

In the Name of God

# Journal of Information Systems & Telecommunication

Vol. 7, No. 3, July-September 2019, Serial Number 27

Research Institute for Information and Communication Technology  
Iranian Association of Information and Communication Technology  
Affiliated to: Academic Center for Education, Culture and Research (ACECR)

**Manager-in-Charge:** Habibollah Asghari, ACECR, Iran

**Editor-in-Chief:** Masoud Shafiee, Amir Kabir University of Technology, Iran

#### Editorial Board

Dr. Abdolali Abdipour, Professor, Amirkabir University of Technology, Iran

Dr. Mahmoud Naghibzadeh, Professor, Ferdowsi University, Iran

Dr. Zabih Ghasemlooy, Professor, Northumbria University, UK

Dr. Mahmoud Moghavvemi, Professor, University of Malaya (UM), Malaysia

Dr. Ali Akbar Jalali, Professor, Iran University of Science and Technology, Iran

Dr. Alireza Montazemi, Professor, McMaster University, Canada

Dr. Ramezan Ali Sadeghzadeh, Professor, Khajeh Nasireddin Toosi University of Technology, Iran

Dr. Hamid Reza Sadegh Mohammadi, Associate Professor, ACECR, Iran

Dr. Sha'ban Elahi, Associate Professor, Tarbiat Modares University, Iran

Dr. Shohreh Kasaei, Professor, Sharif University of Technology, Iran

Dr. Mehrnoush Shamsfard, Associate Professor, Shahid Beheshti University, Iran

Dr. Ali Mohammad-Djafari, Associate Professor, Le Centre National de la Recherche Scientifique (CNRS), France

Dr. Saeed Ghazi Maghrebi, Assistant Professor, ACECR, Iran

Dr. Rahim Saeidi, Assistant Professor, Aalto University, Finland

**Executive Manager:** Shirin Gilaki

**Executive Assistants:** Ali Mokhtarani, Mahdokht Ghahari

**Print ISSN:** 2322-1437

**Online ISSN:** 2345-2773

**Publication License:** 91/13216

**Editorial Office Address:** No.5, Saeedi Alley, Kalej Intersection., Enghelab Ave., Tehran, Iran,

P.O.Box: 13145-799

Tel: (+9821) 88930150 Fax: (+9821) 88930157

E-mail: info@jist.ir , infojist@gmail.com

URL: www.jist.ir

#### Indexed by:

- |   |                         |
|---|-------------------------|
| - SCOPUS  | www.Scopus.com          |
| - Index Copernicus International                                  | www.indexcopernicus.com |
| - Islamic World Science Citation Center (ISC)                     | www.isc.gov.ir          |
| - Directory of open Access Journals                               | www.Doaj.org            |
| - Scientific Information Database (SID)                           | www.sid.ir              |
| - Regional Information Center for Science and Technology (RICEST) | www.ricest.ac.ir        |
| - Iranian Magazines Databases                                     | www.magiran.com         |

#### Publisher:

Regional Information Center for Science and Technology (RICEST)

Islamic World Science Citation Center (ISC)

This Journal is published under scientific support of  
Advanced Information Systems (AIS) Research Group and  
Digital & Signal Processing Research Group, ICTRC

## Acknowledgement

JIST Editorial-Board would like to gratefully appreciate the following distinguished referees for spending their valuable time and expertise in reviewing the manuscripts and their constructive suggestions, which had a great impact on the enhancement of this issue of the JIST Journal.

### (A-Z)

- Amintoosi, Haleh, Ferdowsi University, Mashhad, Iran
- Behkamal, Behshid, Ferdowsi University, Mashhad, Iran
- Abouei, Jamshid, Yazd University, Yazd, Iran
- Alavi, Seyed Enayatollah, Chamran University, Ahvaz, Iran
- Kaebbeh, Yaeghoobi, ManavRachna International University, India
- Mir, Ali, Lorestan University, Iran
- Cheraghchi, Hamideh Sadat, Shahid Beheshti University, Tehran, Iran
- Dehkharghani, Rahim, University of Bonab, East Azerbaijan, Iran
- Asgari, Mohammad Javad, University of Torbat Heydarieh, Khorasan Razavi, Iran
- Feshari, Majid, Kharazmi University, Tehran, Iran
- Ghasemzadeh, Ardalan, Urmia University of Technology, West Azerbaijan, Iran
- Mavadati, Samira, Mazandaran University, Iran
- Nikanjam, Amin, University of Science & Technology, Iran
- Zeynali, Esmail, Qazvin Islamic Azad University, Qazvin, Iran
- Fakhar, Fatemeh, PNU University, Ahvaz, Iran
- Mirroshandel, Seyed Abolghasem, University of Guilan, Rasht, Iran
- Mousavirad, Seyed Jaleleddin, University of Kurdistan, Kurdistan, Iran
- Omid Mahdi, Ebadati, Kharazmi University, Tehran, Iran
- Javan, Mohammadreza, Shahrood University of Technology, Semnan, Iran
- Mirzaei, Abbas, Islamic Azad University, Ardabil, Iran
- Mohammadzadeh, Sajjad, University of Birjand, South Khorasan, Iran
- Pirgazi, Jamshid, Zanjan University, Zanjan, Iran
- Reshadat, Vahideh, Malek-Ashtar University of Technology, Tehran, Iran
- Rasi, Habib, Shiraz University of Technology, Shiraz, Iran
- Jampoor, Mehdi, Quchan University of Technology, Razavi Khorasan, Iran
- Alizadeh, Sasan, Qazvin Islamic Azad University, Qazvin, Iran
- Derhami, Vali, Yazd University, Yazd, Iran
- Teymoori, Mehdi, Zanjan University, Iran
- Shirvanimoghadam, Shahriar, Shahid Rajae Teacher Training University, Iran
- Tahernia, Amirhossein, Ferdowsi University, Mashhad, Iran

## Table of Contents

• <b>A New Non-Gaussian Performance Evaluation Method in Uncompensated Coherent Optical Transmission Systems</b> .....	165
Seyed Sadra Kashef and Paeiz Azmi	
• <b>BSFS: A Bidirectional Search Algorithm for Flow Scheduling in Cloud Data Centers</b> .....	175
Hasibeh Naseri, Sadoon Azizi and Alireza Abdollahpouri	
• <b>Balancing Agility and Stability of Wireless Link Quality Estimators</b> .....	184
Mohamad Javad Tanakian and Mehri Mehrjoo	
• <b>SSIM-Based Fuzzy Video Rate Controller for Variable Bit Rate Applications of Scalable HEVC</b> .	193
Farhad Raufmehr and Mehdi Rezaei	
• <b>DeepSumm: A Novel Deep Learning-Based Multi-Lingual Multi-Documents Summarization System</b> .....	204
Shima Mehrabi, Seyed Abolghasem Mirroshandel and Hamidreza Ahmadifar	
• <b>Social Groups Detection in Crowd by Using Automatic Fuzzy Clustering with PSO</b> .....	215
Ali Akbari, Hassan Farsi and sajad Mohamadzadeh	
• <b>Facial Images Quality Assessment based on ISO/ICAO Standard Compliance Estimation by HMAX Model</b> .....	225
Azamossadat Nourbakhsh , Mohammad Shahram Moin and Arash Sharifi	

# A New Non-Gaussian Performance Evaluation Method in Uncompensated Coherent Optical Transmission Systems

Seyed Sadra Kashef \*

Faculty of Electrical and Computer Engineering, Urmia University, Iran  
s.kashef@urmia.ac.ir

Paeiz Azmi

Faculty of Electrical and Computer Engineering, Tarbiat Modares University, Iran  
pazmi@modares.ac.ir

Received: 04/Jul/2019

Revised: 24/Oct/2019

Accepted: 20/Dec/2019

## Abstract

In this paper, the statistical distribution of the received quadrature amplitude modulation (QAM) signal components is analyzed after propagation in a dispersion uncompensated coherent optical fiber link. Two Gaussian tests, the Anderson-Darling and the Jarque-Bera have been used to measure the distance from the Gaussian distribution. By increasing the launch power, the received signal distribution starts to deviate from Gaussian. This deviation can have significant effects in system performance evaluation. The use of the Johnson  $s_U$  distribution is proposed for the performance evaluation of orthogonal frequency division multiplexing in an uncompensated coherent optical system. Here, the Johnson  $s_U$  is extended to predict the performance of multi-subcarrier and also single carrier systems with M-QAM signals. In particular, symbol error rate is derived based on the Johnson  $s_U$  distribution and performance estimations are verified through accurate Monte-Carlo simulations based on the split-step Fourier method. In addition, a new formulation for the calculation of signal to noise ratio is presented, which is more accurate than those proposed in the literature. In the linear region, the Johnson based estimations are the same as Gaussian; however, in the nonlinear region, Johnson  $s_U$  distribution power prediction is more accurate than the one obtained using the Gaussian approximation, which is verified by the numerical results.

**Keywords:** Coherent optical fiber link; Gaussian distribution; Johnson  $s_U$  distribution; nonlinear transmission performance; Uncompensated Transmission; QPSK.

## 1- Introduction

Modeling of nonlinear propagation in coherent optical communication systems is of fundamental to predict system performance. In particular, [1], [2], [3] present a practical survey on modeling of nonlinear propagation in uncompensated transmission (UT) systems. Although, the propagation of light in optical fiber channels is properly modeled by the non-linear Schrodinger equation, it is difficult to attain an accurate statistical model of nonlinear fiber channel [4] because of the non-Gaussian behavior of noise [5].

In this context, the Gaussian-Noise model (GN-model) is known as an accurate and acceptable reference model that is applied in different system scenarios for coherent optical communication. This model has been successfully used in system design, analysis and network optimization. However, in some scenarios such as strong nonlinear propagation and low dispersion the accuracy of the GN-model is reduced. Therefore, various models with higher accuracy have been suggested in many different scenarios [6], [7], [8].

In long-haul propagation, amplified spontaneous emission (ASE) noise and nonlinear interference (NLI) caused by the Kerr effect are pointed out as the two main system impairments [5]. As demonstrated in [9], the use of an inaccurate signal statistic can give more than 500 km error in reach prediction of a fiber-optic transmission system. This statistical deviation of received signal histogram from Gaussian distribution is measured using two powerful tests named Kolmogorov-Smirnov and Anderson-Darling (AD) [10], [11] and is shown in different scenarios.

To overcome the inaccuracy of Gaussian distribution, some enhanced methods have been introduced in [12] to improve the accuracy of the GN model. In some cases, a correction factor has been applied to achieving more accuracy in different systems, like coherent optical orthogonal frequency-division multiplexing (CO-OFDM) [10], [13], [14].

The much of the literature is focused on calculating the moments of propagated signal based on a simple Gaussian assumption, which requires the estimation of just one moment (the variance). However, the non-Gaussian distribution of the signal after highly nonlinear propagation is not compliant with this assumption and the

\* Corresponding Author

estimation of higher moments of the received signal statistics is required, it as shown in [11]. The Non-Gaussian behavior of propagated signal is included in an enhanced version of the GN model, named EGN, that takes into account for high order moments of the signal [15]. The Gaussian assumption of the received signal starts to be inaccurate at the nonlinear threshold, as shown in [11]. The typical operating parts of optical systems fall around this threshold. Therefore, an accurate model in this region can play an important role.

The use of Johnson  $s_U$  distribution for BER calculation was proposed for the first time in [11], [16] in CO-OFDM system with different modulations. The accuracy of the Johnson  $s_U$  based methods were verified using both analytical and numerical results. The study presented in this article seeks to further extend the use of the Johnson  $s_U$  distribution statistic for the performance evaluation of other coherent optical UT systems with different kinds of modulation and signals ([11] and [16] were only focused on OFDM signals) which is presented in [17], briefly. Single carrier M-QAM signals are the first aim; next, multi-subcarrier (MSC) QPSK transmission systems are considered because of their higher nonlinear robustness [18/], [19/], [20/], [21/]. We also analyzed the performance of proposed method in dual polarization systems.

Therefore, we investigate the Gaussian assumption accuracy of the propagated signal components in different modulation and signals to extend Johnson  $s_U$  distribution applications. The deviation of nonlinear noise pdf from Gaussian is measured using two well-known pdf tests, namely, Jarque-Bera (JB) and AD. Johnson  $s_U$  pdf includes higher-order statistics to achieve a more accurate prediction. We report the obtained SER through direct error counting in the Monte-Carlo (MC) simulation based on the split-step Fourier method (SSFM), and the difference from (a) the GN-model based SER estimation and (b) Johnson  $s_U$  based method.

The rest of the paper is organized as follow. Section 2 describes the system model which section 3 introduces the Mathematical Preliminaries of system performance evaluation through UT over fiber optic. The statistical features of propagated signal is evaluated and Johnson  $s_U$  pdf is suggested for more accuracy. Finally, Johnson  $s_U$  is applied for BER evaluation in Section 4. Conclusions are presented in Section 5.

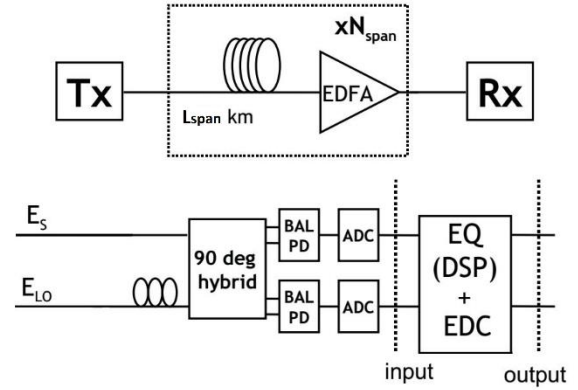


Fig. 1 Top: link structure. Bottom: coherent receiver block diagram.

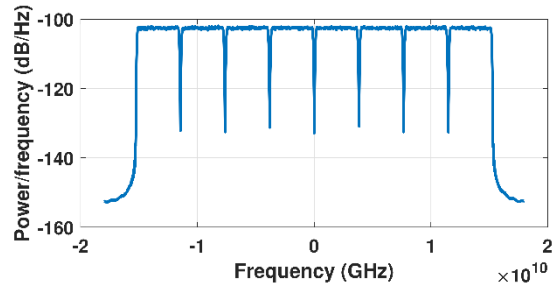


Fig. 2 The spectrum of output signal in a MSC Tx with 32 GHz bandwidth and 250 MHz channel spacing.

## 2- System Model

The diagram of the investigated coherent system is illustrated in the top part of Fig. 1. The optical 32 Gbaud M-QAM signals are generated in electrical domain and then modulated onto an optical carrier at a desired wavelength using an optical modulator Mach-Zehnder modulator (RF/optical converter). Without any loss of generality, calculations are accomplished in baseband. 32-GHz Nyquist-shaped frequency spectrum is divided into  $N_{sc}$  electrical subcarriers for the case of MSC electrical multiplexing. In addition, square-root raised-cosine spectra with  $N_{sc} = 8$  is used as shown in Fig. 2. Independent pseudo-random binary sequences (PRBSs) are used in all MSC channels.

The generated optical signal is fed to the optical channel, consisting of  $N_{span}$  fiber spans with  $L_{span}$  (km) length and each span followed by an erbium-doped fiber amplifier (EDFA), which completely recovers the span loss. On the other hand, EDFA accumulates ASE noise in each span. Therefore, the overall length of the channel is  $L = N_{span}L_{span}$ . Two typical fibers, i.e. non-zero dispersion shifted fiber (NZDSF) and standard single mode fiber (SMF), are used. The fiber parameters (i.e. the

dispersion  $\beta_2$ , the attenuation  $\alpha$ , and the nonlinearity  $\gamma$ ) are reported in Table 1.

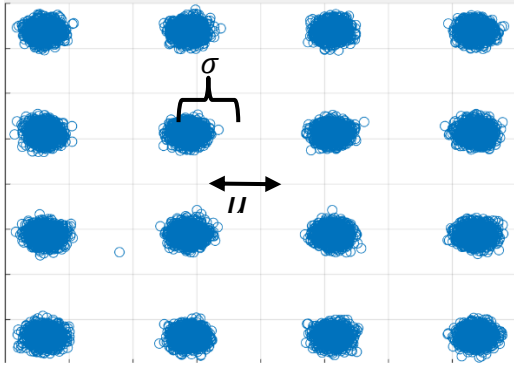


Fig. 3 Typical measured 16-QAM constellation after the dispersion and phase noise compensation in the receiver

The Kerr effect is the origin of nonlinear effects, which are classified as self-channel interference (SPM-like), cross-channel interference (XPM-like) and multi-channel interference (FWM-like). The predominant effect in uncompensated optical systems is FWM or multi-channel interference, which generates new signals through the nonlinear combination of the propagating signals at different frequencies. This nonlinear interference can be modeled as additive noise on the constellation symbols, as shown in Fig. 3 [8].

Table 1: Parameters of simulated systems

Simulated link Parameters	(A)	(B)
Optical fiber	SMF	NZDSF
$\gamma$ (Wkm) <sup>-1</sup>	1.3	2
$\beta_2$ ps <sup>2</sup> /km	-21	-3.38
$\alpha$ dB/km	0.2	0.2288
Signal bandwidth ( $BW$ ) GHz	32	32
Span length ( $L_s$ ) km	100	80
Noise Figure ( $F$ ) dB	5	5
Light frequency ( $\nu$ ) THz	193.1	193.1
PRBS	$2^{18} <$	$2^{18} <$

The ASE noise is the main source of linear noise, which is added by the EDFAs and can be modeled as an additive stationary Gaussian noise with variance [5]:

$$\sigma_{ASE}^2 = N_{span} e^{\alpha L} h \nu F B W \quad (1)$$

where  $F, \nu, h, B W$  are the noise figure of the optical amplifier, the absolute light frequency, the Planck constant, and the reference bandwidth, respectively.

In the bottom part of Fig. 1, a standard coherent receiver is shown, which is used in this paper. The local oscillator (LO) is mixed with the signal in a 90-degree hybrid. In this paper we neglect the effects of LO alignment and linewidth, which is a practical assumption in coherent

system as an advantage of digital signal processing (DSP). Signal components at the output of the two balanced photodetectors are sampled at enough rate for DSP. By using data-aided DSP algorithms linear propagation effects such as CD are completely compensated. Moreover, phase and amplitudes of in-phase and quadrature of received signals are accessible. In dual polarization systems a parallel system is needed at transmitter and receiver, with a different polarization which is completely independent of the other polarization.

### 3- Mathematical Preliminaries

This section presents the principle of signal propagation through the UT optical fiber, together with the relevant basic mathematical equations. Later, the statistical model of the received signal is derived. A new modified signal-to-noise ratio (SNR) is derived based on Johnson  $S_U$  distribution. Then, the deviation of the propagated signal pdf from the Gaussian distribution is measured using the AD and JB Gaussianity tests.

#### 3-1- Gaussianity Tests

For BER evaluation, focusing on the pdf tail is essential and is a better tool for checking pdf. At first glance, the received signal histogram fails accurate Gaussianity tests, particularly when looking at the far tails of the distribution, especially in systems with low BERs. Gaussianity tests are powerful tools for evaluating the extent of deviation of the random sample histogram from Gaussian distribution. The Gaussianity of the symbols pdf is tested by computing the statistics of the JB and AD tests that are described in the following:

AD Test: This test measures the difference between two population sets. The most encouraging point with the AD test is that it focuses upon the difference between the tails of two distributions. If ordered data,  $Y_1, \dots, Y_n$ , come from a distribution with cumulative distribution function  $\Phi$ , the formula for the AD test statistic is,  $A^2 = -n - S$  where [22/]

$$S = \sum_{i=1}^n \frac{2i-1}{n} [\ln(\Phi(Y_i)) + \ln(1 - \Phi(Y_{n+1-i}))] \quad (2)$$

For measuring the distance with Gaussian distribution,  $\Phi(x)$  equals  $1 - Q(x)$

where  $Q(x) = \frac{1}{\sqrt{2\pi}} \int_x^\infty \exp(-\frac{u^2}{2}) du$ .

JB Test: JB test is the second test that is utilized to measure the difference between the histogram of received samples and Gaussianity. Here, the main tools are skewness and kurtosis of the received samples. The test statistic is defined as [23/]:

$$JB = \frac{n}{6} \left( S^2 + \frac{(K - 3)^2}{4} \right), \tag{3}$$

where  $K$  and  $S$  denote kurtosis and skewness of the observed samples, respectively.

It should be mentioned that  $JB$  and  $A^2$  are two metrics and, if the pdf of received samples are Gaussian, the result of both tests will be zero. For non-Gaussian pdf,  $JB$  and  $A^2$  are not zero, and higher values of these metrics mean more distance from Gaussian.

Note that the AD test focuses on the difference between the tails of distributions and the JB test measures Gaussianity using higher order moments.

The system described in Fig. 1 is simulated numerically using the SSFM as a benchmark according to the simulation parameters are reported in Table 1. Linear and nonlinear effects are considered in propagation and then at the receiver linear propagation effects such as CD and constant constellation rotation are compensated. The accuracy of the fit is measured for all symbols of constellation by gathering all symbols to center by subtracting the estimated mean value of each symbols. In each test, the results of real and imaginary parts are summed and plotted for different launch powers.

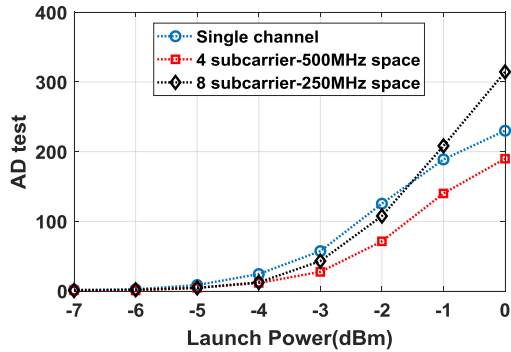


Fig. 4 AD test results for system (B) with  $60 \times 80$  (km) NZDSF link and 4-QAM signals.

The AD and JB tests results are shown in Figs. 4 and 5 for a  $60 \times 80$  (km) NZDSF link (system (B) in Table 1) with 4-QAM signal. Nonlinear threshold is the start point of nonlinear effect where test values diverge from zero (Gaussian pdf). This divergence means that the pdf of received samples changes where nonlinearity begins to affect. The JB and AD test results are again done for system (A) of Table 1 with  $65 \times 100$  (km) SMF link and the results are shown in Figs. 6 and 7. The tails (AD test) and high order statistics (JB test) of the distribution of the received signal are different from Gaussian pdf in nonlinear region and this may significantly affect SER calculations.

