure 6 on the previous page. It aggregates the counts for the base models of each method and sums them up. It hence represents the sum of the counts of the number of times an explanatory KPI is ranked 1, 2, or 3 in the base models obtained by one of the methods M-COL, B-STEP, and LASSO indifferently. A threshold is chosen, here at 45, and the explanatory KPIs that count above this threshold are selected. There are 7 KPIs that count above the threshold, framed in red.

Step 4 considers the 7 "survivor" KPIs as the most relevant for the prediction of the CSR. This is why step 4 reconsiders

every regression problem by restricting explanatory variables to these 7 KPIs. OLS and LASSO methods are run with these explanatory variables alone on another set of training data. Figure 7 shows some examples of the obtained final models on training and test data.

#### VIII. MAKING SENSE OF THE PREDICTIONS

Let us recall that the objective of this work is to design a model that makes it possible to link the CSR indicator with a set of objective performance indicators so that performance engineers better understand customer expectations and act first and foremost on the indicators that give the most dissatisfaction. The results of the prediction problems can be analyzed in two ways: at the level of each French department, and aggregated for the whole France.

# A. Interpretation at the level of each French department

An interpretation at the level of each French department is done by associating a profile to each department. For this purpose, the results of the final B-STEP models (B-STEP method applied to the 7 survivor KPIs) have been used and the department profiles have been obtained by clustering the coefficients of the obtained models. The clustering was carried



Figure 8. Map of French departments colored by profiles given by the weight of top KPIs influencing the CSR. Departements are identified by their name, number and main city in italics. Overseas departments appear in gray and framed and are not included in the analysis.

out using the classical K-means algorithm that consists in iteratively grouping the individuals (here the models) that are the most similar until stability is reached. Thus, 5 groups emerge whose coefficients associated with each KPI are similar. This leads to the map in Figure 8 where the departments that have similar profile are depicted with the same color. A similar profile indicates that the KPIs that must be mainly incriminated are the same, and so are the reasons explaining customer complaints.

#### B. Aggregated interpretation

The aggregated interpretation is at the level of the whole France. It requires an additional analysis based on the final models obtained at step 4 of the ensemble integration method. To complete this analysis, KPIs ranked 1, 2, and 3 over OLS and LASSO final models and all the French departments are determined. These KPIs are shown in Figure 9 on the following page, where the three top KPIs (among the 7 survivors) appear in red, namely: 3G voice traffic, 2G



Figure 9. KPIs ranked 1, 2, and 3 over OLS and LASSO final models and over all the French departments.

availability, 3G voice drop rate. These are the most significant KPIs to explain customer dissatisfaction and they indicate that complaints are highly related to network behavior, which is intuitively understandable.

Among the various metrics used to measure network behavior:

- 3G voice traffic reports about the amount of traffic,
- 2G availability indicates loss of network coverage,
- 3G voice drop rate indicates the rate of call drops.

The KPI 3G voice traffic comes to the first rank. The amount of traffic represented by 3G voice traffic can be related to network unavailability and network engineering issues. It is easy to understand why these problems may be the main cause of the dissatisfaction of customers.

The KPI 2G availability comes to the second rank. Loss of network coverage represented by 2G availability can be associated to network maintenance processes. The fact that this strongly impacts customer dissatisfaction makes sense.

The KPI 3G voice drop rate comes to the third rank, which is not surprising either.

Let us notice that other metrics like accessibility failure rate or mobility issues appear to be less significant than call drops or traffic issues.

To improve client experience, the network operators should therefore prioritize to base their action plans on:

- reducing unavailability periods by, for instance, optimizing the maintenance process,
- improving the call drop rate by modifying network parameter settings, optimizing site engineering, or building new sites.

# IX. CONCLUSION AND PERSPECTIVES

This paper proposes an ensemble learning method to obtain a regression model with explanatory power. In many applications, the number of variables that could be thought to be explanatory for a given dependent variable is huge. However, many of them are correlated or collinear and others do not really impact the predicted variable. The method presented in this paper leverages the benefits of three methods to select relevant explanatory variables and deduce a robust regression model. The originality of the ensemble regression integration phase is to focus the integration on variable selection instead of directly on the prediction of the base models.

The method has been tested on telecom data to obtain a model that indicates the impact of a set of objective performance indicators on the customer complaint rate so that performance engineers better understand customer expectations and act first and foremost on the indicators that give the most dissatisfaction. The final results can be used to cluster French departments according to their profile as a function of the top influencing KPIs. Similar profiles indicate that the reasons to be incriminated to explain customer complaints are close, and so are the actions that should be taken. The final results can also be used on a global scale to exhibit the top KPIs at country level and the high level management strategy to be applied.

Future work will consider mapping the top KPIs returned by the model to actual actions to be performed on the network so that customer satisfaction is increased, i.e., CSR is decreased. This mapping could benefit from ideas coming from the combination of the theories of prospect theory and satisfaction games found in the literature, such as [26].

# ACKNOWLEDGMENT

We would like to thank the telecom expert Didier Rana for his remarks and relevant suggestions throughout this work.