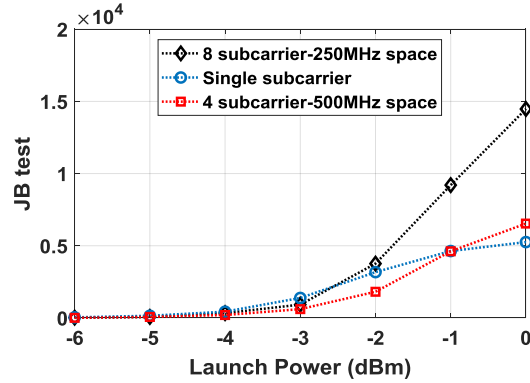


Fig. 5 JB test results for system (B) with  $60 \times 80$  (km) NZDSF link and with 4-QAM signals.

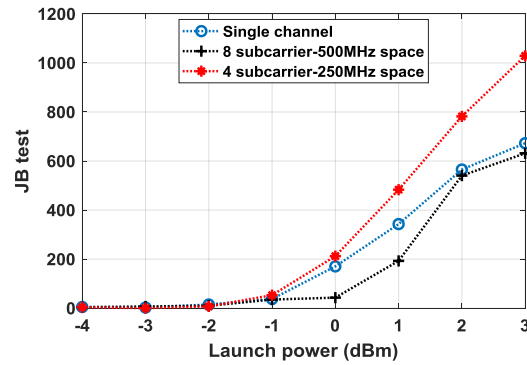


Fig. 6 JB test results for system (A) with  $65 \times 100$  (km) SMF link and with 4-QAM signals.

JB and AD tests are used in the systems with 16-QAM and 64-QAM signals. These two modulations are transmitted over in SMF and NZDSF systems with different span numbers. The results, shown in Figs. 8 and 9, demonstrate that, in a system with NZDSF, received signal distribution deviates from Gaussianity (zero value) at lower powers than SMF, which confirms higher nonlinearity of NZDSF ( $0.002 \text{ (Wkm)}^{-1}$ ) in comparison with SMF ( $0.0013 \text{ (Wkm)}^{-1}$ ). Moreover, over SMF, the results of JB and AD tests show that Gaussianity at the end of span number 8 is more than at span number 6. This can be the effect of high CD in SMF. However, in NZDSF Gaussianity of received signal after 6 spans propagation is more than 8 spans; because nonlinear effect is dominant respect to CD in NZDSF.

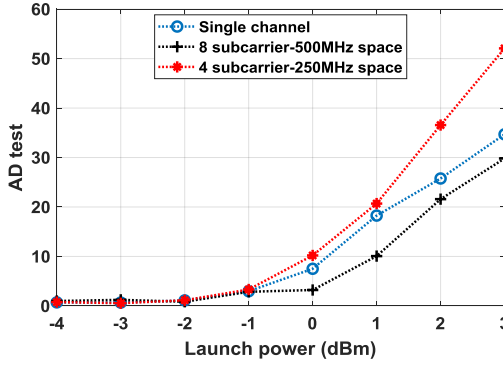


Fig. 7 AD test results for system (A) with  $65 \times 100$  (km) SMF link and with 4-QAM signals.

In the next section, Johnson  $s_U$  distribution is used for SER evaluation, and the obtained results are compared to those achieved using other methods.

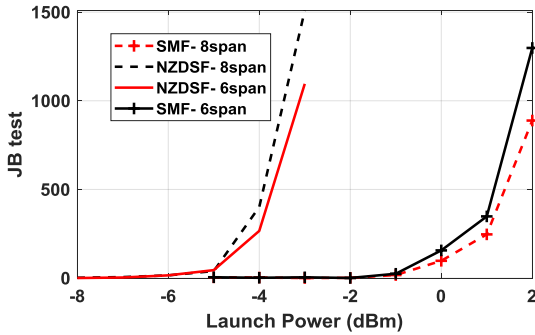


Fig. 8 JB test results for systems (A) and (B) with different span number and dual polarization 16-QAM signals.

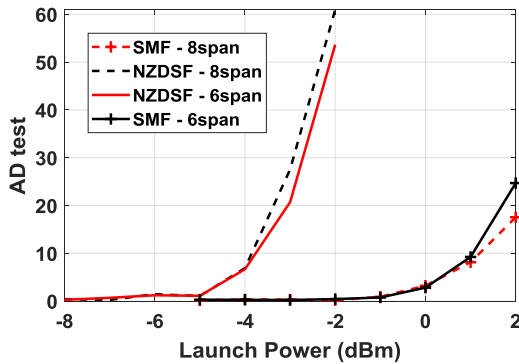


Fig. 9 AD test results for systems (A) and (B) with different span number and dual polarization 16-QAM signals.

### 3-2- New SNR

The BER of any coherent communication system with 4-QAM signals in an AWGN channel is as follows [26/]:

$$BER \approx \frac{1}{2} \operatorname{erfc} \left( \sqrt{\frac{SNR}{2}} \right) \quad (4)$$

and in general for rectangular M-QAM we have

$$SER \approx 2 \left( 1 - \frac{1}{\sqrt{M}} \right) \operatorname{erfc} \left( \sqrt{\frac{SNR}{\frac{2}{3}(M-1)}} \right) \quad (5)$$

It is apparent that using Gray coding:

$$SER = \log_2(M) \times BER \quad (6)$$

SNR in a multi-span optical fiber link has an NLI component which is added to ASE noise and SNR can be written as [24/]:

$$SNR = \frac{P}{\sigma_{ASE}^2 + \sigma_{NLI}^2} \quad (7)$$

where  $P$  is signal power,  $\sigma_{ASE}^2$  and  $\sigma_{NLI}^2$  are ASE noise and nonlinear noise variances, respectively.

Eq. (4) is achieved assuming that the propagated signal in nonlinear optical fiber has a Gaussian pdf. However, it is demonstrated in Figs. 4 - 9 that received signal histogram is not Gaussian after nonlinear threshold and nonlinear region. According to JB test, 3<sup>rd</sup> and 4<sup>th</sup> moments of propagated signal are not zero and we can use them to increase the accuracy of performance estimation, which was proposed also in [25/].

In this way, we proposed in [11] and [16] to use the Johnson  $s_U$  as a four-parameter distribution with two more degrees of freedom with respect to Gaussian pdf for distribution fitting. Johnson  $s_U$  pdf is a transformation of the standard normal pdf [26/] by applying 4<sup>th</sup> and 3<sup>rd</sup> moments of noise statistic. Because of symmetry of FWM effect, 3<sup>rd</sup> moment is set to zero, as mentioned in [11], [16] and we just use the 4<sup>th</sup> moment as an additional degree of freedom to improve the accuracy of system performance prediction. According to Johnson  $s_U$  pdf a close form relationship between BER and SNR can be written as [11]:

$$BER \approx \frac{1}{2} \operatorname{erfc} \left( \delta \sinh^{-1} \left( \frac{\zeta}{\sqrt{2}\lambda} \right) \right) \quad (8)$$

where  $\zeta$ ,  $\lambda$ , and  $\delta$  are parameters of Johnson  $s_U$  distribution that are equal to:

$$\begin{aligned} \zeta &= \hat{\mu}_1 \\ \delta &= \ln(\hat{\omega})^{-\frac{1}{2}} \\ \lambda &= \sqrt{\frac{2\hat{\mu}_2}{\hat{\omega}^2 - 1}} \end{aligned} \quad (9)$$

where,

$$\hat{\omega} = \left[ (2\hat{K} - 2)^{\frac{1}{2}} - 1 \right]^{\frac{1}{2}} \quad (10)$$

and

$$\hat{K} = \frac{\hat{\mu}_4}{\hat{\mu}_2^2} \quad (11)$$



$\hat{\mu}_1$ ,  $\hat{\mu}_2$ , and  $\hat{\mu}_4$  are first, second and fourth central moment estimations, respectively [27/]. The hat sign means numerical approximation of parameters. By comparing Eq.(4) and Eq.(6), the following relationships between the SNR and the parameters of the Johnson  $s_U$  distribution are obtained:

$$\sqrt{\frac{SNR}{2}} = \delta \sinh^{-1}\left(\frac{\zeta}{\sqrt{2}\lambda}\right) \quad (12)$$

Therefore,  $SNR/2$  is equivalent to  $h^{-1}\left(\frac{\zeta}{\sqrt{2}\lambda}\right)$ , which can thus be interpreted as a new modified version of SNR. This modified SNR version depends on  $\delta$  as the representative of 4<sup>th</sup> moment in addition to the  $\frac{\zeta}{\lambda}$ , which is the representative of signal and noise powers.

The new SNR can be extended to higher order modulations M-QAM with rectangular constellation like 16-QAM and 64-QAM as follows:

$$\sqrt{\frac{SNR}{\frac{2}{3}(M-1)}} = \delta \sinh^{-1}\left(\frac{\zeta}{\sqrt{\frac{2}{3}(M-1)}\lambda}\right) \quad (13)$$

and the symbol error rate (SER) can be written as:

$$SER \approx 2\left(1 - \frac{1}{\sqrt{M}}\right) \operatorname{erfc}\left(\delta \sinh^{-1}\left(\frac{\zeta}{\sqrt{\frac{2}{3}(M-1)}\lambda}\right)\right) \quad (14)$$

The method of moments is applied to calculate the Johnson  $s_U$  distribution parameters from 2<sup>nd</sup> and 4<sup>th</sup> moments of received signal, numerically.

#### 4- Results and Discussion

Here, three methods are used for estimating the SER as an important parameter of system performance:

- SSFM: as a benchmark to measure the accuracy of the two other methods.

- Gaussian numerical: uses the estimated variance  $\sigma_y^2$ ,  $\sigma_x^2$  and mean  $\mu_y$ ,  $\mu_x$  of the constellation points (see Fig. 3). We assume that in-phase and quadrature parts of received signal are i.i.d [6]. If  $\sigma_y^2$ ,  $\sigma_x^2$  and mean  $\mu_y$ ,  $\mu_x$  belongs to the corners of the constellation for rectangular M-QAM. Based on the Gaussian approximation, Eq. (5) can be used for SER calculation.

- Johnson numerical: pdf parameters  $(\zeta, \lambda, \delta)$  of received sampled signal are needed to calculate SER semi-analytically using the Johnson  $s_U$  distribution. Estimating  $\mu_1$  and  $\mu_2$  is straightforward, similar to the second method, and fourth central moments  $\mu_4$ , is estimated numerically as follows:

$$\hat{\mu}_4 = \frac{1}{N} \sum_{i=1}^N (y_i - \mu_1)^4. \quad (15)$$

where  $y_i$ s are the received samples.

Johnson  $s_U$  distribution parameters can be calculated using Eqs. (9) - (11) according to estimated  $\mu_1$ ,  $\mu_2$ , and  $\mu_4$ . In this method, SER is calculated using Eq. (14) for rectangular distributions.

It should be mentioned that there is an assumption in Eq. (14) and (5) that is the symmetry of distribution which is correct for both Johnson  $s_U$  and Gaussian. Therefore, SER calculation is done just for in-phase or quadrature part.

SER versus launch power estimated using the three methods described above shown in Fig. 10. A 40×80 (km) NZDSF system was simulated, using the parameters reported in Table 1, column (B). According to Fig. 10, in nonlinear region, Johnson  $s_U$  distribution will achieve about 1 dB better power prediction at  $SER = 10^{-3}$  than Gaussian. The simulation of Fig. 10 is repeated for the system (A) with 60×100 (km) SMF spans. This means that there is 60 numbers of 100 km SMF spans. The results shown in Fig. 11 confirm the higher accuracy of proposed method over SMF, as well.

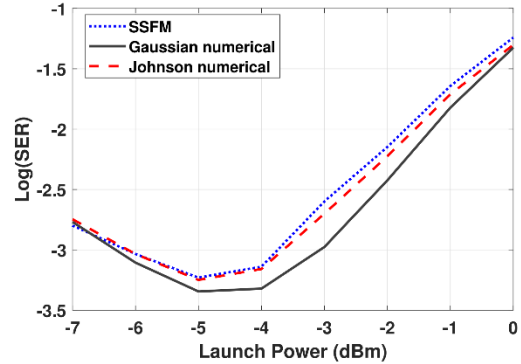


Fig. 10 Performance comparison among three methods in 40×80 (km) NZDSF UT link with parameters of System (B) according to Table 1 and 4-QAM signals.

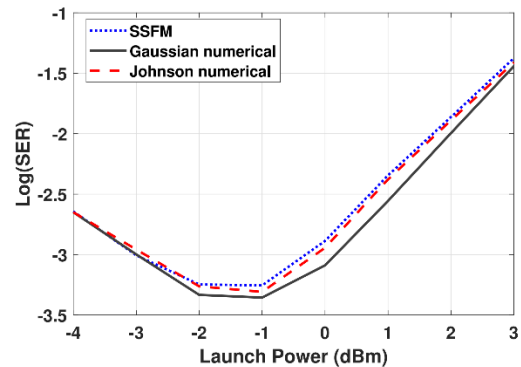


Fig. 11 Performance comparison among three methods in 60×100 (km) SMF UT link with parameters of System (A) according to Table 1 and 4-QAM signals.

SMF has a higher CD and a lower nonlinearity in comparison with NZDSF which results in more Gaussian

received samples and therefore Gaussian based method in Fig. 11 is closer to MC than in Fig. 10.

The accuracy of the proposed method was also assessed in a WDM scenario with 50 GHz spacing. The results over NZDSF (see Table 1, column B) are shown in Fig. 12 where the higher accuracy with respect to the Gaussian based method is confirmed.

In Figs. 13, 14, 15 and 16 Johnson  $s_U$  distribution is used for performance evaluation in MSC systems with different subcarrier number and frequency spacing. The performance improvement is clear in this four figures when the spectrum is divided into 8 and 4 parts with 4 GBaud and 8 GBaud symbol-rate each, respectively. It means that by dividing the spectrum into 4 or 8 parts, nonlinear effects fall down and system can reach more distances with desired performance.

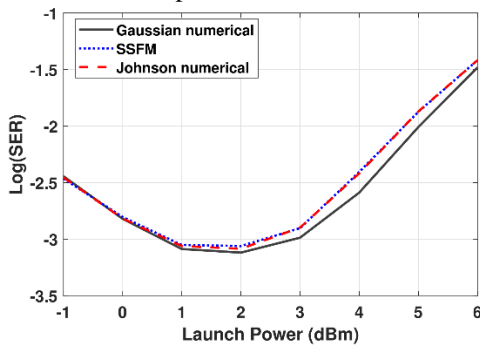


Fig. 12 Performance comparison among three methods in 40x80 (km) NZDSF UT 5-channel WDM link with parameters of System (B) according to Table 1 for each channel and 4-QAM signals.

By comparing Figs.13 and 14, it can be concluded that 8-subcarrier spectrum have more robustness than 4-subcarrier spectrum against nonlinearity. It is also apparent that the proposed Johnson  $s_U$  based method is more accurate than the Gaussian method in all analyzed configurations and SER prediction error is reduced to less than 0.1 order of magnitude in SMF links. According to Figs. 13-16, system performance can improve of about one and more than one order of magnitude in SMF and NZDSF links, respectively.

The performance evaluation of MSC is again done for system (B) scenario, as is shown in Figs. 15 and 16. In MSC system with 8 subcarriers and 500 MHz spacing of Figs. 15 and 16, there is more than 0.5 order of magnitude SER deviation at -2 dBm launch power and more than 1.5 dB power prediction error using Gaussian based method at  $SER = 10^{-3.5}$ . The frequency spacing in Fig.16 decreased to 250 MHz, which can increase nonlinear effect. It is also shown that the proposed Johnson  $s_U$  based method is more accurate than the Gaussian methods in NZDSF based links.

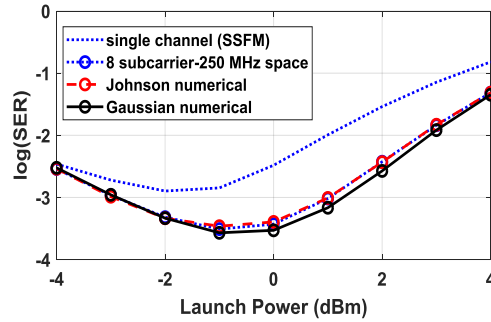


Fig. 13 Performance comparison among three methods in 65x100 (km) SMF UT link with parameters of System (A) according to Table1 and 4-QAM signals. Blue dotted curves belong to numerical SSFM.

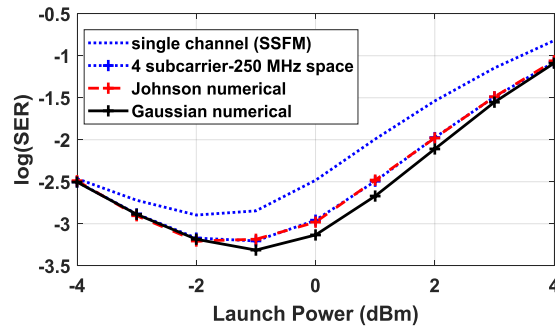


Fig. 14 Performance comparison among three methods in 65x100 (km) SMF UT link with parameters of System (A) according to Table1 and 4-QAM signals. Blue dotted curves belong to numerical SSFM.

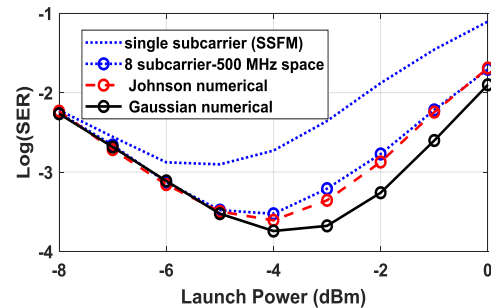


Fig. 15 Performance comparison among three methods in 45x80 (km) NZDSF UT link with parameters of System (B) according to Table1 and 4-QAM signals. Blue dotted curves belong to numerical SSFM.

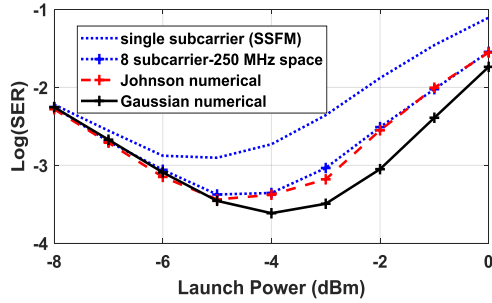


Fig. 16 Performance comparison among three methods in 45×80 (km) NZDSF UT link with parameters of System (B) according to Table1 and 4-QAM signals. Blue dotted curves belong to numerical SSFM.

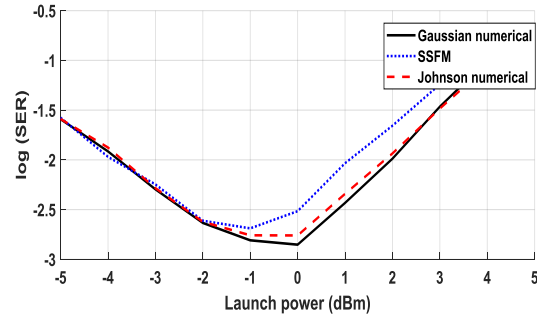


Fig. 18 Performance comparison among three methods in 3×100 (km) SMF UT link with parameters of System (A) according to Table1 with dual polarization 64-QAM signals.

SER versus launch power estimated using the three methods, described above, is shown in Figs. 17 and 18. A 10×100 (km) and a 3×100 SMF based link is simulated, using the parameters reported in Table 1, column (A). It can be seen from Fig. 17 and 18 that, in nonlinear region, Johnson  $s_{\mathcal{U}}$  distribution works better than Gaussian, and the performance fits to the SSFM results more accurately. The nonlinear effects of high order modulation cause an equalization error that corrupts the moment estimation and that causes a gap between SSFM and two other methods. For tackling with this problem, analytical methods can help, which will be done in future works.

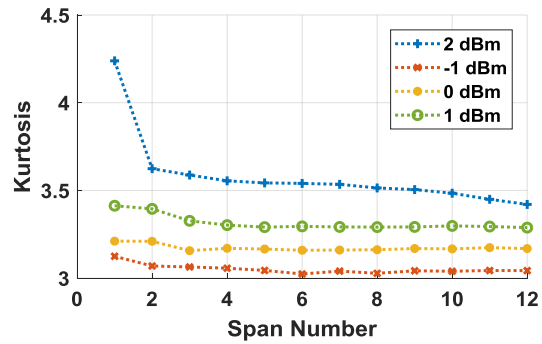


Fig. 19 Kurtosis comparison in 10×100 (km) SMF UT link with parameters of System (A) according to Table1 with dual polarization 16-QAM signals. Kurtosis for a Gaussian data set is 3.

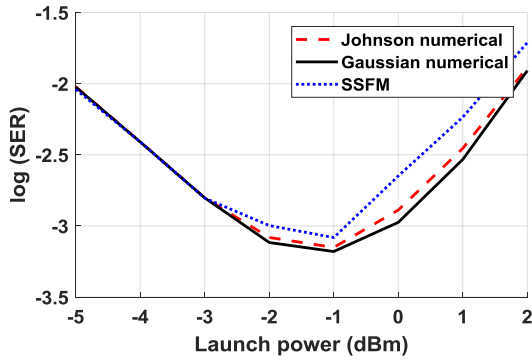


Fig. 17 Performance comparison among three methods in 10×100 (km) SMF UT link with parameters of System (A) according to Table1 with dual polarization 16-QAM signals.

In addition, CD is another parameter that affects Gaussianity of signal i.e. the signal is dispersed by increasing the span number. By increasing the span number, received signal kurtosis starts to converge to 3 (kurtosis of a Gaussian data set), which is shown in Fig.19 for 4 different powers. At high powers, kurtosis value is higher due to nonlinear effect.

### 5- Conclusion

The statistics of the propagated signal in optical fiber transmission system can be affected by nonlinear effects which have a critical role in system performance calculation. Nonlinear effects deviate propagated signal probability distribution from the Gaussian distribution (which is a typical assumption in modeling of optical systems) when the launch power increases. In this paper, two JB and AD tests have been used to measure the deviation of propagated signal at different powers, modulations, and span numbers. As a result of the mentioned tests, propagated signal distribution starts to deviate from Gaussian after the nonlinear threshold. Therefore, the performance prediction methods based on the Gaussian assumption are not accurate in the nonlinear region. This paper extended the use of the Johnson  $s_{\mathcal{U}}$  distribution for performance evaluation of coherent optical systems with different kinds of signals; single carrier M-QAM signals, multi-subcarrier QPSK transmission systems are considered because of their higher nonlinear robustness in comparison with single-carrier systems. We also analyzed the performance of proposed method in dual

polarization systems. Monte-Carlo simulation results were used for verification of the proposed semi-analytical approach, which is more accurate in different scenarios in

### Acknowledgments

We would like to thank Prof. Gabriella Bosco from Technical University of Turin for her time, technical supports and valuable comments.

### References

- [1] P. Poggiolini and Y. Jiang, "Recent advances in the modeling of the impact of nonlinear fiber propagation effects on uncompensated coherent transmission systems," *Journal of Lightwave Technology*, Vol. 35, No 3, 2017, pp. 458-480.
- [2] P. Poggiolini, G. Bosco, A. Carena, V. Curri, Y. Jiang, F. Forghieri, "The GN model of fiber non-linear propagation and its applications," *Journal of Lightwave Technology*, Vol. 32, No. 4, 2014, pp. 694-721.
- [3] P. Poggiolini, Y. Jiang, A. Carena, F. Forghieri, Analytical Modeling of the Impact of fiber Non-Linear Propagation on Coherent Systems and Networks, in *Enabling Technologies for High Spectral-efficiency Coherent Optical Communication Networks*, Xiang Zhou, Chongjin Xie editors, chapter 7, pp. 247-310, ISBN: 978-1-118-71476-8, Wiley, Hoboken (New Jersey), 2016.
- [4] R.-J. Essiambre, G. Kramer, P. J. Winzer, G. J. Foschini, and B. Goebel, "Capacity limits of optical fiber networks," *Journal of Lightwave Technology*, Vol. 28, No 4, 2010, pp. 662-701.
- [5] G. P. Agrawal, *Nonlinear fiber optics*, Academic Press, 2007.
- [6] L. Beygi, E. Agrell, P. Johannisson, M. Karlsson, and H. Wymeersch, "A discrete-time model for uncompensated single-channel fiber-optical links," *IEEE Transaction on Communication*, Vol. 60, No 11, 2012, pp. 3440-3450.
- [7] R. Dar, M. Feder, A. Mecozzi, and M. Shtaif, "Properties of nonlinear noise in long, dispersion-uncompensated fiber links," *Optics Express*, Vol. 21, No 22, 2013, pp. 25685-25699.
- [8] A. Carena, V. Curri, G. Bosco, P. Poggiolini, and F. Forghieri, "Modeling of the impact of nonlinear propagation effects in uncompensated optical coherent transmission links," *Journal of Lightwave Technology*, Vol. 30, No 10, 2012, pp.1524-1539.
- [9] P. Serena, A. Bononi, and N. Rossi, "The impact of the modulation dependent nonlinear interference missed by the gaussian noise model," *ECOC 2014*.
- [10] S. T. Le, K. J. Blow, V. K. Mezentsev, and S. K. Turitsyn, "Bit error rate estimation methods for QPSK CO-OFDM transmission," *Journal of Lightwave Technology*, Vol. 32, No 17, 2014, pp. 2951-2959.
- [11] S. S. Kashef, P. Azmi, G. Bosco, M.D. Matinfar, and D. Pileri, "NonGaussian Statistics of CO-OFDM Signals after Non-Linear Optical Fiber Transmission," *IET optoelectronics*, Vol. 12, No 3, 2017, pp. 150 – 155.
- [12] A. Carena, G. Bosco, V. Curri, Y. Jiang, P. Poggiolini, and F. Forghieri, "On the accuracy of the GN-model and on analytical correction terms to improve it," *arXiv preprint arXiv:1401.6946*, 2014.
- [13] G. Gao, X. Chen, and W. Shieh, "Analytical expressions for nonlinear transmission performance of coherent optical OFDM systems with frequency guard band," *IEEE Photonics Technology Letters*, Vol. 30, No 15, 2012, pp. 2447-2454.
- [14] D. Uzunidis, C. Matrakidis, and A. Stavdas, "An improved model for estimating the impact of FWM in coherent optical systems," *Optics Communications*, Vol. 378, 2016, pp. 22-27.
- [15] A. Carena, G. Bosco, V. Curri, Y. Jiang, P. Poggiolini, and F. Forghieri, "EGN model of non-linear fiber propagation," *Optics Express*, Vol. 22, No 13, 2014, pp. 16335-16362.
- [16] S. S. Kashef and P. Azmi, "Performance Analysis of Nonlinear Fiber Optic in CO-OFDM Systems with High Order Modulations," *IEEE Photonics Technology Letters*, Vol. 30, No 8, 2018, pp. 696-699.
- [17] S. S. Kashef, G. Bosco, and P. Azmi, "Johnson S U Distribution in Uncompensated QPSK Coherent Optical Transmission Systems." *ICEE*, 2019, pp. 1284-1288.
- [18] D. Uzunidis, C. Matrakidis, and A. Stavdas, "Analytical FWM expressions for coherent optical transmission systems," *Journal of Lightwave Technology*, Vol. 35, No13, 2017, pp. 2734-2740.
- [19] D. Uzunidis, C. Matrakidis, and A. Stavdas, "Closed-form FWM expressions accounting for the impact of modulation format," *Optics Communications*, Vol. 440, 2019, pp.132-138.
- [20] F. P. Guiomar, A. Carena, G. Bosco, L. Bertignono, A. Nespola, and P. Poggiolini, "Nonlinear mitigation on subcarrier-multiplexed PM-16QAM optical systems," *Optics Express*, Vol. 25, No 4, 2017, pp. 4298-4311.
- [21] F. Buchali, W. Idler, K. Schuh, L. Schmalen, T. Eriksson, G. Bcherer, P. Schulte, and F. Steiner, "Study of electrical subband multiplexing at 54 GHz modulation bandwidth for 16QAM and probabilistically shaped 64QAM," *ECOC*, 2016, pp. 4951.
- [22] R. D'Agostino, *Goodness-of-fit-techniques*, Routledge, 2017.

- [23] C.M. Jarque, and A. K. Bera, "A test for normality of observations and regression residuals," *International Statistical Review/Revue Internationale de Statistique*, 1987, pp. 163-172.
- [24] A. Carena, V. Curri, G. Bosco, P. Poggiolini, F. Forghieri, "Modeling of the Impact of Non-Linear Propagation Effects in Uncompensated Optical Coherent Transmission Links," *Journal of Lightwave Technology*, Vol. 30, No 10, 2012, pp. 1524-1539.
- [25] P. Jenneve, P. Ramantanis, J.C. Antona, G. de Valicourt, M. Mestre, H. Mardoyan, and S. Bigo, "Pitfalls of error estimation from measured nongaussian nonlinear noise statistics over dispersion-unmanaged systems," *ECOC 2014*.
- [26] W. P. Elderton and N. L. Johnson, *Systems of frequency curves*, Cambridge University Press London, 1969.
- [27] I. D. Hill, R. Hill, R. L. Holder, "Algorithm AS 99: Fitting Johnson curves by moments," *J. the royal statistical society. Series C (Applied statistics)*, vol. 25, no. 2, 1976, pp 180-189.

**Seyed Sadra Kashef** received the B.S. degree in Telecommunication Engineering from Tabriz University, Tabriz, Iran in 2012, M.S. degree Ph.D. in Telecommunication Engineering from Tarbiat Modares University, Tehran, Iran, in 2014 and 2018, respectively. Currently he is assistant professor in Urmia University, Iran. Her research interests include Optical communications, machine learning, communication theory, estimation theory.

**Paeiz Azmi** was born in Tehran, Iran, in April 1974. He received the B.Sc., M.Sc., and Ph.D. degrees in electrical engineering from the Sharif University of Technology (SUT), Tehran, in 1996, 1998, and 2002, respectively. From 1999 to 2001, he was with the Advanced Communication Science Research Laboratory, Iran Telecommunication Research Center (ITRC), Tehran, where he was with the Signal Processing Research Group, from 2002 to 2005. Since September 2002, he has been with the Electrical and Computer Engineering Department, Tarbiat Modares University, Tehran, where he became an Associate Professor, in January 2006, and is currently a Full Professor. His current research interests include modulation and coding techniques, digital signal processing, wireless communications, radio resource allocation, molecular communications, and estimation and detection theories.

# BSFS: A Bidirectional Search Algorithm for Flow Scheduling in Cloud Data Centers

Hasibeh Naseri

Department of Computer Engineering and IT, University of Kurdistan, Sanandaj, Iran  
hasibehnaseri@gmail.com

Sadoon Azizi\*

Department of Computer Engineering and IT, University of Kurdistan, Sanandaj, Iran  
s.azizi@uok.ac.ir

Alireza Abdollahpouri

Department of Computer Engineering and IT, University of Kurdistan, Sanandaj, Iran  
abdollahpouri@uok.ac.ir

Received: 04/Jul/2019

Revised: 24/Aug/2019

Accepted: 08/Dec/2019

## Abstract

To support high bisection bandwidth for communication intensive applications in the cloud computing environment, data center networks usually offer a wide variety of paths. However, optimal utilization of this facility has always been a critical challenge in a data center design. Flow-based mechanisms usually suffer from collision between elephant flows; while, packet-based mechanisms encounter packet re-ordering phenomenon. Both of these challenges lead to severe performance degradation in a data center network. To address these problems, in this paper, we propose an efficient mechanism for the flow scheduling problem in cloud data center networks. The proposed mechanism, on one hand, makes decisions per flow, thus preventing the necessity for rearrangement of packets. On the other hand, thanks to SDN technology and utilizing bidirectional search algorithm, our proposed method is able to distribute elephant flows across the entire network smoothly and with a high speed. Simulation results confirm the outperformance of our proposed method with the comparison of state-of-the-art algorithms under different traffic patterns. In particular, compared to the second-best result, the proposed mechanism provides about 20% higher throughput for random traffic pattern. In addition, with regard to flow completion time, the percentage of improvement is 12% for random traffic pattern.

**Keywords:** Cloud Computing; Data Center Networks; Flow Scheduling; Routing Algorithm; Load Balancing; Bidirectional Search.

## 1- Introduction

Over the past few years, several companies and organizations have shifted their services such as large scale computing, web search, online gaming, and social networking to cloud computing environment [1]. Recently, with the emergence of IoT-based applications and massive data processing, the demand for cloud resources has increased dramatically. In order to meet these needs, various data center networks are deployed around the world, including hundreds of servers and large amounts of traffic are exchanged between them.

Today's data center networks often use multi-rooted tree topologies such as Fat-tree [2-4] and Clos [5, 6]. These topologies provide multiple paths at an equal cost between each pair of end hosts, and thus significantly increase

bisection bandwidth. However, given the burstiness and unpredictable nature of the traffic matrix and the flow pattern generated by virtual machines on hosts, achieving load balancing in a data center network is not a trivial task.

Over the past few years, network researchers and traffic engineers have proposed various algorithms and mechanisms to provide load balancing in cloud data center networks [7-17]. Although these efforts are valuable steps towards improving the efficiency of data center networks, there exist still some challenges and issues in this regard. The mechanisms that use *per-packet* approach to manage network traffic, although provide good load balancing across the network, but they are faced with the phenomenon of packet re-ordering. Packet re-ordering not only affects TCP throughput but also imposes significant computational overhead on hosts [12]. On the other hand, *flow-based* mechanisms usually suffer from the

\* Corresponding Author

phenomenon of collision between the elephant flows, which leads to network performance reduction. Therefore, the issue of load balancing in data center networks is still challenging and needs further research efforts [1].

In this paper, we aim to design an efficient flow-based mechanism to achieve load balancing in data center networks. Given the fact that most flows in data centers are only a few kilobytes in size (i.e., mice flows) and a very small percentage of them are large-sized flows (i.e., elephant flows), we take advantage of a hybrid mechanism for flow scheduling in a data center network. For this purpose, we use the distributed ECMP algorithm for mice flows, while a central controller is used for elephant flows. When an elephant flow is detected by a host, it sends the flow to the controller for routing the first packet. In the controller, based on the defined cost matrix, an optimal bidirectional search is performed on the network to find and select the best route for that flow.

Our proposed mechanism has three major advantages. First, it prevents the packet re-ordering phenomenon; because it performs per flow. Second, since the controller is used only for elephant flows, it does not become a bottleneck. Third, because the central controller provides a macroscopic view of the network traffic, our algorithm is able to distribute the elephant flows smoothly across the network. We have compared our approach with various mechanisms such as Static [2], ECMP [18] and DiFS [19]. The results of the experiments clearly show the superiority of our algorithm in terms of *delay*, *throughput* and *flow completion time* in comparison with other mentioned approaches.

The rest of the paper is organized as follows. In Section 2, related works are reviewed. Background and problem definition are described in Section 3. The proposed mechanism is presented in Section 4. In Section 5, we describe the simulation and evaluate the performance of the proposed method. Finally, Section 6 concludes the paper.

## 2- Related Works

In general, flow scheduling algorithms are divided into two main categories [1]: distributed and centralized. On the other hand, in terms of how the flows are handled, they can be classified into three categories [1]: packet-based, flow-based and flowlet-based. Below, we review some of the most important works performed on flow scheduling in cloud data center networks.

ECMP [18] is the most common routing algorithm for flow scheduling in data center networks. It is a distributed and flow-based algorithm. When a flow enters a switch for the first time, the ECMP performs the routing operation by

applying the hash function to the header of the packet. Although the implementation of this algorithm is very simple, it does not differentiate between mice and elephant flows; and therefore, collisions between elephant flows is inevitable.

DARD [4] is another flow-based distributed algorithm. In this algorithm, end-hosts are responsible for monitoring the status of network traffic. Based on the network feedback received from the probe packets, each host moves the flows from high-traffic routes to low-traffic routes. However, injecting a large number of probe packets into the network puts considerable overhead on it. In addition, since this algorithm is host-based, all hosts need to be upgraded, which imposes a lot of administrative costs.

Cui et al. [14, 19] have recently proposed an adaptive distributed mechanism, called DiFS, for flow scheduling in data center networks with Fat-tree topology. They use ECMP to forward mice flows, while for scheduling the elephant flows each switch greedily distributes them to the output ports. In order to prevent over-utilized links, DiFS may change the path of some flows based on the collaboration between switches. Simulation results show that DiFS performs far better than ECMP. However, since this algorithm has no macroscopic view of the network, it has to transmit a large number of messages between the switches to provide load balancing. This, on the one hand, leads to overhead on the network, and on the other hand, as some flows change their direction, packets may need to be re-ordered.

DRILL [20] is another distributed mechanism which is inspired by the idea of "the power of two random choices". In each switch and for each packet, this approach decides which output port to send the packet to, based on local information about the queue length. The port selection mechanism in DRILL is simple and easy to implement. However, due to the fact that DRILL operates on a per-packet basis, packet out-of-ordering is inevitable.

Some other works [2, 21, 22] use Static or deterministic routing to forward packets over the network. In Static routing, the path between each host pair is determined in advance and remains unchanged. Although in practice the implementation of such mechanism is very simple, it cannot well take advantage of the multi-path benefit provided by the topology of data center networks, and usually provides very low performance.

Hedera [7] is a dynamic flow scheduling system in which mice flows are separated from elephant flows using a specified threshold value. By default, Hedera uses ECMP to transmit flows on the network. However, when a

large flow is detected, the system generates a demand matrix for active flows. Therefore, it proposes two schedulers “Global First Fit” and “Simulated Annealing” to send the flows. The results show that Hedera achieves a significant performance improvement compared to ECMP for the moderate cost. On the one hand, the demand estimation matrix in Hedera is only run once per scheduling period, which takes about 200 milliseconds for a data center network with 27,648 hosts and 250,000 large flows. On the other hand, the execution time of the scheduler is in order of tens of milliseconds. Putting these two points together, it can be seen that it takes several hundred milliseconds to find a suitable approximate path for guiding large flows in a data center network, which is a considerable time. In addition, given the drastic changes in traffic patterns in data centers, Hedera has to build demand matrixes over and over again in a short period of time, which imposes a significant overhead on the system.

Wang et al. in [23] proposed an adaptive mechanism called Freeway for flow scheduling in data center networks. This mechanism partitions the paths between hosts into *low latency* and *high throughput* paths. It then transmits mice flows through low latency paths and elephant flows through high throughput paths. Although Freeway performs better than ECMP and Hedera, in practice it may leave much of the network’s capacity unused [1].

Based on the ant colony optimization algorithm, the authors of [16] proposed a centralized scheduling mechanism for transmitting the flows in data center networks. Their algorithm divides the elephant flows into  $k$  segments and sends them through  $k$  edge-disjoint paths. Since the flows are broken in this algorithm, the problem of re-ordering packets arises.

Authors in [24] have modeled the elephant flow scheduling problem as a multi-knapsack problem and proposed a mechanism based on hybrid Genetic and Simulated Annealing algorithm to solve it. Simulation results confirm that their algorithm provides higher bisection bandwidth and lower latency in comparison with similar methods. However, since their approach is similar to Hedera, it has same drawbacks.

### 3- Background

The topologies proposed for data center networks over the past few years provide multiple paths between each pair of hosts [2, 5, 21, 25]. Although our proposed mechanism can be applicable to any topology, in this paper, we focus on the well-known and common Fat-tree topology [2]. In

this section, we first briefly describe the Fat-tree topology. Then, we will focus on the characteristics of flows in data center networks. We then investigate the mechanisms for detection of mice flows and elephant flows. Finally, we illustrate the problem of collision of elephant flows in data center networks with a detailed example.

#### 3-1- Fat-tree topology

Fat-tree is a hierarchical multi-root tree topology that contains three layers of switches called Top of Rack (ToR), aggregation, and core. In this topology, the switches are homogeneous and the degree of each switch is determined by the parameter  $n$ . Fat-tree consists of  $n$  pods, each pod having two layers and each layer has  $n/2$  switches that form a complete bipartite graph. Figure 1 shows an example of a Fat-tree topology with 4-port switches ( $n = 4$ ).

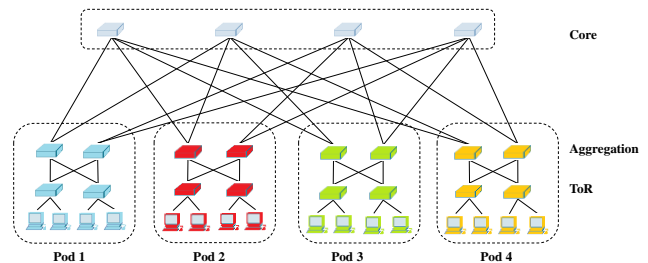


Fig. 1 Fat-tree topology with 4-port switches

In this work, we define the Fat-tree topology as a directed graph  $G = (V, E)$  where,  $V$  represents the switches and  $E$  represents the links. Also, links that connect lower layer switches to higher layer switches are called *uphill* and links that connect higher layer switches to the lower layer switches are called *downhill* links.

#### 3-2- Flow properties in data center networks

Each flow contains several packets that are chained together. In data center networks, if a flow contains a lot of packets, or it takes a long period of time, or its traffic is more than a threshold value, it is known as an elephant flow. On the other hand, flows with low information volumes or low number of packets are called mice flows [26]. In terms of number of flows in a data center network, typically more than 90% of them are mice, while only less than 10% of them are elephant flows. Nevertheless, on the other hand, more than 90% of the data volume belongs to the elephant flows and only 10% to the mouse flows [5]. This paradox highlights the importance of elephant flows.



### 3-3- Mechanisms to detect Elephant and mice flows

The mechanisms for detecting elephant flow from mice are divided into two main categories [27]:

*Detection by edge switches:* In this case, edge switches are responsible for detecting elephant flows. Hedera [7] and DiFS [19] use this method.

*Detection by hosts:* In this method, the detection of elephant from the mice flows is the responsibility of the host itself. Mahout [27] and DARD [4] are some of the algorithms that use this method.

### 3-4- The problem of collision between elephant flows

As previously mentioned in Section 2, DiFS is a greedy distributed mechanism for scheduling elephant flows in a data center network. In Fig. 2, suppose that hosts A and B produce 16 elephant flows in total, where the destination of 8 of these flows is host C or D (within the pod) while, other 8 flows are toward outside the pod. In switch SW1, DiFS distributes the flows quite evenly. But in the worst case, all of the eight flows that enter SW3 may be intra-pod flows. This situation leads to improper equilibrium of the elephant flows in the links between aggregation and core layers. Having a macroscopic view, one can easily achieve the proper load balance in the network (see Fig. 3).

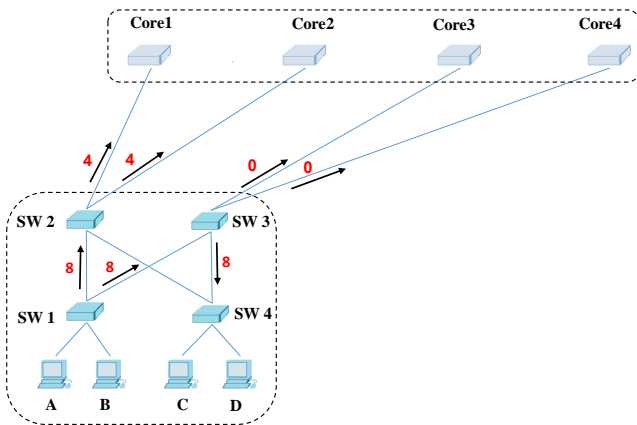


Fig. 2 The problem of load-balance in DiFS

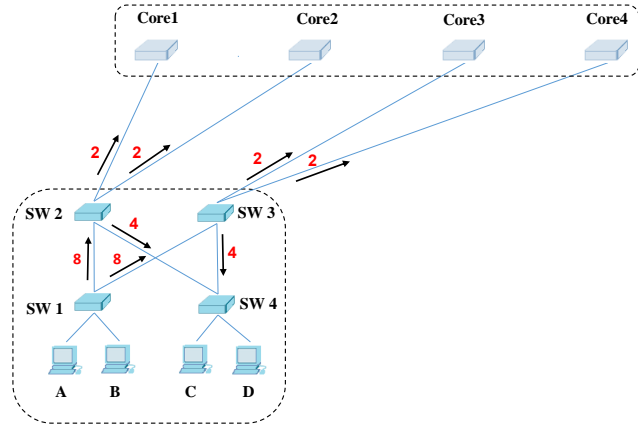


Fig. 3 Achieving a proper load balancing using macroscopic view

## 4- Proposed Method

In this section, we present an efficient mechanism for the flow scheduling problem in data center networks. To this end, we first give an overview of the proposed mechanism and then describe it.