#### REFERENCES

- A. Schaffner, L. Travé-Massuyès, S. Pachy, and B. Le Marec, "Explaining radio access network user dissatisfaction with multiple regression models," in 15th International Conference on Communication Theory, Reliability, and Quality of Service (CTRQ 2022), Barcelona, Spain. IARIA, April 2022, pp. 11–17.
- [2] J. Jacoby and J. J. Jaccard, "The sources, meaning, and validity of consumer complaint behavior: A psychological analysis." *Journal of Retailing*, vol. 57, no. 3, pp. 4–24, 1981.
- [3] J. Singh and R. E. Wilkes, "When consumers complain: A path analysis of the key antecedents of consumer complaint response estimates," *Journal of the Academy of Marketing Science*, vol. 24, no. 4, pp. 350– 365, 1996.
- [4] J. Goodman and S. Newman, "Understand customer behavior and complaints," *Quality Progress*, vol. 36, no. 1, pp. 51–55, 2003.

- [5] W.-H. Au, K. C. Chan, and X. Yao, "A novel evolutionary data mining algorithm with applications to churn prediction," *IEEE Transactions on Evolutionary Computation*, vol. 7, no. 6, pp. 532–545, 2003.
- [6] A. K. Ahmad, A. Jafar, and K. Aljoumaa, "Customer churn prediction in telecom using machine learning in big data platform," *Journal of Big Data*, vol. 6, no. 1, pp. 1–24, 2019.
- [7] A. Amin, F. Al-Obeidat, B. Shah, A. Adnan, J. Loo, and S. Anwar, "Customer churn prediction in telecommunication industry using data certainty," *Journal of Business Research*, vol. 94, pp. 290–301, 2019.
- [8] B. Huang, M. T. Kechadi, and B. Buckley, "Customer churn prediction in telecommunications," *Expert Systems with Applications*, vol. 39, no. 1, pp. 1414–1425, 2012.
- [9] I. Ullah, B. Raza, A. K. Malik, M. Imran, S. U. Islam, and S. W. Kim, "A churn prediction model using random forest: analysis of machine learning techniques for churn prediction and factor identification in telecom sector," *IEEE Access*, vol. 7, pp. 60134–60149, 2019.
- [10] Q. Yang, G. Ji, and W. Zhou, "The correlation analysis and prediction between mobile phone users complaints and telecom equipment failures under big data environments," in 2nd International Conference on Advanced Robotics and Mechatronics (ICARM 2017). IEEE, 2017, pp. 201–206.
- [11] C. Choi, "Predicting customer complaints in mobile telecom industry using machine learning algorithms," Ph.D. dissertation, Purdue University, USA, 2018.
- [12] R. Polikar, "Ensemble based systems in decision making," *IEEE Circuits and Systems Magazine*, vol. 6, no. 3, pp. 21–45, 2006.
- [13] O. Sagi and L. Rokach, "Ensemble learning: A survey," Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, vol. 8, no. 4, p. e1249, 2018.
- [14] D. Opitz and R. Maclin, "Popular ensemble methods: An empirical study," *Journal of Artificial Intelligence Research*, vol. 11, pp. 169–198, 1999.
- [15] J. Mendes-Moreira, C. Soares, A. M. Jorge, and J. F. D. Sousa, "Ensemble approaches for regression: A survey," ACM Computing Surveys (CSUR), vol. 45, no. 1, pp. 1–40, 2012.
- [16] N. Rooney, D. Patterson, S. Anand, and A. Tsymbal, "Dynamic integration of regression models," in *International Workshop on Multiple Classifier systems (MCS 2004), Cagliari, Italy.* Springer, June 2004, pp. 164–173.
- [17] D. C. Montgomery, E. A. Peck, and G. G. Vining, *Introduction to linear regression analysis*. John Wiley & Sons, 2021.
- [18] G. S. Watson, "Linear least squares regression," *The Annals of Mathe-matical Statistics*, pp. 1679–1699, 1967.
- [19] I. J. Myung, "Tutorial on maximum likelihood estimation," Journal of Mathematical Psychology, vol. 47, no. 1, pp. 90–100, 2003.
- [20] B. Craven and S. M. Islam, "Ordinary least-squares regression," *The SAGE Dictionary of Quantitative Management Research*, pp. 224–228, 2011.
- [21] R. Tibshirani, "Regression shrinkage and selection via the lasso: a retrospective," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 73, no. 3, pp. 273–282, 2011.
- [22] —, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.
- [23] P. Ruengvirayudh and G. P. Brooks, "Comparing stepwise regression models to the best-subsets models, or, the art of stepwise," *General Linear Model Journal*, vol. 42, no. 1, pp. 1–14, 2016.
- [24] T. A. Craney and J. G. Surles, "Model-dependent variance inflation factor cutoff values," *Quality Engineering*, vol. 14, no. 3, pp. 391–403, 2002.
- [25] G. Smith, "Step away from stepwise," *Journal of Big Data*, vol. 5, no. 1, pp. 1–12, 2018.
- [26] S. Papavassiliou, E. E. Tsiropoulou, P. Promponas, and P. Vamvakas, "A paradigm shift toward satisfaction, realism and efficiency in wireless networks resource sharing," *IEEE Network*, vol. 35, no. 1, pp. 348–355, 2020.