### 4-1- Overview of the proposed method

The proposed algorithm has two main objectives. First, it aims to evenly distribute the load across the network. Second, it does not impose too much overhead on the central controller to achieve the first goal and can schedule the elephant flows at an acceptable speed. To manage traffic on a data center network, the proposed mechanism uses per-flow approach. This approach prevents out-of-ordering of packets in end-hosts. As a result, we will not confront a degradation in TCP performance and end-host memory usage. Our mechanism combines the advantages of both distributed and centralized systems. Due to their global view, centralized systems are very suitable for routing elephant flows, while distributed systems are the best option for routing mice flows to avoid overloading the central controller.

For the centralized system, we use a bidirectional search algorithm for scheduling elephant flows, which we describe in the following subsection. While for the distributed system, we use a simple yet efficient ECMP algorithm for mice flows. It is worth noting that in the proposed method, such as the mechanism presented in [27], the elephant flows are detected in the end-hosts. Similar to many of the existing works [4, 19, 27], we consider flows with a volume less than 100KB as mice flows and assume the others as elephant flows. In this work, we use the number of elephant flows as a load balancing parameter and the goal is to keep the number of active elephant flows on the network links as equal as

possible. Although other parameters such as *current bandwidth consumption* can be used for this purpose, the results presented in [19] show that this parameter would give us similar performance in practice. Fig. 4 shows the flowchart of the proposed method.

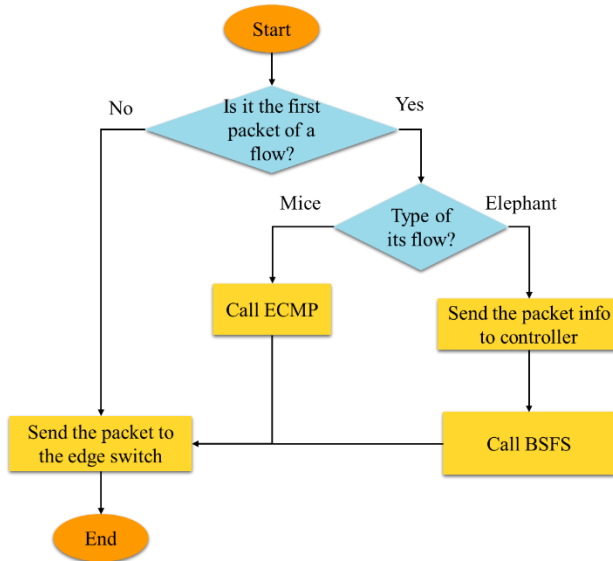


Fig. 4 The flowchart of the proposed method

#### 4-2- BSFS

For elephant flows, we use bidirectional search algorithm to find an optimal path for each of them. When an elephant flow is detected by the source host and its destination host does not have the same edge switch as the source host, the packets of that flow are labeled with an “E”, indicating that the flow is elephant. Upon arrival of the first packet from an elephant flow to the edge switch, that switch sends the source and destination address of the packet to the controller to find the appropriate route. The controller executes the proposed BSFS algorithm and, through the OpenFlow protocol, installs routing information on the switches in the path suggested by the algorithm. On the other hand, when the last packet of a flow is processed by the source edge switch, a request to update the network traffic information is sent to the controller. The details of the proposed algorithm are discussed below.

As mentioned previously in subsection 3.1, we use a directed graph to represent the Fat-tree topology. Based on this graph, we create a cost matrix  $M_{m \times 2n}$ , where  $m$  is the number of network switches and  $n$  is the number of ports per switch. Each element of this matrix represents the number of active elephant flows on each network link. The reason for using  $2n$  numbers for each switch is that we use a directed graph to model the topology;  $n$  numbers for uphill links and  $n$  numbers for downhill links.

When a packet from an elephant flow is sent to the controller for routing, depending on the source and destination address of the packet, the controller can determine whether the two hosts are in the same pod or they are located in separate pods. If two hosts are in the same pod, the BSFS algorithm selects the best aggregation switch as the intermediate switch using a simple bidirectional search. But if the two hosts are located in two separate pods, the proposed BSFS algorithm starts two searches simultaneously; first one from the source edge switch to the core switches, searching between uphill links, and the other one from the destination edge switch to the core switches, searching between the downhill links. The aggregated result of these two searches is obtained for each of the core switches, and finally, using a simple linear search, the core switch that gives us a smaller value is chosen for routing. It is worth mentioning that in the Fat-tree topology, when the core switch is specified, there will be only one path between each pair of hosts [2]. It is also important to note that since the two searches are completely independent of each other, they can be run in parallel, which significantly reduces the execution time. On the other hand, when the last packet of an elephant flow is reported to the controller, the cost matrix is immediately updated; That is, one unit is reduced from the cost of all links that were along that flow. **Algorithm 1** shows the pseudo-code of the proposed BSFS.

The algorithm takes the cost matrix, packet  $p$  (the first or last packet of a flow), the source address, and the destination address of the host as input. If  $p$  is the first packet of a flow, using the bidirectional search method in the cost matrix, the path with the lowest cost is found for that flow (lines 1 to 6). Otherwise, if  $p$  is the last packet of a flow, the cost matrix is updated (lines 7 to 9); that means the cost of all the links along that flow is decreased by one.

---

**Algorithm 1.** BSFS: Bidirectional Search Algorithm for Flow Scheduling

---

**Input:** Cost Matrix, packet  $p$ ,  $p.src$ ,  $p.des$

**Output:** Optimal Path

1. **if**  $p$  is the first packet of a flow **then**
2.     **if**  $p.src$  and  $p.des$  belong to the same pod **then**
3.         Find an aggregation switch with minimum cost using BS // BS stands for Bidirectional Search
4.     **else**
5.         Find a core switch with minimum cost using BS
6.     **end if**
7. **else if**  $p$  is the last packet of a flow **then**
8.     Update the Cost Matrix and the flow table of related switches

9. **end if**

---

### 4-3- Time Complexity Analysis

Here, we analyze the time complexity of our proposed BSFS method. For the first packet of each elephant flow, the time complexity of the proposed algorithm is as follows. If the source and destination host of a packet have the same pod, our algorithm searches among the aggregation switches inside that pod to find a switch with minimum cost (lines 2 and 3). Since the number of aggregation switches is equal to  $n/2$ , this step needs  $O(n)$ . However, when the source and destination host of a packet are within different pods, our algorithm must find a core switch with minimum cost (lines 4 and 5). Regarding to the fact that the number of core switches in a Fat-tree topology is  $n^2/4$  [2], so the time complexity is  $O(n^2)$ . As a result, the time complexity of the proposed algorithm is  $O(n^2)$  for a Fat-tree topology with  $n$ -port switches. We should mention that in updating the cost matrix (lines 7 and 8), only three (in intra-pod case) or six switches (in inter-pod case) are involved for each flow, which is constant numbers.

It is worth to note that our BSFS method runs only for the first packet of elephant flows. Since the number of elephant flows in a data center usually is very low, the time complexity of our algorithm is reasonable.

## 5- Performance Evaluation

In this section, we evaluate the performance of our proposed BSFS algorithm. We compare it with Static [2], ECMP [18] and DiFS [19] in various respects. It should be noted that in this work we neglect DiFS performance degradation due to packet re-ordering.

### 5-1- Simulation settings

In this work, evaluation of the proposed algorithm on Fat-tree topology with 8-port switches is performed. C++ programming language has been used for simulation of the proposed method. In the literature, there are many works that use custom simulators [7, 9, 28, 29]. Experiments have been performed on a computer having Intel® Core™ i5 CPU 2.3 GHz and 16 GB of memory.

The event-driven simulation is developed on a packet level. The length of each packet is assumed to be 1KB. For each port, a buffer of size 64KB is assumed. The capacity of all network links is equal and set to 1Gbps. For the transmission delay, we consider  $8\mu s$  while the propagation delay is ignored. In this work, the queuing delay has been considered.

In our experiments, each server generates 20 flows continuously. We consider each flow with the probability of 90% as a mice flow and with the probability of 10% as an elephant flow. The size of mice flows is chosen randomly from the values 2KB, 10KB or 100KB. For the elephant flows, we assume the fixed size of 10MB.

### 5-2- Traffic patterns

We have used the following synthetic traffic patterns to perform the experiments [7, 13, 19]:

*Stride(i)*: This pattern sends a flow from host  $x$  to another host with the number  $(x + i)\% N$ ; where,  $N$  represents the number of hosts in the network.

*Random*: In this traffic pattern, a host with index  $x$  sends a flow randomly with uniform probability to another host  $y$  anywhere in the network, such that,  $x \neq y$ .

*Staggered( $P_{host}, P_{pod}$ )*: In this pattern, a host sends its flows with the probability of  $P_{host}$  to another host connected to the same edge switch, and with the probability of  $P_{pod}$  to another host in the same pod. It also sends the flows to other hosts with different pods with the probability of  $1 - P_{host} - P_{pod}$ .

### 5-3- Evaluation criteria

We use the following criteria to evaluate and compare our proposed method with other mechanisms:

*Flow Completion Time (FCT)*: This criterion specifies the end time of a flow. In fact, it indicates the time when all packets of a flow are received by the destination.

*Delay*: Indicates average network latency. This criterion tells us that how long it takes in average for a packet to reach its destination.

*Aggregate Throughput*: This criterion measures the utilization of network links. In fact, it indicates the average rate at which the network delivers the packets.

### 5-4- Simulation results

Here, we evaluate the results of simulations. **Fig. 5(a)** and **Fig. 5(b)** show the average delay and network aggregate throughput under different traffic patterns, respectively. As can be seen, for *Stride(2)* and *Staggered(0.5,0.3)* almost all methods have low delay and high throughput. However, in *Stride(i)*, by increasing the value of  $i$  and in *Staggered( $P_{host}, P_{pod}$ )* by decreasing the values of  $P_{host}$  and then  $P_{pod}$ , BSFS performs much better. For the *Random* traffic pattern, since it is more likely for elephant flows to collide and most of the flows will be out of the pods, our proposed algorithm performs best by establishing a proper balance between the elephant flows.

Table 1. The Flow Completion Time of different algorithms under Stride (4)

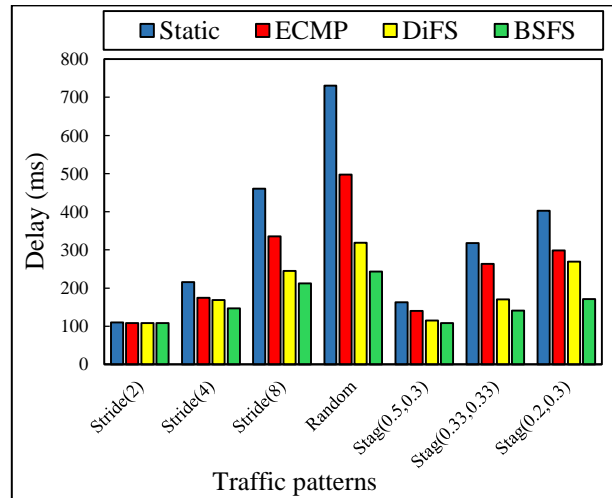
Algorithm	Time (in millisecond)										
	0	50	100	150	200	400	600	800	1000	1200	1400
Static	0	969	969	969	1087	1899	2088	2413	2528	2541	2560
ECMP	0	1011	1011	1011	1421	1913	2120	2454	2549	2551	2560
DiFS	0	1031	1031	1031	1561	2005	2299	2479	2551	2560	-
<b>BSFS</b>	<b>0</b>	<b>1030</b>	<b>1030</b>	<b>1030</b>	<b>1695</b>	<b>2120</b>	<b>2356</b>	<b>2560</b>	-	-	-

Table 2. The Flow Completion Time of different algorithms under Random

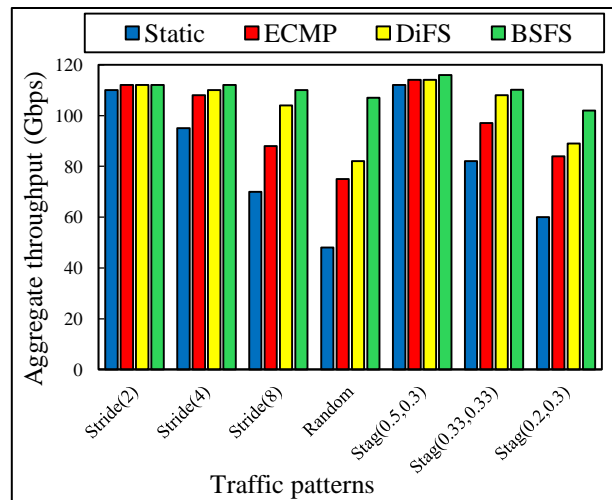
Algorithm	Time (in millisecond)															
	0	50	100	150	200	400	600	800	1000	1200	1400	1600	1800	2000	2200	2400
Static	0	976	980	984	1129	1310	1582	1671	1860	2088	2329	2434	2490	2519	2540	2560
ECMP	0	983	988	989	1207	1441	1740	1949	2081	2219	2416	2480	2529	2539	2560	-
DiFS	0	971	973	974	1439	1741	1961	2134	2262	2386	2481	2537	2549	2560	-	-
<b>BSFS</b>	<b>0</b>	<b>905</b>	<b>916</b>	<b>918</b>	<b>1393</b>	<b>1968</b>	<b>2389</b>	<b>2555</b>	<b>2560</b>	-	-	-	-	-	-	-

In particular, for the Random traffic pattern, BSFS provides about 20% higher throughput than DiFS. Therefore, we claim that our proposed mechanism performs better for non-local traffic than other mechanisms.

Table 1 and Table 2 illustrate the cumulative distribution function of number of completed flows for different algorithms under two traffic patterns Stride(4) (Table 1) and Random (Table 2). It can be clearly seen that BSFS terminates the flows earlier. For the Random traffic pattern, this superiority is much more impressive. This is due to the balanced distribution of elephant flows by the proposed method. While, in other algorithms, there is a longer flow completion time due to the frequent collisions between the elephant flows. For the Random traffic pattern, our BSFS delivers all the flows to their destination hosts below one second, while DiFS delivers about 88%, ECMP 81% and Static only deliver 72% of flows at this time.



(a) delay



Traffic patterns

(b) Aggregate throughput

**Fig. 5** Performance comparison of algorithms under different traffic patterns

## 6- Conclusion and Future Work

In this paper, we proposed an efficient mechanism to achieve load balancing in data center networks. The proposed mechanism uses the ECMP algorithm to send mice flows, while it takes advantage of the bidirectional search algorithm in the central controller to schedule the elephant flows. Simulation results under various traffic patterns show that the proposed mechanism can balance the network load more efficiently and provide better performance in comparison with the Static, ECMP and DiFS mechanisms. The less locality in network traffic, the higher the advantage of our approach is. Specifically, for the Random traffic model, our mechanism provides 20% higher throughput than DiFS. As a possible future research direction, one can take into account the priority of flows when scheduling them. Furthermore, the proposed mechanism can be extended by considering failures in data center.

## References

- [1] Zhang, J., Yu, F. R., Wang, S., Huang, T., Liu, Z., Liu, Y.: Load Balancing in Data Center Networks: A Survey, *IEEE Communications Surveys & Tutorials*, vol. 20, no. 3, pp. 2324-52, 2018.
- [2] Al-Fares, M., Loukissas, A., Vahdat, A.: A scalable, commodity data center network architecture, In *ACM SIGCOMM Computer Communication Review*, vol. 38, no. 4, pp. 63-74, 2008.
- [3] Niranjana Mysore, R., Pamboris, A., Farrington, N., Huang, N., Miri, P., Radhakrishnan, S., Subramanya, V. and Vahdat, A.: Portland: a scalable fault-tolerant layer 2 data center network fabric, In *ACM SIGCOMM Computer Communication Review*, vol. 39, no. 4, pp. 39-50, 2009.
- [4] Wu, X., Yang, X.: Dard: Distributed adaptive routing for datacenter networks, In *32nd International Conference on Distributed Computing Systems (ICDCS)*, pp. 32-41, 2012, IEEE.
- [5] Greenberg, A., Hamilton, J.R., Jain, N., Kandula, S., Kim, C., Lahiri, P., Maltz, D.A., Patel, P. and Sengupta, S.: VL2: a scalable and flexible data center network, In *ACM SIGCOMM computer communication review*, vol. 39, no. 4, pp. 51-62, 2009.
- [6] Zahavi, E., Keslassy, I., Kolodny, A.: Distributed adaptive routing for big-data applications running on data center networks, In *Proceedings of the eighth ACM/IEEE Symposium on Architectures for networking and communications systems*, pp. 99-110, 2012.
- [7] Al-Fares, M., Radhakrishnan, S., Raghavan, B., Huang, N., Vahdat, A.: Hedera: Dynamic Flow Scheduling for Data Center Networks, In *NSDI*, vol. 10, pp. 19-19, 2010.
- [8] Zats, D., Das, T., Mohan, P., Borthakur, D., Katz, R.: DeTail: reducing the flow completion time tail in datacenter networks, In *Proceedings of the ACM SIGCOMM conference on Applications, technologies, architectures, and protocols for computer communication*, pp. 139-150, 2012.
- [9] Sen, S., Shue, D., Ihm, S., Freedman, M.J.: Scalable, optimal flow routing in datacenters via local link balancing, In *Proceedings of the ninth ACM conference on Emerging networking experiments and technologies*, pp. 151-162, 2013.
- [10] Modi, T., Swain, P., FlowDCN: Flow Scheduling in Software Defined Data Center Networks. In *IEEE International Conference on Electrical, Computer and Communication Technologies (ICECCT)*, pp. 1-5, 2019.
- [11] Alizadeh, M., Edsall, T., Dharmapurikar, S., Vaidyanathan, R., Chu, K., Fingerhut, A., Matus, F., Pan, R., Yadav, N. and Varghese, G.: CONGA: Distributed congestion-aware load balancing for datacenters, In *ACM SIGCOMM Computer Communication Review*, vol. 44, no. 4, pp. 503-514, 2014.
- [12] He, K., Rozner, E., Agarwal, K., Felter, W., Carter, J. and Akella, A.: Presto: Edge-based load balancing for fast datacenter networks, In *ACM SIGCOMM Computer Communication Review*, vol. 45, no. 4, pp. 465-478, 2015.
- [13] Zhang, J., Ren, F., Huang, T., Tang, L., Liu, Y.: Congestion-aware adaptive forwarding in datacenter networks, *Computer Communications*, vol. 62, pp. 34-46, 2015.
- [14] Cui, W., Qian, C.: Difs: Distributed flow scheduling for adaptive routing in hierarchical data center networks, In *Proceedings of the tenth ACM/IEEE Symposium on Architectures for networking and communications systems*, pp. 53-64, 2014.
- [15] Ghorbani, S., Yang, Z., Godfrey, P., Ganjali, Y. and Firoozshahian, A.: DRILL: Micro load balancing for low-latency data center networks, In *Proceedings of*

*the Conference of the ACM Special Interest Group on Data Communication*, pp. 225-238, 2017.

- [16] Wang, C., Zhang, G., Chen, H. and Xu, H.: An ACO-based elephant and mice flow scheduling system in SDN, In *2nd International Conference on Big Data Analysis (ICBDA)*, pp. 859-863, 2017, IEEE.
- [17] Perry, Perry, Balakrishnan, H., Shah, D.: Flowtune: Flowlet Control for Datacenter Networks, In *NSDI*, pp. 421-435, 2017.
- [18] Hopps, C.E.: Analysis of an equal-cost multi-path algorithm, 2000.
- [19] Cui, W., Yu, Y., Qian, C.: DiFS: Distributed Flow Scheduling for adaptive switching in FatTree data center networks, *Computer Networks*, vol. 105, pp. 166-179, 2016.
- [20] Ghorbani, S., Godfrey, B., Ganjali, Y. and Firoozshahian, A.: Micro load balancing in data centers with DRILL, In *Proceedings of the 14th ACM Workshop on Hot Topics in Networks*, p. 17, 2015.
- [21] Azizi, S., Hashemi, N., Khonsari, A.: A flexible and high-performance data center network topology, *The Journal of Supercomputing*, vol. 73, no. 4, pp. 1484-1503, 2017.
- [22] Guo, C., Wu, H., Tan, K., Shi, L., Zhang, Y. and Lu, S.: Dcell: a scalable and fault-tolerant network structure for data centers, In *ACM SIGCOMM Computer Communication Review*, vol. 38, no. 4, pp. 75-86, 2008.
- [23] Wang, W., Sun, Y., Salamatian, K., Li, Z.: Adaptive Path Isolation for Elephant and Mice Flows by Exploiting Path Diversity in Datacenters, *IEEE Transactions on Network and Service Management*, vol. 13, no. 1, pp. 5-18, 2016.
- [24] Li, P., Xu, H., Wang, R., Luo, B.: Data center network flow scheduling mechanism based on HGSAFS algorithm, In *Proceedings of the High Performance Computing Symposium*, 2019.
- [25] Li, D., Wu, J.: On data center network architectures for interconnecting dual-port servers, *IEEE Transactions on Computers*, vol. 64, no. 11, pp. 3210-3222, 2015.
- [26] Marron, J., Hernandez-Campos, F., Smith, F.: Mice and elephants visualization of internet traffic, In *Proceedings of Computational Statistics*, pp. 47-54, 2002.
- [27] Curtis, A.R., Kim, W., Yalagandula, P.: Mahout: Low-overhead datacenter traffic management using

end-host-based elephant detection, In *Proceedings IEEE INFOCOM*, pp. 1629-1637, 2011.

- [28] Wischik, D., Raiciu, C., Greenhalgh, A., Handley, M.: Design, Implementation and Evaluation of Congestion Control for Multipath TCP, In *NSDI*, vol. 11, pp. 8, 2011.
- [29] Singla, A., Hong, C.-Y., Popa, L., Godfrey, P.B.: Jellyfish: Networking Data Centers, Randomly, In *NSDI*, vol. 12, pp. 1-6, 2012.

**Hasibeh Naseri** is a master graduated in artificial intelligence from University of Kurdistan, Sanandaj, Iran. Currently, she is a member of the Internet of Things research laboratory at the University of Kurdistan. Her research interests include Cloud computing, Data center networks, Flow scheduling, Heuristic and meta-heuristic algorithms.

**Sadoon Azizi** is an assistant professor in the department of computer engineering and IT, University of Kurdistan, Sanandaj, Iran. He received his Ph.D degree in computer science, with focus on Cloud Data Centers, from Amirkabir University of Technology, Tehran, Iran, in 2016. He also received his M.Sc. degree in computer science, with focus on High Performance Computing (HPC), from Amirkabir University of Technology, Tehran, Iran, in 2012. His main research interests include Cloud computing, Fog computing, Internet of Things, Data center networks, and Design and analysis of algorithms. He is the director of the Internet of Things research laboratory and manager at the High Performance Computing (HPC) center at the University of Kurdistan.

**Alireza Abdollahpouri** received his Ph.D. from University of Hamburg, Germany in 2012, and now he is an associate professor of computer networks at the Department of Computer Engineering, University of Kurdistan, Sanandaj, Iran. His main research interests are in the field of social network analysis, IPTV modeling, and quality of service in wireless networks.

# Balancing Agility and Stability of Wireless Link Quality Estimators

Mohamad Javad Tanakian

Department of Telecommunication Engineering, Faculty of Electrical & Computer Engineering,  
University of Sistan and Baluchestan, Zahedan, Iran  
mj.tanakian@pgs.usb.ac.ir

Mehri Mehrjoo\*

Department of Telecommunication Engineering, Faculty of Electrical & Computer Engineering,  
University of Sistan and Baluchestan, Zahedan, Iran  
mehrjoo@ece.usb.ac.ir

Received: 04/Mar/2019

Revised: 17/Oct/2019

Accepted: 02/Dec/2019

## Abstract

The performance of many wireless protocols is tied to a quick Link Quality Estimation (LQE). However, some wireless applications need the estimation to respond quickly only to the persistent changes and ignore the transient changes of the channel, i.e., be agile and stable, respectively. In this paper, we propose an adaptive fuzzy filter to balance the stability and agility of LQE by mitigating the transient variation of it. The heart of the fuzzy filter is an Exponentially Weighted Moving Average (EWMA) low-pass filter that its smoothing factor is changed dynamically with fuzzy rules. We apply the adaptive fuzzy filter and a non-adaptive one, i.e., an EWMA with a constant smoothing factor, to several types of channels from short-term to long-term transitive channels. The comparison of the filters outputs shows that the non-adaptive filter is stable for large values of the smoothing factor and is agile for small values of smoothing factor, while the proposed adaptive filter outperforms the other ones in terms of balancing the agility and stability measured by the settling time and coefficient of variation, respectively. Notably, the proposed adaptive fuzzy filter performs in real time and its complexity is low, because of using limited number of fuzzy rules and membership functions.

**Keywords:** Link quality estimation; adaptive fuzzy filter; agility; stability; wireless channel.

## 1- Introduction

Telecommunication network is deployed in smart grid to exchange the measurement status and instructions of numerous widely distributed control devices of power grid. Among different telecommunication networks, Wireless Networks (WNs), because of the low cost and flexibility of installation and maintenance are more probable to be deployed for monitoring, collecting data and controlling smart grid assets [1],[2]. The success of WN applications depends on the reliable transmission of sensory data. Reliability is defined as the success rate of source to destination data transmission in the network within its required latency. Accordingly, research community has been paying significant attention to design and implementation of reliable data transmission protocols in WNs [3],[4],[5],[6],[7],[8]. The performance of a large number of these protocols is highly dependent on Link Quality Estimation (LQE) [3],[4],[6],[7],[8]. Poor LQE may lead to an unstable network with high packet loss and/or high delay.

Performance of LQE is assessed in terms of accuracy, cost, agility and stability [9]. Accuracy is quantified by comparing the measured link quality and the estimated

link quality using the Mean Square Error (MSE) metric. Consuming the energy by excessive re-transmissions, occurred by imperfect link estimation, over low quality links is inferred as cost. Agility is the ability to react quickly to persistent changes in link quality. Agility is measured by settling time, defined as the time needed by the estimator to reach the measured value within an error bound of  $\epsilon$ . Finally, stability is the ability to resist the short-term variations, a.k.a, fluctuations, in link quality. Stability is assessed quantitatively, by the Coefficient of Variation (CV) defined as the ratio of the standard deviation to the mean of variations. Balancing between agility and stability is of paramount importance in a deployed LQE in WNs. In general, whenever the overhead of signaling is high and the decision is made based on the channel status the balancing between stability and agility becomes critical. For example handover and scheduling in cellular network or routing in wireless local area network. Routing protocols do not have to reroute information when a link quality shows transient degradation, because rerouting is a very energy and time-consuming operation. Too frequent protocol updates may cause unexpected network problems, such as, routing loops and routing shocks [10].

\* Corresponding Author

LQEs are classified in two categories, hardware and software based LQEs [9], [11]. In the hardware based LQEs, the estimation is based on the measurement of a dedicated signal either on the transmitter or the receiver side and do not require any further computation; however, they are not as good as software based LQEs [9],[13]. The received signal strength (RSS), link quality indicator (LQI), and signal-to-noise ratio (SNR) are primary metrics used in hardware-based LQE [11]. In particular, none of these metrics by itself is sufficient to accurately characterize the quality of a link because [11],[13]: (i) the RSS is not sensitive to changes in link quality; (ii) the variance of LQI readings is significantly increased for transitional links, and (iii) the SNR rapidly and randomly fluctuates. Software based LQEs are divided into two categories: (1)Single metric based,(2)Hybrid metric based. Single metric based estimators count or approximate either (i) the reception rate or (ii) the average number of packet transmissions/retransmissions required before its successful reception. For instance, Packet Reception Rate (PRR) of a wireless link over an estimation window consisting of  $w$  instances of communication and Acquitted Reception Rate (ARR) count the reception rate at receiver side and sender side, respectively [9],[11]. Required Number of Packet transmissions (RNP) counts the average number of packet transmissions/retransmissions required before its successful reception within a window of  $w$  communication instances [9]. Furthermore, the expected transmission count (ETX) takes into account link asymmetry by estimating the uplink quality and downlink quality using both forward and backward PRR values, respectively [12].

Hybrid metric based LQEs consider a number of link quality metrics. For instance, Four-bit LQE combines individual estimations of uplink and downlink qualities based on measured RNP and PRR, respectively [14]. Stable Link Quality Estimation (SLQE) combines active probing with passive snooping to make a stable estimation [10]. In this estimator, an active node sends control packets periodically and uses long period active detection mechanisms to detect quality of the link, while a passive node listens RSS Indicator mean and perceives links in sudden changes effectively. Fuzzy Link Quality Estimator (F-LQE) deploys four link quality properties, namely, packet delivery, link quality difference of forward and backward direction (asymmetry), stability, and SNR of a transceiver [15]. Each of the link properties is considered as a different fuzzy variable. Opt-FLQE (Optimized FLQE) [16] is a modification of F-LQE that aims to improve its reactivity and to reduce its computational complexity. A method that uses fuzzy logic to combine LQI, SNR and PRR metrics is proposed in [17] to improve the accuracy rate for evaluating a link quality. Fuzzy logic based link quality indicator (FLI) uses the PRR, the coefficient of variance of PRR, and a metric to assess the

burstiness of packet loss, to estimate link quality [18]. Remarkably, all the research works on fuzzy link quality estimator have limited the application of the fuzzy system to combine some ELQ metrics. Kalman filter based LQE approximates the packet reception ratio based on RSSI and a pre-calibrated PRR/SNR curve [19].

In the cases of the PRR, RNP, and ETX, there is a tradeoff between estimation accuracy and latency. For example, the estimation latency can be improved by shortening the window size  $w$ , but at the cost of increased fluctuation in the estimation results and degraded estimation accuracy [11],[17]. To address this problem, some LQEs apply Exponentially Weighted Moving Average (EWMA) filter on the estimated link quality (ELQ) which smooth the variations of it to turn them robust against the fluctuations [9],[11],[15]. In these LQEs, the EWMA, a non adaptive Infinite Impulse Response (IIR) filter, is tuned by a constant smoothing factor  $\alpha$ , where  $0 \leq \alpha \leq 1$ . A stricter smoothing filter, to remove the transient variations, is needed when fluctuation amplitude is high. On the other hand, the strict smoothing filter prevents following link quality status when it has a relatively high persistent change in link quality [9],[10],[15],[17]. Therefore, the smoothing factor of filter should be tuned carefully proportional to the amount of fluctuation and the persistent changes in link quality.

In this paper, inspired by our previous research work on video stabilization [20], we propose an adaptive fuzzy system to tune the EWMA filter, named adaptive fuzzy filter. The fuzzy system has two inputs and one output, so it requires low computation resources and responds in real-time. As the inputs, the fuzzy system uses quantitative representations of the transient variations and the persistent changes in estimated link quality. The fuzzy inputs are defined according to the few numbers of the last estimated link qualities. The output of fuzzy system calculates the best value of smoothing factor to tune the EWMA filter adaptively. The performance of the proposed adaptive fuzzy filter in terms of stability and agility is compared with the results provided by an EWMA filter with a three different constant smoothing factors. Numerical results show that our proposed adaptive fuzzy filter provide balanced stable and agile estimation results, while the ones of a constant smoothing factor filter are either stable or agile, depending on the value of the filter smoothing factor.

The remainder of this paper is organized as follows. The basic concept and details of the proposed filter are described in Section 2. The numerical results are presented in Section 3, and the paper is concluded in Section 4.

## 2- Proposed Filter

The adaptive fuzzy filter consists of a fuzzy system and an EWMA filter whose smoothing factor is tuned by the



fuzzy system. In this section, first we briefly explain fuzzy system, and then we describe the proposed filter.

### 2-1- Fuzzy System

Fuzzy logic is an approach to computing based on "degrees of truth", rather than the usual "true or false"(1 or 0) Boolean logic on which the modern computer is based. Fuzzy logic starts with the concept of a fuzzy set. A fuzzy set is a set without a crisp, clearly defined boundary. It can contain elements with only a partial degree of membership. The fuzzy logic system incorporates five steps as shown in Fig.1. It starts with fuzzification process, then the inference system comes, including: application of the operators, implication methods, and aggregation all outputs to one fuzzy output, finally defuzzify the fuzzy output to numerical values [21],[22].

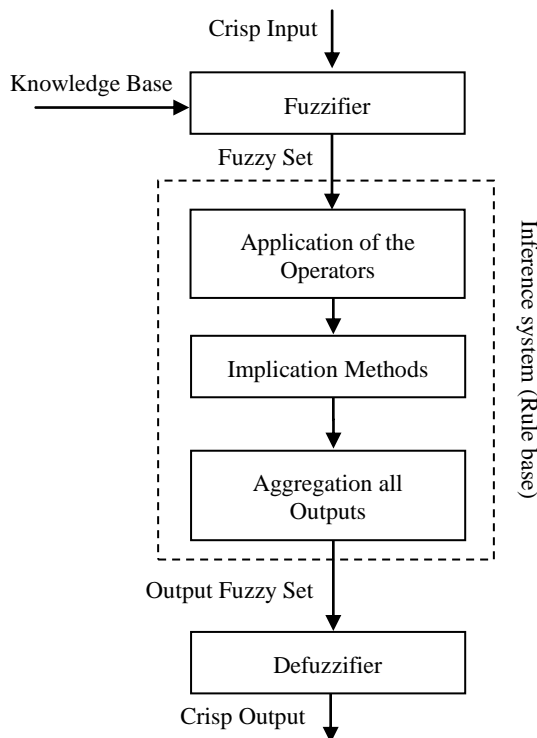


Fig. 1 The architecture of fuzzy logic system [22]

The fuzzy knowledge base includes rule base and the database. The rule base contains a number of IF-THEN rules, and the database defines the membership functions (MF) of the fuzzy sets. Fuzzifier converts the crisp input to a linguistic variable using the MF stored in the fuzzy knowledge base. Inference system converts the fuzzy input to the fuzzy output using IF-THEN fuzzy rules. Defuzzifier converts the fuzzy output of the inference system to crisp using membership functions analogous to the ones used by the Fuzzifier. The logic operators that combine the sets in the antecedent define the relationships

between input sets. This process includes three steps based on the rules of the fuzzy logic to be followed [22]: i)applying the operators of the rules when there is more than one part for the antecedent of the rule. This step results in one number (between 0 and represents all parts of the antecedent based on the operator of the rule .ii) finding the consequence of the rules by combining the rule strength and the output membership function which is defined as implication and iii) combining the consequences to get an output distribution which is defined as aggregation.

### 2-2- Adaptive Fuzzy Filter

The ELQ of a wireless link fluctuate over time due to many factors, principally related to the physical environment and the nature of low-power radios. Assuming that ELQ variation corresponds to its high-frequency components; we smoothen ELQ using a low-pass filter tuned by a fuzzy system to achieve Fuzzy Filtered Estimated Link Quality (FFELQ).The EWMA, first-order IIR filter, as the low-pass filter is applied to ELQ, at time interval  $n$ , and the FFELQ is resulted:

$$FFELQ(n) = \alpha(n) \times FFELQ(n-1) + (1 - \alpha(n)) \times ELQ(n) \quad (1)$$

The parameter  $\alpha$ ,  $0 \leq \alpha \leq 1$ , is regarded as the smoothing factor of the filter and adjusted by the fuzzy system in each time interval. The fuzzy system has two inputs (Input1, Input2) and one output. The Input1 and Input2, calculated during link quality status estimation, represent the amount of fluctuations and persistent changes in link quality at time interval  $n$ , respectively. The output of fuzzy system defines the smoothing factor  $\alpha$  of the EWMA filter. The block diagram of the proposed adaptive fuzzy filter is depicted in Fig. 2.

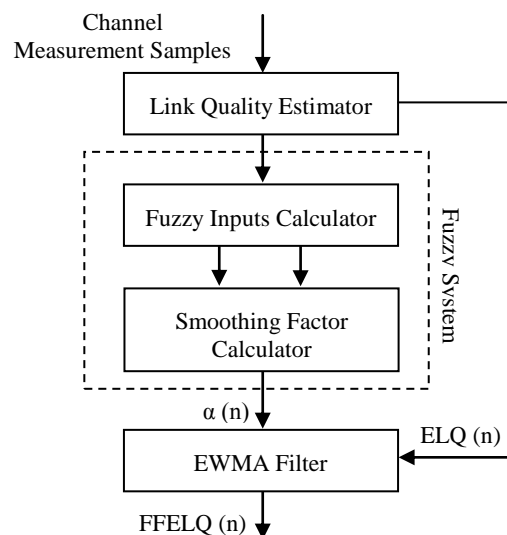


Fig. 2 Block diagram of adaptive fuzzy filter

We define the fuzzy inputs as

$$Input1(n) = \frac{1}{M} \sum_{i=n-M+1}^n |ELQ(i) - ELQ(i - 1)| \quad (2)$$

$$Input2(n) = |\sum_{i=n-M+1}^n (FFELQ(i - 1) - ELQ(i - 1))| \quad (3)$$

where  $M+1$  is the number of the last ELQs deployed in computation. Input1 is the average of absolute differences between consequent ELQs. Input2 is the absolute sum of difference between ELQ and FFELQ.

To justify Input1 and Input2 definitions, consider the two scenarios shown in Fig.3. The last four samples are considered for inputs computing. In Fig. 3(a) and 3(b), the amount of amplitude changes in the ELQ (solid line) are within the range of (0.47~0.50) and (0.54~0.65), respectively. The total amount of link quality fluctuations in Fig. 3(b) is more than the ones of Fig. 3(a). In addition, the FFELQ (dash line) follows ELQ path direction in Fig. 3(a), while the FFELQ in Fig. 3(b) is moving away from the ELQ path direction. Therefore, the Input2 is defined to reduce the deviation. The values of Input1 and Input2 are derived with (2) and (3) shown in Table1. The output of fuzzy system defines the smoothing factor of the EWMA filter, i.e.,  $\alpha$ .

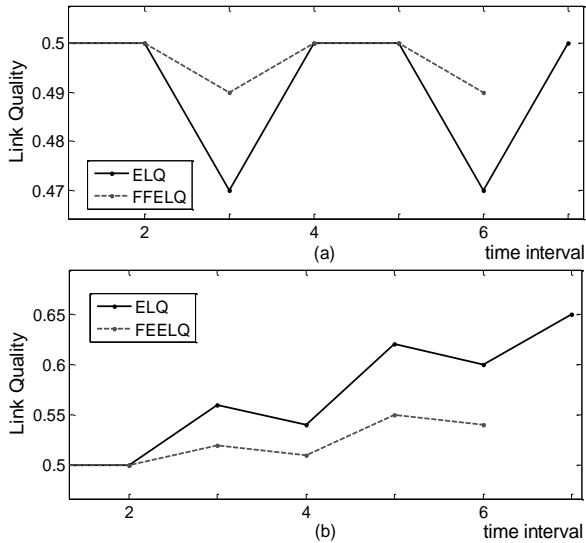


Fig. 3 (a) transient degradation with no persistent changes in link quality  
(b) fluctuation and the persistent changes in link quality

Table1: Values of fuzzy inputs for the two scenarios in Fig.3

	Input1	Input2
Scenario1 (Fig. 3(a))	0.02	0.02
Scenario2 (Fig. 3(b))	0.05	0.16

In the proposed fuzzy system, trapezoidal and triangular MFs are used for the inputs and the outputs, respectively. Trial and error method is used for MF shape of the inputs and the output. Type of MF doesn't play a crucial role in shaping how the model performs. However, the number of

MF has greater influence as it determines the computational time [23]. We select as few MFs as possible to maintain low system complexity while we obtain decent performance. The experimentally designed inputs and output MFs as well as the surface of the desired output, which is a graphical interface that allows you to examine the FIS output surface for two inputs, are shown in Fig.4.

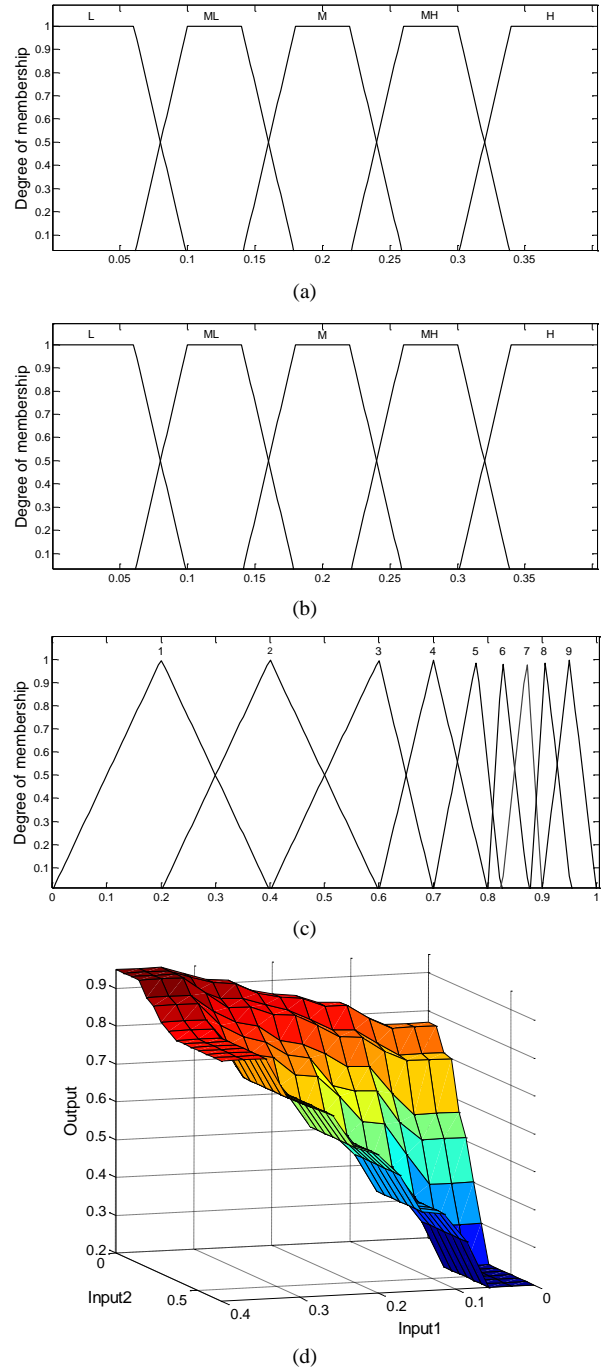


Fig. 4 (a) MFs of fuzzy Input1 (b) MFs of fuzzy Input2, (c) MFs of fuzzy output, (d) surface of desired outputs

Table2: Central Values of fuzzy system output

		Input 2				
		L	ML	M	MH	H
Input1	L	0.8	0.7	0.6	0.4	0.2
	ML	0.825	0.8	0.7	0.6	0.4
	M	0.875	0.825	0.8	0.7	0.6
	MH	0.9	0.875	0.825	0.8	0.7
	H	0.95	0.95	0.9	0.875	0.825

\* L=Low, ML=Medium Low, M=Medium, MH=Medium High, H=High.

According to experimental results, the performance of the EWMA filter is more sensitive to larger values of  $\alpha$  [20]. Therefore, more MFs of the fuzzy output are concentrated in this operating area. The constructed rule base is containing 25 rules as presented in Table 2. The proposed fuzzy system is implemented while the min function is used for the fuzzy implication and the max function is used for the fuzzy aggregation. Furthermore, the centroid defuzzification method is applied. After computing the smoothing factor  $\alpha$  (n) by the fuzzy system, FFELQ is calculated by Equation (1).

### 3- Numerical Results

In this section, the performance of the proposed adaptive fuzzy filter in terms of stability and agility of the LQE is compared with the results provided by a non-adaptive EWMA filter with a three different constant smoothing factors  $\alpha = 0.9$  [16],  $\alpha = 0.5$ , and  $\alpha = 0.2$ . The performance has been evaluated with several different recognized scenarios extracted from link quality status curves published in the literature [9], [10], [24], [25] and some synthetic link quality status trajectory: (a) link quality mutation frequently occurs in short times, (b) link quality remains unstable for a long time, (c) link quality is relatively stable, (d) link quality has persistent changes. To adjust the fuzzy system inputs, the initial value  $\alpha = 0.1$  is chosen for the first four time intervals. To keep the computation delay low, the filter window size  $M=3$  is chosen. The simulation tools is MATLAB 8.2.0.701 (R2013b, 32-bit) and the hardware configuration are: Intel(R) Core(TM) Duo CPU T9300 2.5GHz and 3.00GBRAM. The average simulation time for our adaptive fuzzy filter and non-adaptive filter with a constant smoothing factor are 3.4 msec and 10  $\mu$ sec, respectively. Therefore, to apply our adaptive fuzzy filter, the time interval between link quality calculations should be longer than 3.5 msec.

It is observed that a large smoothing factor, e.g.  $\alpha = 0.9$ , increases the stability of estimators at the expense of a relatively large delay when there are persistent changes in

link quality. Similarly, a small smoothing factor, e.g.  $\alpha = 0.2$ , closely tracks the persistent changes in calculated link quality at the cost of slightly reduced smoothing capabilities. In fact, a small smoothing factor just follows the original ELQs. Comparing the graphs shows that the proposed fuzzy filter provides expanded smoothing while enables the close tracking of the persistent changes in ELQs, i.e., the proposed method provides both agility and stability. The temporal behavior in Fig. 5(a), has sudden drop at  $t=7$ ,  $t=13$ , and  $t=25$ . The results demonstrate adaptive fuzzy filter and the non-adaptive filter with  $\alpha = 0.9$  smooth the ELQ and resist these temporary changes. While the non-adaptive filters with lower value of  $\alpha$  does not perform well against the transient fluctuations. The same observation can be seen in Fig. 5(b) and 5(c).

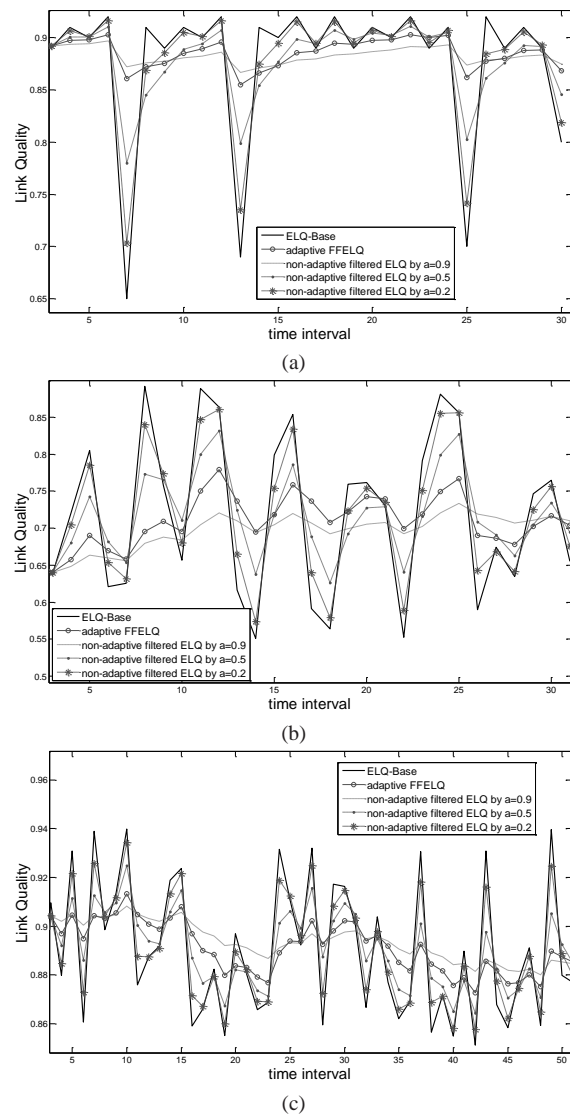


Fig. 5 Comparison results of EWMA filtering of low persistent ELQs with adaptive fuzzy filter and non-adaptive filter with different constant smoothing factors

In Fig. 6(a), (b), (c) and (d) the link quality curves have a high persistent change. The results show that the non-adaptive filter with  $\alpha = 0.9$  has a long delay to track the persistent changes. On the contrary, the delay is low when the smoothing factor is low. The adaptive fuzzy filter has a moderate delay in tracking the persistent change with the price of being stable in transient changes. In other words, the results shown in Fig.5 and Fig.6 indicate that the adaptive fuzzy filter can distinguish well between transient and non-transient changes of the link quality.

The Coefficient of Variation (CV), defined as the ratio of the standard deviation to the mean, shows the performance of the LQEs in terms of stability [9],[15]. The CV for the low persistent link quality curves shown in Fig.5, are presented in Table3. The lower CV represents the more stable estimation. The CV values of the filtered ELQs by the adaptive fuzzy filter and non-adaptive fuzzy with  $\alpha = 0.9$  are low compared to the two others. Hence the formers are more stable estimation.

Agility is measured by settling time (ST), defined as the time needed by the estimator to reach the measured value within an error bound of  $e$  [9]. The lower ST represents the more agile estimation. The CV and ST for the four link quality curves shown in Figure 6, are presented in Table4. The value of ST is in terms of time interval and  $e$  is about 5%. The numerical results show the adaptive fuzzy filter provide a balanced stable and agile estimation results, while the ones of constant smoothing factor filters are either stable or agile, depending on the value of the smoothing factor.

The empirical Cumulative Distribution Function (CDF) of two different links, which are shown in Figure 5(a) and Figure 6(d), is presented in Figure 7 for proposed adaptive fuzzy filter and non-adaptive filter with a three different constant smoothing factors  $\alpha = 0.9$ ,  $\alpha = 0.5$  and  $\alpha = 0.2$ . At the same time that the adaptive filter tries to balance between agility and stability, it should be confident to real quality of the link as much as possible. In other words, the proportions of link quality in terms of poor, moderate, or high quality in the CDF of basic ELQ, not filtered one, should remain almost the same in the CDF of the filtered ELQ. The presented scenario in Figure 5(a) shows a link with constant qualities equal to 0.9, and the presented scenario in Figure 6(d) shows that almost 28% of the link is in near to high quality; about 60% of the links is in intermediate quality; and about 12% of the link is in poor quality. According to the results shown in Fig. 7, adaptive LQE classify the link qualities close to the proportions set in these scenarios. The comparison results show that the non-adaptive filter with a smaller and larger constant smoothing factor over estimate and underestimate the link quality, respectively.

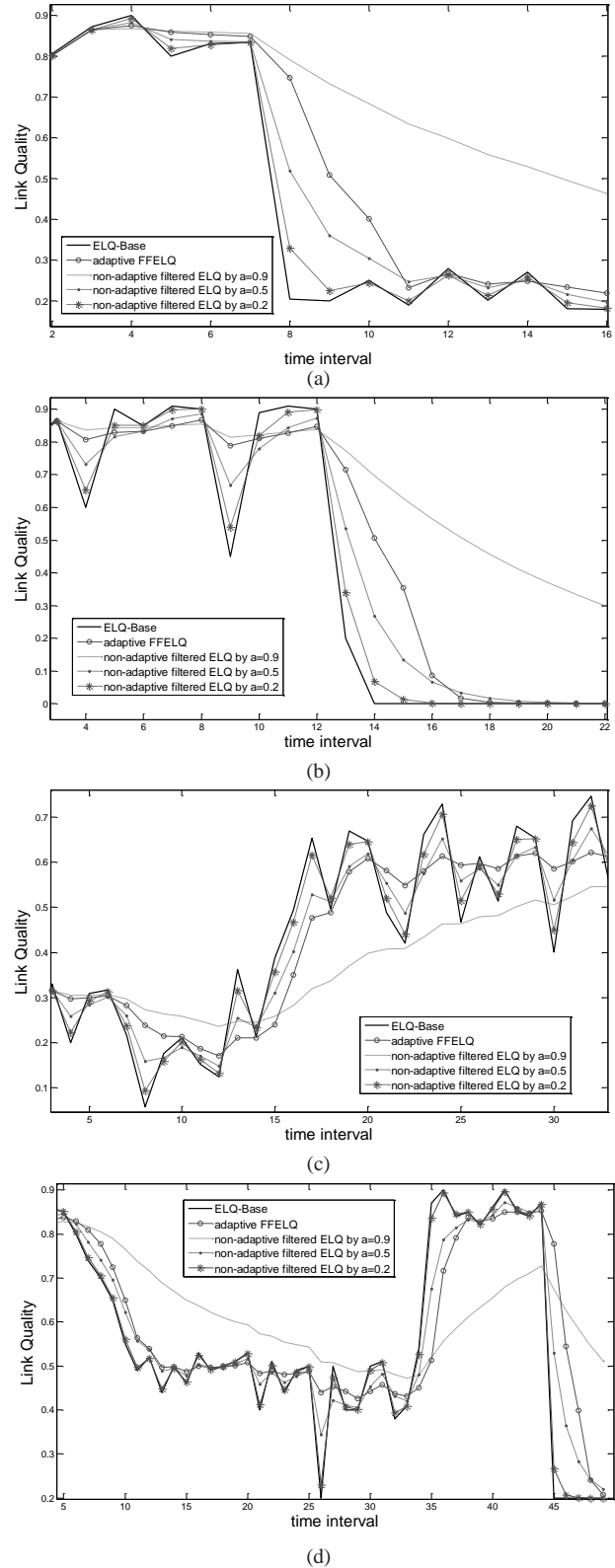


Fig. 6 Comparison results of EWMA filtering of high persistent ELQs with adaptive fuzzy filter and non-adaptive filter with different constant smoothing factors

Table 3: The coefficient of variation for presented results in Fig. 5

Link Quality	adaptive fuzzy filter	CV		
		non-adaptive filter		
		$\alpha = 0.9$	$\alpha = 0.5$	$\alpha = 0.2$
Link1 (Fig. 5(a))	0.0225	0.0184	0.0422	0.0647
Link2 (Fig. 5(b))	0.0496	0.0378	0.0781	0.1198
Link3(Fig. 5(c))	0.0230	0.0221	0.0261	0.0311

Table 4: The coefficient of variation and settling time for presented results in Fig.6

Link Quality	Criterion	adaptive fuzzy filter	non-adaptive filter		
			$\alpha = 0.9$	$\alpha = 0.5$	$\alpha = 0.2$
Link1 Fig. 6(a)	CV	0.5095	0.2039	0.5494	0.6202
	ST	4	more than 9	9	2
Link2 Fig. 6(b)	CV	0.7140	0.2901	0.7723	0.8791
	ST	4	more than 10	4	3
Link3 Fig. 6(c)	CV	0.4147	0.2954	0.4316	0.4576
	ST	4	more than 15	4	4
Link4 Fig. 6(d)	CV	0.3034	0.1791	0.3250	0.3748
	ST	4	more than 6	4	2

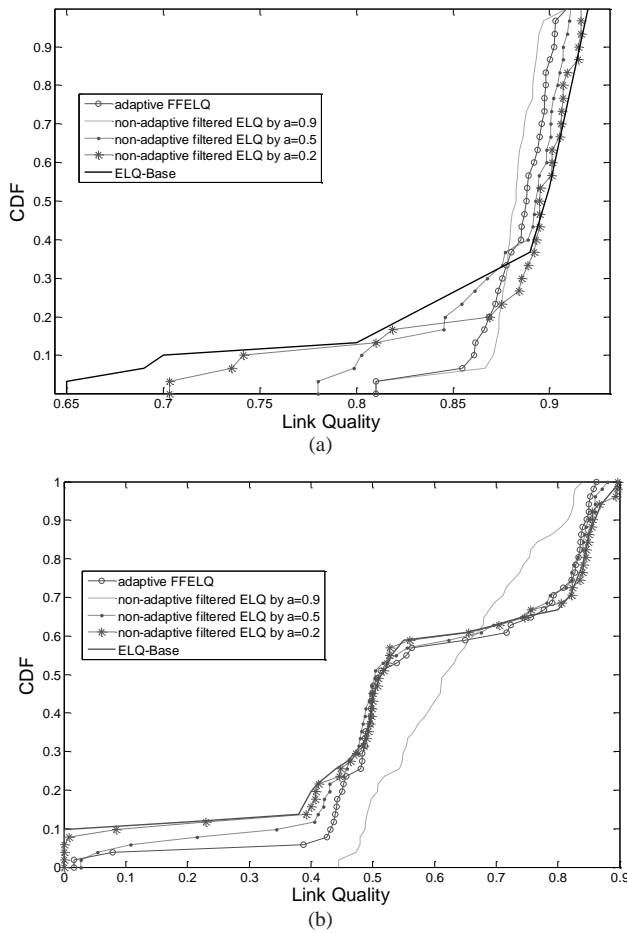


Fig 7. Empirical CDFs of link quality estimators for two scenarios are presented graphically in (a) Fig 5(a) and (b) Fig 6(d).

When measuring the quality of a link over a given period, all sorts of different scenarios may occur in the combination of noise and persistent changes in link quality. Generally, as shown in the Fig. 8, selecting a constant value for  $\alpha$  causes the EWMA filter output to work fine in either noise-canceling or in tracking the persistent changes, not necessarily both. Therefore, a dynamic value for  $\alpha$  is required. The fuzzy system adapts the system to different scenarios and chooses an appropriate value for  $\alpha$  at any given moment.

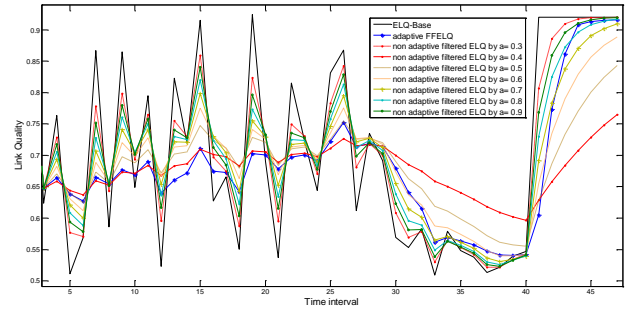


Fig 8. Comparison results of EWMA filtering with adaptive fuzzy filter and non-adaptive filter with different constant smoothing factors

### 4- Conclusion

An adaptive fuzzy filter to smooth transient variations of LQE has been proposed in this paper. The filter makes a balance between stability and agility in LQE. The proposed filter consists of an EWMA filter and a fuzzy system. The performance of the EWMA filter depends on the value of the smoothing factor which is tuned by the fuzzy system. The fuzzy system uses two inputs which are quantitative representations of the transient and the persistent changes in link quality status. In the fuzzy inputs, we selected as few MFs as possible to obtain decent performance with low system complexity. We have evaluated the filter in terms of stability and agility with CV and ST metrics. Numerical results show that our proposed adaptive fuzzy filter provides balanced, stable and agile, estimation results, while the ones of a constant smoothing factor filter are either stable or agile, depending on the value of the filter smoothing factor. The adaptive fuzzy filter is more complex with respect to the non-adaptive EWMA. However, the complexity cost is negligible with respect to the resource utilization improvement and/or signaling overhead reductions (e.g., rerouting).

## References

- [1] V.C. Gungor, D. Sahin, T. Koçak, S. Ergüt, C. Buccella, C. Cecati, and G.P. Hancke, "Smart grid technologies: Communication technologies and standards", *IEEE Transactions on Industrial Informatics*, Vol. 7, No. 4, 2011, pp. 529–539.
- [2] H.M.Nejad,N.Movahhedinia, and M.R.Khayyambashi, "Provisioning required reliability of wireless data communication in smart grid neighborhood area networks",*The Journal of Supercomputing*, Vol.73, No.2, 2017, pp. 866–886.
- [3] J. Cai, X.Song,J.Wang, and M.GU, "Reliability analysis for chain topology wireless sensor networks with multiple-sending transmission scheme" ,*EURASIP Journal on Wireless Communications and Networking*, Vol. 2014,No. 1, pp,156-169.
- [4] T.Korkmaz, and K.Sarac, "Characterizing link and path reliability in large-scale wireless sensor networks", *Proceedings of the 2010 IEEE 6th International Conference on Wireless and Mobile Computing, Networking and Communications (WiMob)*, 2010, pp. 217–224.
- [5] M. Kumar, R. Tripathi, and S. Tiwari, "QoS guarantee towards reliability and timeliness in industrial wireless sensor networks",*Multimedia Tools and Applications*, Vol. 77, No. 4, 2018, pp. 4491–4508.
- [6] W. Sun, W. Lu, Q. Li, L. Chen, D. Mu, and X. Yuan, "WNN-LQE: Wavelet-neural-network-based link quality estimation for smart grid WSNs", *IEEE Access*, Vol. 5, 2017, pp. 12788–12797.
- [7] F. Aalamifar, and L. Lampe, "Cost-efficient QoS-Aware Data Acquisition Point Placement for Advanced Metering Infrastructure", *IEEE Transactions on Communications (Early Access)*, Vol.66, No. 12, 2018, [Online]. Available:<https://arxiv.org/abs/1802.06656>.2018.
- [8] X.Zhu,Y.Lu,J.Han,andL.Shi, "Transmission Reliability Evaluation for Wireless Sensor Networks", *International Journal of Distributed Sensor Networks*,Vol. 2016,DOI:<http://dx.doi.org/10.1155/2016/1346079>.
- [9]N. Baccour, A.Koubaa,M.B.Jamaa,H.Youssef,M.Zuniga,andM.Alves, "A comparative simulation study of link quality estimators in wireless sensor networks", *IEEE International Symposium on Modeling, Analysis and Simulation of Computer and Telecommunication Systems (MASCOTS'09)*,2009,DOI: 10.1109/MASCOT.2009.5366798
- [10]A. Zhou,B.Wang,S.Xingming,X.You,H.Sun,andT.Li, "SLQE: an improved link quality estimation based on four-bit LQE",*International Journal of Future Generation Communication and Networking*, Vol.8, No.1, 2015, pp. 149–160.
- [11] W. Sun, X. Yuan, J. Wang, Q. Li, L. Chen, and D. Mu, "End-to-end data delivery reliability model for estimating and optimizing the link quality of industrial WSNs", *IEEE Transactions on Automation Science and Engineering*, Vol. PP, No. 99, 2017, pp. 1–11.
- [12] D. S. De Couto, D. Aguayo, J. Bicket, and R. Morris, "A high-throughput path metric for multi-hop wireless routing", *Wireless Networks.*, Vol. 11, No. 4, 2005, pp. 419–434.
- [13] A. Vlavianos, L.K. Law, L. Broustis, S. Krishnamurthy, and M. Faloutsos, "Assessing link quality in wireless networks: Which is the right metric?" ,*IEEE 19th International Symposium on Personal, Indoor and Mobile Radio Communications*, 2008, pp. 1–6.
- [14] R.Fonseca,O.Gnawali,K.Jamieson,andP.Lewis, "Four-Bit Wireless Link Estimation", In *Proceedings of the 6th Workshop on Hot Topics in Networks (HotNetsVI)*, 2007.
- [15] N. Baccour, A.Koubaa, H.Youssef, and M.Alves, "Reliable link quality estimation in low-power wireless networks and its impact on tree-routing", *Ad Hoc Networks*,Vol. 27, 2015, pp.1–25.
- [16] S. Rekik, N. Baccour, M. Jmaiel, and K. Drira , " Low-power link quality estimation in smart grid environments ", *International Wireless Communications and Mobile Computing Conference (IWCMC 2015)*, pp. 1211-1216, 2015.
- [17] Z. Huang, L. Y. Por, T. F. Ang, M. H. Anisi, and M. S. Adam , "Improving the Accuracy Rate of Link Quality Estimation Using Fuzzy Logic in Mobile Wireless Sensor Network ", *Advances in Fuzzy Systems*, vol. 2019, Article ID 3478027, 13 pages, <https://doi.org/10.1155/2019/3478027>, 2019.
- [18] Z. Q. Guo, Q. Wang, M. H. Li, and J. He, " Fuzzy logic based multidimensional link quality estimation for multi-hop wireless sensor networks ", *IEEE Sensors Journal*, vol.13, no.10, pp. 3605-3615, 2013.
- [19] M. Senel, K. Chintalapudi, D. Lal, A. Keshavarzian, and E. J. Coyle, "A kalman filter based link quality estimation scheme for wireless sensor networks", *IEEE Global Telecommunications Conference (GLOBECOM07)*, 2007.
- [20] M.J. Tanakian, M. Rezaei, and F.Mohanna , "Digital video stabilizer by adaptive fuzzy filtering", *Journal Image Video Proc* , vol.21, doi:10.1186/1687-5281-2012-21,2012
- [21] L. A. Zadeh, "Fuzzy sets," *Information and Control*, Vol. 8, 1965, pp. 338-353.
- [22] A. Abdel-Aleem, M.A.El-sharief, M.A. Hassan, and M.G. El-sebaie, "Implementation of fuzzy and adaptive neuro-fuzzy inference systems in optimization of production inventory problem", *Applied Mathematics & Information Sciences*, Vo.11, No.1, 2017, pp.289–298

- [23] A. Sadollah, *Fuzzy Logic Based in Optimization Methods and Control Systems and Its Applications*, London: IntechOpen, 2018, DOI: 10.5772/intechopen.73112.
- [24] S. Woo and H. Kim, "An empirical interference modeling for link reliability assessment in wireless networks", *IEEE/ACM Transactions on Networking*, Vol. 21, Issue 1, 2013, pp. 272 – 285.
- [25] P. Millan, C. Molina, E. Medina, D. Vega, R. Meseguer, B. Braem, and C. Blondia, "Tracking and predicting link quality in wireless community networks", In *IEEE 10th international conference on wireless and mobile computing, networking and communications (WiMob)*(pp. 239–244).

**Mohamad Javad Tanakian** received the B.S. and M.S. degrees in Electrical engineering in 2007 and 2011, from University of Sistan and Baluchestan(USB), Zahedan, Iran, respectively. He is now Ph.D. candidate in telecommunication engineering in the USB. He has been working as a Fiber optic expert at Sistan and Baluchestan Regional Electric Company (SBREC) since 2014. His research interests are in the area of wireless communication, smart grid communication, signal processing, image and video processing and fuzzy systems

**Mehri Mehrjoo** received the B.A.Sc. and the M.A.Sc. degrees from Ferdowsi University, Mashhad, Iran, and Ph.D. from the University of Waterloo, Waterloo, Canada in 1993, 1996, and 2008, respectively. From 2008 to 2009, she has been a post-doctoral fellow at the University of Waterloo. She is an IEEE senior member. Currently, she is an associate professor in the Department of Telecommunications, University of Sistan and Baluchestan, Zahedan, Iran. Her research interests are in the areas of resource allocation and performance analysis of broadband wireless protocols.

# SSIM-Based Fuzzy Video Rate Controller for Variable Bit Rate Applications of Scalable HEVC

Farhad Raufmehr

Department of Communications Engineering, University of Sistan and Baluchestan, Zahedan, Iran  
farhad.raufmehr@gmail.com

Mehdi Rezaei\*

Department of Communications Engineering, University of Sistan and Baluchestan, Zahedan, Iran  
mehdi.rezaei@ece.usb.ac.ir

Received: 24/May/2019

Revised: 24/Aug/2019

Accepted: 20/Nov/2019

## Abstract

Scalable High Efficiency Video Coding (SHVC) is the scalable extension of the latest video coding standard H.265/HEVC. Video rate control algorithm is out of the scope of video coding standards. Appropriate rate control algorithms are designed for various applications to overcome practical constraints such as bandwidth and buffering constraints. In most of the scalable video applications, such as video on demand (VoD) and broadcasting applications, encoded bitstreams with variable bit rates are preferred to bitstreams with constant bit rates. In variable bit rate (VBR) applications, the tolerable delay is relatively high. Therefore, we utilize a larger buffer to allow more variations in bitrate to provide smooth and high visual quality of output video. In this paper, we propose a fuzzy video rate controller appropriate for VBR applications of SHVC. A fuzzy controller is used for each layer of scalable video to minimize the fluctuation of QP at the frame level while the buffering constraint is obeyed for any number of layers received by a decoder. The proposed rate controller utilizes the well-known structural similarity index (SSIM) as a quality metric to increase the visual quality of the output video. The proposed rate control algorithm is implemented in HEVC reference software and comprehensive experiments are executed to tune the fuzzy controllers and also to evaluate the performance of the algorithm. Experimental results show a high performance for the proposed algorithm in terms of rate control, visual quality, and rate-distortion performance.

**Keywords:** Fuzzy Control, Quality, Rate, Scalable high-efficiency video coding (SHVC), SSIM, Variable bit rate (VBR).

## 1- Introduction

Multimedia and video technology improvements resulted in a new video coding standard which is called High Efficiency Video Coding (H.265/HEVC). The first version of HEVC was completed in January 2013. This new video coding standard has about 50% bit-rate reduction compared to the previous standard (H.264/AVC) [1, 2]. But video transmission over various networks usually faces many challenges such as diverse end-users and different connections quality. The solution to this problem is scalable video coding (SVC) [3]. So the need for a scalable video coding standard motivated the joint collaborative team on video coding (JCT-VC) to propose a new scalable extension in the second version of HEVC which was published in January 2015 [4]. This newly published version also includes the 3D/Multi-view and Range extensions [5]. The scalable extension of HEVC which name is Scalable High-efficiency Video Coding (SHVC) not only supports the conventional scalable features such as temporal, spatial, and quality scalabilities but also supports new scalability features such as hybrid

codec, bit depth, and color gamut. This new scalable extension was proposed by only modifying the first version of HEVC at high-level syntax and the encoding core is unchanged [4].

The available bandwidth and buffering constraints are other challenges in video transmission. Rate control algorithms (RCA) are utilized to solve this problem. According to the tolerable delay, video transmission applications are divided into constant bit rate (CBR) and variable bit rate (VBR) applications. In CBR applications such as conversational applications, the end-to-end delay is crucial, therefore; a CBR RCA with a relatively small buffer is required. The CBR rate control algorithms try to control the short-term average bit rate strictly in order to prevent the small buffer from overflow and underflow. Strict rate control means high fluctuations in the quantization parameter and thereafter a low-level perceptual quality. In VBR applications such as video streaming and broadcasting a relatively higher delay is tolerable so a larger buffer and a VBR RCA can be used in which a loose control over the long-term average bit rate is imposed. The initial buffering delay in VBR applications is higher than CBR applications [6]. For many VBR

\* Corresponding Author:



applications, a higher bit rate is usually used to guarantee acceptable video quality. Also, the rate-distortion performance of VBR is much better than CBR [7]. Considering these practical applications, three encoding configurations, including All Intra, Random Access, and Low Delay configurations, are introduced for HEVC. In these configurations, selected Level and Tier determine the buffering capabilities of the decoder. The standard coded picture buffer (CPB) size, decoded picture buffer (DPB) size, and maximum bit rate are sample parameters related to the buffering capabilities. According to these, the random access (RA) configuration is appropriate for VBR applications. The compression performance of the RA configuration is relatively high, but it includes structural coding delay [6].

### 1-1- Related Works

There are several RCAs including CBR and VBR algorithms proposed for HEVC and SHVC which are reviewed here. Li et al. proposed a rate- $\lambda$  (Lagrange multiplier) based RCA for the first version of HEVC which enables the encoder to select among coding parameters to achieve the target rate as well as to minimize distortion [8]. Marzuki et al. take the advantages of the rate- $\lambda$  model to propose a tile-level rate controller for HEVC on the tile parallelization case [9]. Wang et al. proposed a distortion model, a rate model, and a mixed distribution model for residual signal and developed a  $\rho$ -domain RCA [10]. Choi et al. proposed a precise RCA based on a rate-quantization model [11]. Seo et al. considered the bandwidth and buffering constraints and proposed a video quality controller based on a distortion-quantization model and a rate-quantization model [12]. Lee et al. developed a rate-quantization model for each Coding Unit (CU) depth based on texture and non-texture models and proposed a frame-level rate control scheme [13]. Wang et al. improved the performance of the conventional rate- $\lambda$  model by proposing a gradient-based rate- $\lambda$  model for intra-frame rate control [14].

All the algorithms mentioned above operate as CBR algorithms which are not suitable for the VBR applications. Lopez et al. proposed a VBR control algorithm based on long-term and short-term sliding windows [15]. Inspiring from the idea proposed by Rezaei et al. in [16] as the semi-fuzzy rate controller, Fani et al. and Kamran et al. proposed fuzzy rate controllers for GOP-level and frame-level, respectively [17, 18]. Fani et al. utilized the proportional, integral and derivative components of the GOP bit error as the inputs of a fuzzy system to propose a novel PID-fuzzy video rate controller [19]. Although these RCAs are targeted for VBR applications they have been designed for the non-scalable version of HEVC and it is essential to design appropriate ones for the SHVC.

According to [20, 21], Li et al. extended the idea of the rate- $\lambda$  model-based rate controller of HEVC to SHVC. Biatek et al. proposed an adaptive rate controller for SHVC which dynamically adjusts the bit rate ratio between the base layer (BL) and an enhancement layer (EL) to optimize the coding performance under the global bitrate constraint [22]. However, these two algorithms fall into the CBR category too. Considering high delay applications of SHVC we proposed a fuzzy-logic-based scalable video rate controller, which falls into VBR [23].

The distortion models which are used in most of the previously discussed RCAs are based on the error-sensitive metrics such as mean square error (MSE) and peak signal to noise ratio (PSNR). The simplicity of calculation is the main popularity reason for these metrics. These metrics are purely mathematical and they do not consider the characteristics of the human visual system (HVS) so they have less correlation with HVS. Since the video quality is ultimately judged by human eyes it is better to utilize metrics with more adaptation to characteristics of HVS. The HVS reacts quickly to the structural information in the field of viewing, so it is better to use structural similarity-based metrics such as the well-known structural similarity index (SSIM) which exploits structural information to estimate the quality of a compressed video [24]. There are several RCAs that utilized the SSIM as a distortion metric in order to increase perceptual video quality. Zhao et al. incorporated the SSIM into the HEVC rate-distortion optimization (RDO) framework and a CU-level RCA [25]. Zeng et al. utilized the SSIM in the video quality assessment and bit allocation scheme and improved the rate- $\lambda$  model to control the bit rate [26]. Gao et al. considered the bit allocation as a resource allocation problem and by defining an SSIM-based utility function proposed a Nash bargaining solution [27]. These algorithms are CBR but Wang et al. proposed an SSIM-motivated two-pass VBR rate controller for HEVC in which collected information during the first pass is used for bit allocation and bit rate control in the second pass [28]. Zupancic et al. proposed a two-pass rate controller that occupies a fast encoder in the first pass to collect necessary information for rate allocation and model parameter estimation and use the collected information in the second pass [29].

In this paper, inspiring from the presented algorithm in [23], we propose an SSIM-based fuzzy video rate controller for VBR applications of the SHVC which is able to improve the SSIM-based quality measures for compressed video. It controls the bit rate of several temporal, spatial, and quality layers at the same time. The proposed algorithm tries to achieve long-term average target rates for the SHVC video layers by smooth changing of QP used for encoding each layer. In the proposed algorithm, a fuzzy controller is used for each layer to minimize the changes of QP at the frame level while the

buffering constraints are obeyed. Moreover, an SSIM-based quality controller in cooperation with the fuzzy controller is used in each layer in order to improve and smooth the SSIM metric over encoded video frames. The SSIM-based quality controller suppresses the unnecessary QP fluctuation allowing more fluctuation in bit rate and buffer occupancy. The proposed RCA provides encoded videos with smooth and high visual quality. All conventional rate controllers use a target bit rate as the main reference point in the control process and therefore the bit rate is pushed toward a constant value which is unwelcoming for the VBR applications. In our proposed algorithm, in fact, the QP and SSIM are used as references and the attempt is to prevent unnecessary changes of QP and SSIM. This leads to controlled variations in bit rate and smooth visual quality of the compressed video. Using these references enables the rate controller to operate in a wide rate-distortion (R-D) range, which is the main characteristic of VBR rate controllers.

The rest of this paper is organized as follows. First, the details of our proposed RCA are explained in Section 2. Then, some experimental results are reported in Section 3. Finally, conclusions are given in Section 4.

## 2- Proposed Rate Control Algorithm

The block diagram of our proposed RCA is shown in Fig. 1. A fuzzy rate controller, an SSIM-based quality controller, a number of virtual buffers, and a number of multiplexers are the main parts of the diagram. The algorithm operates at GOP (groups of pictures) level. It computes a base QP for each GOP and the well-known QP cascading technique is used to calculate a QP for each frame in the GOP. In the base QP calculation process, we consider the correlation between coding complexities of consequent GOPs in a scene. So the coding complexity of previous GOP is used as an estimate for that of the current GOP. Therefore, the base QP of previous encoded GOP is used as an estimate for that of the current GOP and then the fuzzy rate controller and the SSIM-based quality controller adjust the base QP by:

$$BaseQP_b^d = BaseQP_{b-1}^d + \Delta QPF_b^d + \Delta QPQ_b^d, \quad (1)$$

where  $BaseQP_b^d$  is the calculated QP for the  $b^{th}$  (current) GOP.  $\Delta QPF_b^d$  and  $\Delta QPQ_b^d$  denote the base QP changes calculated by the fuzzy rate controller and the SSIM-based quality controller, respectively.  $b$  and  $d$  denote the indices of GOP and layer, respectively. In fact, we can say the base QP of current GOP, consist of the delayed version of the base QP used for previous GOP plus the base QP changes that are calculated by the fuzzy rate controller and the SSIM-based quality controller. The details of the proposed RCA are discussed in the following subsections.

### 2-1- Virtual Buffer

As shown in Fig. 1 we employed a virtual buffer denoted by ( $Buffer^d$ ) for each layer in order to simulate the buffering process at the decoder side. The buffer size ( $BS^d$ ) and the target rate ( $TR^d$ ) are determined by the users according to the bandwidth, buffering, and delay constraints.

Inspiring from the fact that the sub-stream of each layer is multiplexed with those of other layers in the scalable video encoder in order to produce the scalable bitstream, so in the buffering process, the consumed bits of each layer are aggregated with those of lower layers. This is done by layer multiplexers ( $MUX^d$ ) which is used to simulate the encoder multiplexing process as:

$$MB_a^d = \sum_{j=0}^d B_a^j, \quad (2)$$

where  $B_a^j$  denotes the consumed bits for the  $a^{th}$  frame in the  $j^{th}$  layer and  $MB_a^d$  is the multiplexed consumed bits for the  $a^{th}$  frame in the  $d^{th}$  layer. Then, the output of the  $d^{th}$  multiplexer is used to update the occupancy of the  $d^{th}$  virtual buffer after encoding  $a^{th}$  frame according to (3):

$$BO_a^d = BO_{a-1}^d - MB_a^d + \frac{1}{F} \sum_{j=0}^d TR^j, \quad (3)$$

Here,  $BO_a^d$  denotes the buffer occupancy of the  $d^{th}$  layer after encoding the  $a^{th}$  frame. Also  $TR^j$  is the target rate of the  $j^{th}$  layer and  $F$  stands for the frame rate. It is notable that we assume %60 of the virtual buffer size ( $BS^d$ ) is initially occupied by initial buffering i.e.

$$BO_0^d = 0.6 \times BS^d, \quad (4)$$

### 2-2- Fuzzy Rate Controller

As a conventional fuzzy controller, our fuzzy rate controller contains a fuzzifier, a defuzzifier, a fuzzy interface engine, and a fuzzy rule base. The fuzzifier maps the crisp inputs to the input fuzzy sets. The fuzzy interface engine maps the input fuzzy set to the output fuzzy set according to the fuzzy rule base and the defuzzifier maps the output fuzzy set into a crisp output. The fuzzy rule base is presented with linguistic variables to appropriately link daily conversations to a mathematical framework [30]. We choose the fuzzy logic controller because many non-linear relations that exist in video rate control can be easily included in the fuzzy rules and membership functions (MSF).



Table 2 Desired Central Values of the Fuzzy System Output

	VH	6	6	6	5	4	3	2	1	0	
	H	6	6	5	4	3	2	1	0	-1	
	MH	6	5	4	3	2	1	0	-1	-2	
$\bar{z}_p^d$	M	5	4	3	2	1	0	-1	-2	-3	
	ML	4	3	2	1	0	-1	-2	-3	-4	
	L	3	2	1	0	-1	-2	-3	-4	-5	
	VL	2	1	0	-1	-2	-3	-4	-5	-6	
	UL										
		EL	VL	L	ML	M	MH	H	VH		
					$x_1^d$						

In MSFs of input1, the sets in the middle ranges such as ML and M cover a wider area than the others since in the middle ranges the buffer occupancy is far from critical conditions and so the QP is kept unchanged or change slowly. On the other hand, where the buffer status is critical, the sets close to zero or one cover narrower ranges to allow faster changes of QP. In other words, while the normalized buffer occupancy is about 0.6 and the normalized consumed bits is close 1 so the inputs are close to the ideal condition and there is no need to change the QP. As the result, the desired central value corresponding to the such area is set equal to 0 and whatever we select the MSFs in such regions wider, so the QP will be kept unchanged in wider region. However, when the inputs are far from ideal condition and close to the critical regions such as normalized buffer occupancy close to 1 and 0, the QP should be changed abruptly to prevent buffer overflow and underflow. Since the abrupt QP change, increases the unwelcomed quality fluctuation, the MSFs should be narrow in order to limit abrupt QP fluctuation to the critical regions.

Finally, by using a singleton fuzzifier, a product interface engine, and a center average defuzzifier the output of the fuzzy system is computed by (8):

$$f^d(x_1^d, x_2^d) = \frac{\sum_{i_1=1}^{N_1} \sum_{i_2=1}^{N_2} \bar{y}^{i_1 i_2} \mu_{A_1^{i_1}}(x_1^d) \mu_{A_2^{i_2}}(x_2^d)}{\sum_{i_1=1}^{N_1} \sum_{i_2=1}^{N_2} \mu_{A_1^{i_1}}(x_1^d) \mu_{A_2^{i_2}}(x_2^d)}, \quad (8)$$

where  $f^d$  denotes the output of the fuzzy system for the  $d^{th}$  layer.  $\{A_1^1, A_1^2, A_1^3, \dots, A_1^{N_1}\}_{i=1,2}$  stands for the input fuzzy set and  $\bar{y}^{i_1 i_2}$  stands for the central desired value.  $N_1$  and  $N_2$  are the number of fuzzy sets for input  $x_1^d$  and  $x_2^d$  respectively. It is emphatic that all the aspects of the fuzzy controller are considered based on the expert experiences and experiments execution without performing an optimization process. However, the experimental results confirm that the performance of our proposed RCA is better than the others. Readers are referenced to [30] for more detailed information about fuzzy design and derivations. The output of the fuzzy system will be passed

through a content-adaptive gain ( $G_f^d$ ) which can be tuned in the range of (0.5 ~ 1) in order to adjust the control intensity accordingly to the video content as:

$$\Delta QPF_b^d = G_f^d \times f^d(x_1^d, x_2^d), \quad (9)$$

A higher gain is suitable for a video sequence with a lot of heterogeneous scenes and vice versa. The output of the fuzzy controller is used for computing the base QP as presented in equation (1).

### 2-3- SSIM-Based Quality Controller

The distortion model used in the HEVC framework is based on the error-sensitive metrics such as PSNR and MSE that are full-reference (FR) objective metrics. The main popularity reason for them is the calculation simplicity. They do not take the human visual system (HVS) characteristics into account while the output video quality is ultimately judged by human eyes. Researchers show that HVS is very sensitive to the structural information in the field of viewing. Therefore, it is better to use structural similarity-based metrics which have more correlation with HVS and estimate the output video quality more efficiently. SSIM is a structural similarity-based metric which attempts to extract structural information to evaluate the video quality.

To define SSIM, let  $\alpha$  and  $\beta$  be the original signal and distorted signal respectively.  $\mu_\alpha$  and  $\mu_\beta$  denote the mean of  $\alpha$  and  $\beta$ , respectively which estimate the luminance.  $\sigma_\alpha$  and  $\sigma_\beta$  stand for the variance that estimate the contrast, and  $\sigma_{\alpha\beta}$  is the covariance of  $\alpha$  and  $\beta$  which measures the non-linear similarity of  $\alpha$  and  $\beta$ . By utilizing these parameters, the luminance ( $l$ ), the contrast ( $c$ ) and structure ( $s$ ) comparison measures are defined as follow:

$$l(\alpha, \beta) = \frac{2\mu_\alpha \mu_\beta}{\mu_\alpha^2 + \mu_\beta^2},$$

$$c(\alpha, \beta) = \frac{2\sigma_\alpha \sigma_\beta}{\sigma_\alpha^2 + \sigma_\beta^2},$$

$$s(\alpha, \beta) = \frac{\sigma_{\alpha\beta}}{\sigma_\alpha \sigma_\beta}, \quad (10)$$

Then, the structural similarity measure is yielded (11) by combining these measures:

$$S(\alpha, \beta) = l(\alpha, \beta) c(\alpha, \beta) s(\alpha, \beta) = \frac{4\mu_\alpha \mu_\beta \sigma_{\alpha\beta}}{(\mu_\alpha^2 + \mu_\beta^2)(\sigma_\alpha^2 + \sigma_\beta^2)}, \quad (11)$$

The equation (11) is unstable and so it is modified to a new measure named SSIM as:

$$SSIM(\alpha, \beta) = \frac{(2\mu_\alpha \mu_\beta + C_1)(2\sigma_{\alpha\beta} + C_2)}{(\mu_\alpha^2 + \mu_\beta^2 + C_1)(\sigma_\alpha^2 + \sigma_\beta^2 + C_2)}, \quad (12)$$

where  $C_1$  and  $C_2$  are given by (13):

$$C_1 = (k_1 L)^2, \quad C_2 = (k_2 L)^2, \quad (13)$$

where  $L$  is the dynamic range of pixel values set to 255 for 8-bit videos.  $k_1$  and  $k_2$  are two constants set to 0.01 and 0.03, respectively. Readers are referenced to [24, 31, 32] for more details.

In this paper, we take the advantages of SSIM as a quality metric and propose an SSIM-based quality controller to improve the performance of our algorithm. Our quality controller uses the SSIM of the compressed video in each layer as a feedback signal and calculates a QP change ( $\Delta QPQ_b^d$ ) in the range of (-2 ~ 2) for each layer as output. The relation between the input and output of our quality controller is represented in (14):

$$\Delta QPQ_b^d = G_q^d \times \overline{QP^d} \left( \overline{SSIM^d} - SSIM_{b-1}^d \right), \quad (14)$$

where the  $\overline{QP^d}$  and  $\overline{SSIM^d}$  denote the average QP and the average SSIM, respectively for all previously encoded frames at the same layer.  $SSIM_{b-1}^d$  stand for the average SSIM of previously encoded GOP.  $G_q^d$  is a constant gain which can be used to adjust the control intensity. The output of the quality controller is used for computing the base QP as presented in equation (1).

### 3- Experimental Results

To evaluate the performance of proposed RCA, we implemented our proposed algorithm on the SHVC standard reference software SHM-12.1 [33] and executed a set of experiments. The random access scalable configuration with three layers including a base layer, a spatial 2.x layer, and an SNR layer is used for the experiments. The reason for using a spatial and an SNR layer as enhancement layers is to show the performance of our algorithm on both types of scalable layers. Each enhancement layer uses only the previous layer in interlayer processing. In order to configure the RCA, the size of each buffer is chosen equal to 1.5 seconds buffering of a bitstream with the aggregated target bit rate. Moreover, the gains of the fuzzy rate controller and the quality controller were set to 0.65 and 0.7, respectively.

We used a set of well-known sequences such as Keiba, RaceHorses, BQMall, BasketballDrill, PartyScene, Kimono, and ParkScene in our experiments. The proposed RCA is targeted for long-term rate-controlling, so we concatenated short test sequences to make longer sequences, suitable for our experiments. KR, BP, and KP are the abbreviations for the name of concatenated sequences Keiba to RaceHorses, BasketballDrill to PartyScene, and Kimono to ParkScene, respectively.

In video encoding with a constant QP (CQP), there is no control over the bit rate and therefore, there is no guarantee for the buffer constraint to be obeyed especially for a long time. However, CQP encoding provides smooth and higher visual quality for compressed video. On the other hand, the  $\lambda$ -domain RCA implemented in the reference software is supposed to produce a constant bit rate suitable for low-delay applications. From the operating region point of view, a high-delay RCA should operate in a region between low-delay algorithms and CQP case. Hence, we selected the  $\lambda$ -domain RCA and CQP cases in order to compare our proposed algorithm with them from the rate control and video quality points of view. These two algorithms are compared with the proposed algorithm in terms of mean QP, mean PSNR and mean SSIM. PSNR is a signal fidelity metric that measures the correlation between the original video and the encoded one. The higher PSNR means higher quality. As discussed in the previous section, SSIM measures the structural similarity between the original video and the processed one and reports the similarity value in the range of (0~1). The closer SSIM to 1 means higher quality. Moreover, we introduced a metric namely Mean Absolute Gradient (MAG) as a fluctuation metric to compare the algorithms in terms of fluctuations on QP, PSNR, and SSIM. The MAG on the variable  $\theta^d$  is defined as:

$$MAG(\theta^d) = \frac{1}{M-1} \sum_{a=0}^{M-1} |\theta_{a+1}^d - \theta_a^d|, \quad (15)$$

where the variable  $\theta^d$  can be substituted by QP, PSNR or SSIM in order to compute the MAG of these metrics.  $M$  denotes the number of encoded frames and stands for the index of each frame in display order. The MAG of PSNR and SSIM measure that how much the quality and structural similarity of consequent frames are correlated. Small MAG of PSNR and SSIM means that the quality of the output video has smooth fluctuation and a more pleasant video display is provided for the user.

To evaluate the performance of our RCA from the buffering constraints point of view, the encoded sequences are compared in term of minimum initial buffering delay which is formulated in (16):

$$Delay^d = \frac{0.6 \times (BO_{\max}^d - BO_{\min}^d)}{\sum_{j=0}^d TR^j}, \quad (16)$$

where a higher delay means more variations in bitrate. A high delay value can be interpreted as overflow or underflow or both of them. However, low delay values cannot be necessarily interpreted as perfect control. In other words, (16) measures the time that should be passed until 60% of the buffer space be filled with the incoming bits. For perfect control, the following constraints must be obeyed for all GOPs in all layers:

$$BO_{\max}^d \leq BS^d \quad \text{AND} \quad BO_{\min}^d \geq 0, \quad (17)$$

From the rate control point of view, it is notable that in all experiments, the proposed RCA completely obeyed the buffering constraints with neither buffer overflow nor underflow and successfully achieved the target rate. On the other hand, the virtual buffer simulated for CQP and  $\lambda$ -domain RCAs has overflow or underflow in several cases. The test sequences were encoded by the algorithms in four operation points according to the SHM common test conditions [34] and for the lack of space, only a part of numerical experimental results are represented in Table 3. According to Table 3 by averaging PSNR mean over all layers of the test sequences, the values of 37.38, 37.46, and 37.29 are resulted by the CQP, the proposed (S-F), and the  $\lambda$ -domain (LAD) RCAs, respectively.

Table 3(a) Comparison Simulation Results of S-F with CQP and LAD

Sequences	Layer ID	RCA	PSNR (dB)		SSIM	
			Mean	MAG	Mean	MAG
KR	BL	CQP	34.74	1.47	0.932	0.013
		S-F	34.91	1.44	0.926	0.013
		LAD	34.75	1.81	0.921	0.016
		LAD	34.75	1.81	0.921	0.016
	EL1	CQP	35.32	1.32	0.928	0.011
		S-F	35.44	1.24	0.922	0.011
		LAD	35.35	1.81	0.918	0.016
		LAD	35.35	1.81	0.918	0.016
	EL2	CQP	37.33	1.68	0.948	0.011
		S-F	37.42	1.55	0.945	0.011
		LAD	37.27	2.20	0.941	0.015
		LAD	37.27	2.20	0.941	0.015
BQMall	BL	CQP	32.64	0.58	0.928	0.006
		S-F	32.94	0.60	0.929	0.006
		LAD	32.80	0.64	0.925	0.007
		LAD	32.80	0.64	0.925	0.007
	EL1	CQP	33.65	0.47	0.915	0.005
		S-F	33.82	0.48	0.914	0.005
		LAD	33.50	0.58	0.907	0.007
		LAD	33.50	0.58	0.907	0.007
	EL2	CQP	35.75	0.61	0.94	0.005
		S-F	35.87	0.61	0.939	0.005
		LAD	35.59	0.69	0.935	0.006
		LAD	35.59	0.69	0.935	0.006
BP	BL	CQP	39.75	1.44	0.980	0.004
		S-F	39.84	1.49	0.980	0.005
		LAD	39.67	1.87	0.979	0.006
		LAD	39.67	1.87	0.979	0.006
	EL1	CQP	39.38	1.20	0.970	0.004
		S-F	39.52	1.14	0.970	0.004
		LAD	39.23	1.75	0.967	0.007
		LAD	39.23	1.75	0.967	0.007
	EL2	CQP	41.80	1.70	0.981	0.004
		S-F	41.84	1.75	0.981	0.005
		LAD	41.62	2.48	0.979	0.007
		LAD	41.62	2.48	0.979	0.007
KP	BL	CQP	38.45	1.00	0.958	0.007
		S-F	38.46	0.97	0.957	0.007
		LAD	38.39	0.85	0.956	0.006
		LAD	38.39	0.85	0.956	0.006
	EL1	CQP	39.04	0.73	0.946	0.006
		S-F	38.91	0.73	0.944	0.006
		LAD	38.86	0.73	0.943	0.006
		LAD	38.86	0.73	0.943	0.006
	EL2	CQP	40.70	0.84	0.958	0.006
		S-F	40.57	0.82	0.957	0.006
		LAD	40.47	0.98	0.957	0.007
		LAD	40.47	0.98	0.957	0.007
Total-Average	CQP	37.38	1.09	0.949	0.007	
	S-F	37.46	1.07	0.947	0.007	
	LAD	37.29	1.37	0.944	0.009	

The results show that our algorithm provides a higher video quality level than the CQP and the  $\lambda$ -domain algorithms in terms of PSNR. Also, by averaging the MAG of PSNR over the tested sequences the values of 1.09, 1.07, and 1.37 are resulted by the CQP, the S-F, and the  $\lambda$ -domain RCAs, respectively. According to these results, our algorithm has provided less fluctuation in PSNR than the anchors and so it provides more constant visual quality. Moreover, by averaging the SSIM values over the test sequences, the values of 0.949, 0.947, and 0.944 are resulted by the CQP, the S-F, and the  $\lambda$ -domain RCAs, respectively. That means the performance of the

Table 3(b) Comparison Simulation Results of S-F with CQP and LAD

Sequences	Layer ID	RCA	QP		Delay (Sec)	Average Bit-Rate (kbps)
			Mean	MAG		
KR	BL	CQP	33.10	1.83	1.61	297.81
		S-F	32.86	1.83	0.46	294.89
		LAD	36.05	7.27	0.53	298.02
		LAD	36.05	7.27	0.53	298.02
	EL1	CQP	33.10	1.83	2.12	876.02
		S-F	32.79	1.82	0.59	858.18
		LAD	35.80	7.38	0.65	876.44
		LAD	35.80	7.38	0.65	876.44
	EL2	CQP	29.10	1.82	2.24	972.15
		S-F	28.76	1.81	0.47	962.25
		LAD	31.43	7.31	0.63	973.07
		LAD	31.43	7.31	0.63	973.07
BQMall	BL	CQP	37.10	1.84	0.82	252.94
		S-F	36.91	1.85	0.37	257.38
		LAD	39.77	5.11	0.24	253.66
		LAD	39.77	5.11	0.24	253.66
	EL1	CQP	37.10	1.84	0.91	570.83
		S-F	36.87	1.85	0.44	577.78
		LAD	39.19	5.52	0.27	571.75
		LAD	39.19	5.52	0.27	571.75
	EL2	CQP	33.10	1.83	0.88	704.08
		S-F	32.77	1.84	0.35	705.15
		LAD	35.35	5.88	0.25	704.71
		LAD	35.35	5.88	0.25	704.71
BP	BL	CQP	25.10	1.80	1.30	1461.06
		S-F	24.91	1.80	0.39	1472.64
		LAD	26.79	6.18	0.36	1461.01
		LAD	26.79	6.18	0.36	1461.01
	EL1	CQP	25.10	1.80	2.02	4804.95
		S-F	24.62	1.80	0.32	4810.19
		LAD	26.68	6.51	0.63	4809.08
		LAD	26.68	6.51	0.63	4809.08
	EL2	CQP	21.10	1.80	2.01	5358.12
		S-F	20.72	1.80	0.34	5288.33
		LAD	22.63	6.32	0.66	5359.77
		LAD	22.63	6.32	0.66	5359.77
KP	BL	CQP	29.10	1.83	1.22	1089.25
		S-F	29.08	1.81	0.54	1073.22
		LAD	30.54	3.95	0.79	1089.88
		LAD	30.54	3.95	0.79	1089.88
	EL1	CQP	29.10	1.83	1.20	2779.86
		S-F	29.19	1.81	0.60	2723.22
		LAD	30.60	4.31	0.83	2782.28
		LAD	30.60	4.31	0.83	2782.28
	EL2	CQP	25.10	1.82	1.36	3599.17
		S-F	25.10	1.8	0.54	3585.63
		LAD	26.75	4.90	0.87	3601.02
		LAD	26.75	4.90	0.87	3601.02
Total-Average	CQP	29.77	1.82	1.47	1897.19	
	S-F	29.55	1.82	0.45	1884.07	
	LAD	31.80	5.89	0.56	1898.39	

proposed algorithm in terms of the visual quality of the compressed video is between those of the CQP and the  $\lambda$ -domain RCA as expected. The CQP algorithm uses a constant QP over the whole sequence while our proposed algorithm needs to vary the QP in order to control the bit rate and buffer state so the CQP outperforms the proposed algorithm. However, smooth control of QP by the proposed algorithm provides high visual quality for compressed video. Furthermore, by averaging the MAG of SSIM the values of 0.007, 0.007, and 0.009 are resulted by the CQP, the S-F, and the  $\lambda$ -domain RCAs, respectively. That means a similar performance in terms of visual quality smoothness for the proposed and the CQP algorithms. According to the average values for QP mean (29.77, 29.55, 31.80) and QP MAG (1.82, 1.82, 5.89) in the table, the proposed RCA provided a QP mean lower than those of CQP and  $\lambda$ -domain RCAs while in term of QP MAG the proposed RCA performs similar to the CQP and much better than the  $\lambda$ -domain RCA. From the initial buffering delay point of view, the average values of 1.47, 0.45, and 0.56 are resulted by the CQP, the S-F, and the  $\lambda$ -domain RCAs, respectively. According to these results, the proposed algorithm provided a lower initial buffering delay than that of the CQP case and even lower than that of the  $\lambda$ -domain RCA as a constant bit rate algorithm. The point is that the  $\lambda$ -domain RCA failed to obey the buffer constraints in several cases in the experiments while and the buffering constraints are completely obeyed by the proposed RCA. For more investigations, we compared the proposed RCA with the CQP and the  $\lambda$ -domain RCAs in terms of the rate-distortion performance utilizing the Bjøntegaard metrics. The test sequences were encoded by the algorithms in four operation points according to the SHM common test conditions [34] and the Bjøntegaard Delta PSNR (BDPSNR) and Bjøntegaard Delta Bit Rate (BDBR) are computed between the proposed algorithm and the anchor algorithms. The BDPSNR measures the average PSNR between two rate-distortion curves. Positive BDPSNR means quality enhancement and negative BDPSNR means quality degradation. The BDBR measures the average bit rate between two rate-distortion curves. Negative BDBR means bit rate saving and positive BDBR is interpreted as bit rate wasting. The comparison results are presented in Table 4.

Also, the algorithms were compared in terms of the Bjøntegaard Delta SSIM (BDSSIM) and Bjøntegaard Delta Bit Rate (BDBR) and provided results are presented in Table 5. According to the average results reported in Table. 4, our RCA performs better than the anchor algorithms in terms of Rate-PSNR. Moreover, the reported average results in Table. 5 show that our RCA performs close to the CQP case (-0.0011, 3.8004) and better than the  $\lambda$ -domain RCA (0.0023, -6.0229) in terms of Rate-SSIM. As sample graphical results, Fig. 3-5 show the buffer

occupancy (BO), PSNR, SSIM and QP graphs on GOP-based for the three layers of the KP test sequence. As presented in the figures, the buffer occupancy (BO) graphs resulted by the CQP and  $\lambda$ -domain RCAs show several buffer overflows and underflows while for the proposed RCA the buffer has neither overflow nor underflow. As shown in the graphs, the proposed algorithm efficiently uses the buffer size in order to reduce the QP and perceptual quality fluctuations. Also, strong correlations between the graphs of proposed RCA and corresponding graphs of CQP can be seen in the figures that mean a high performance for the proposed RCA close to the CQP encoding.

Table 4 Rate-Distortion Performance Comparison between S-F, CQP, and  $\lambda$ -Domain in Term of PSNR

Sequence Name	Layer ID	PSNR			
		S-F Vs. CQP		S-F Vs. LAD	
		BDPSNR	BDBR	BDPSNR	BDBR
KR	BL	0.188	-3.782	0.215	-4.113
	EL1	0.087	-2.631	0.088	-2.693
	EL2	0.058	-1.695	0.145	-4.071
BQMall	BL	0.168	-3.048	0.087	-1.612
	EL1	0.019	-0.583	0.164	-4.415
	EL2	-0.061	1.995	0.117	-3.590
BP	BL	0.126	-2.486	0.061	-1.195
	EL1	0.212	-5.132	0.244	-5.873
	EL2	0.158	-3.444	0.194	-4.227
KP	BL	0.077	-1.833	0.145	-3.457
	EL1	-0.046	1.552	0.133	-4.518
	EL2	-0.105	4.387	0.092	-3.299
Total-Average		0.073	-1.392	0.140	-3.589

Table. 5 Rate-Distortion Performance Comparison between S-F, CQP, and  $\lambda$ -Domain in Term of SSIM

Sequence Name	Layer ID	SSIM			
		S-F Vs. CQP		S-F Vs. LAD	
		BDSSIM	BDBR	BDSSIM	BDBR
KR	BL	-0.0044	8.0664	0.0042	-7.2240
	EL1	-0.0038	11.1120	0.0027	-7.2125
	EL2	-0.0026	10.8361	0.0025	-8.8341
BQMall	BL	-0.0010	2.6213	0.0013	-3.5838
	EL1	-0.0011	3.6879	0.0023	-6.7775
	EL2	-0.0010	5.2484	0.0009	-4.5998
BP	BL	0.0001	-0.2108	0.0012	-2.8484
	EL1	0.0022	-4.4932	0.0045	-8.8604
	EL2	0.0012	-3.2940	0.0024	-6.5988
KP	BL	-0.0008	1.8461	0.0023	-4.8108
	EL1	-0.0012	3.6716	0.0025	-7.0958
	EL2	-0.0013	6.5138	0.0010	-3.8290
Total-Average		-0.0011	3.8004	0.0023	-6.0229

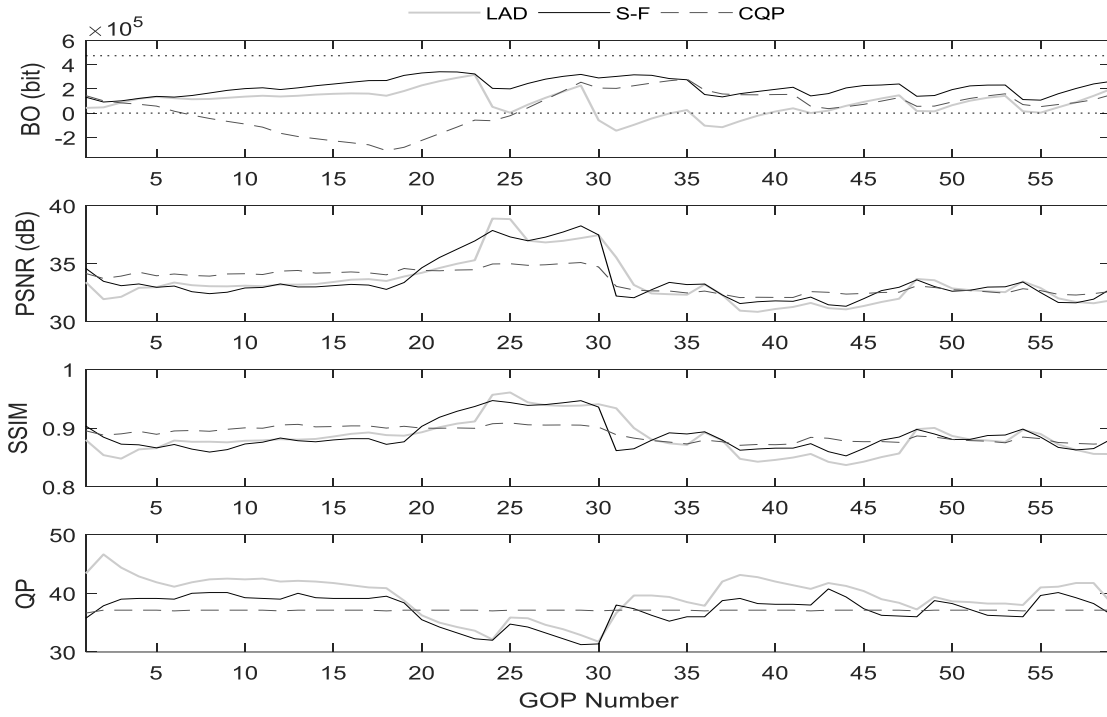


Fig. 3 Buffer Occupancy, PSNR, SSIM and QP Graphs of Base Layer

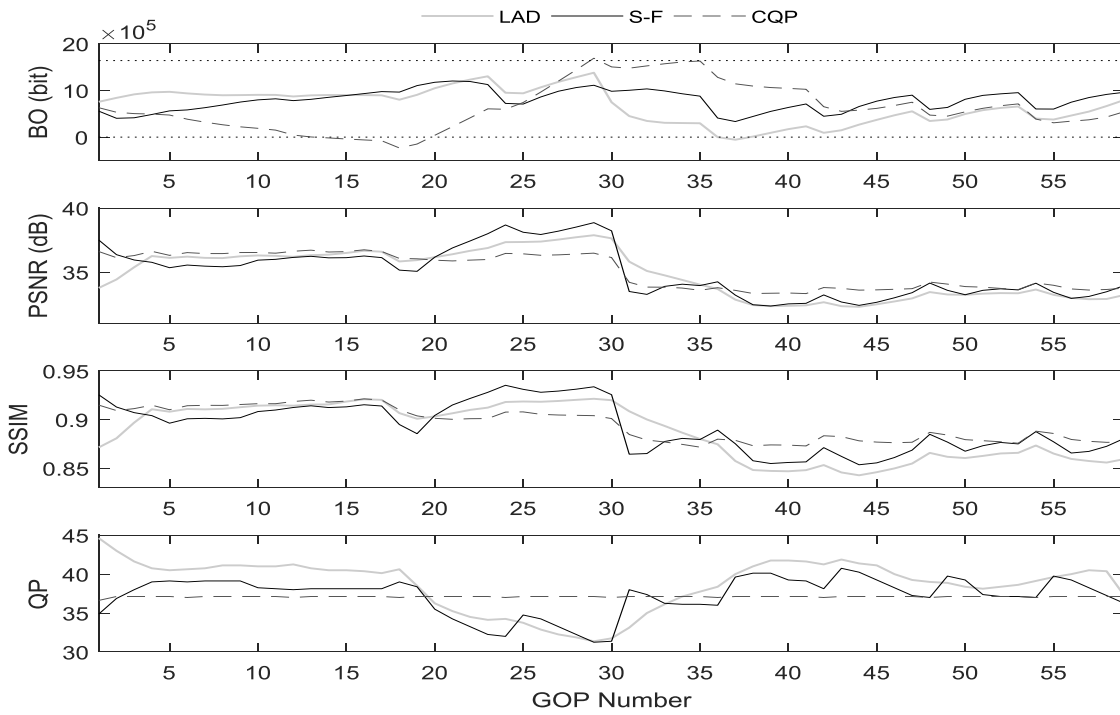


Fig. 4 Buffer Occupancy, PSNR, SSIM, and QP Graphs of Enhancement Layer1



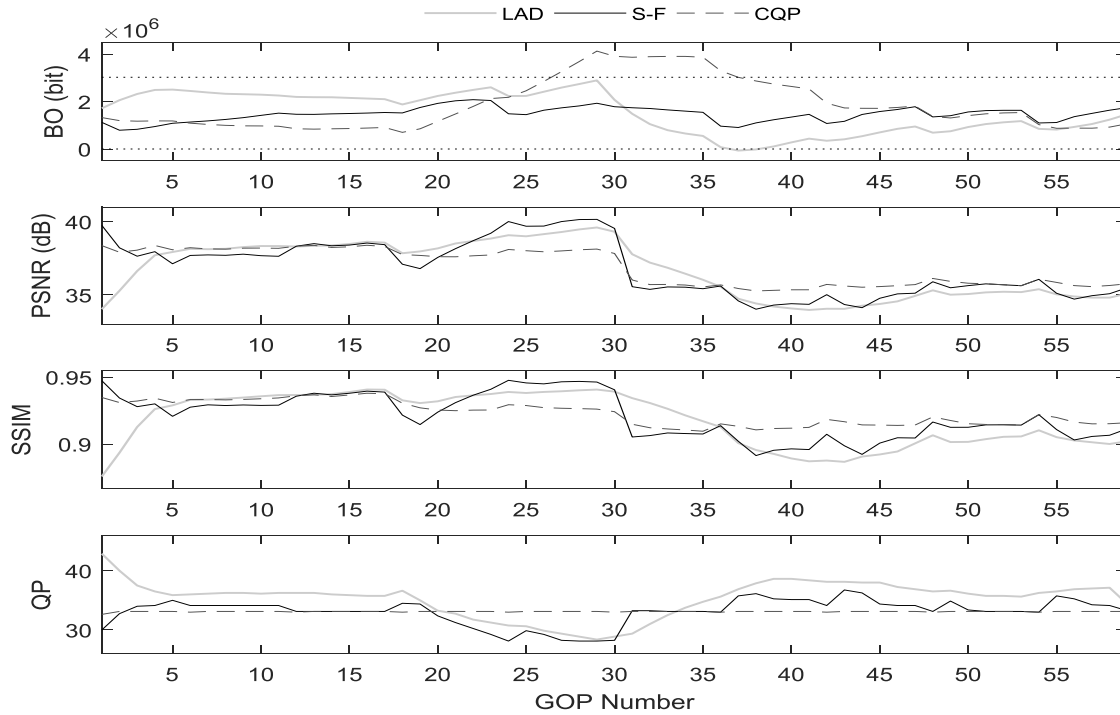


Fig. 5 Buffer Occupancy, PSNR, SSIM, and QP Graphs of Enhancement Layer2

## 4- Conclusions

In this paper, we proposed a video rate controller targeted for VBR applications of the scalable extension of HEVC standard, which is able to control the bit rate and buffer state in all types of scalable layers and it consists of a fuzzy rate controller and an SSIM-based quality controller. The fuzzy controller controls the bit rate while the attempt is to minimize the QP fluctuations and the quality controller smooths the SSIM quality metric over video frames.

Both controllers operate on the GOP level. The proposed rate control algorithm was implemented in the standard reference software and a comprehensive set of experiments was executed. According to the experimental results, the rate and buffering constraints are completely obeyed by the proposed algorithm. Also, from the rate-PSNR performance point of view, the proposed algorithm performs better than the CQP and  $\lambda$ -domain RCAs. Moreover, from the rate-SSIM performance point of view, the proposed algorithm performs better than the  $\lambda$ -domain RCA and close to the CQP case. Furthermore, from the QP and SSIM smoothness point of view the proposed RCA performs better than the  $\lambda$ -domain RCA and very close to the CQP encoding.

## References

- [1] V. Sze, M. Budagavi, and G.J. Sullivan, High Efficiency Video Coding (Hvc), Integrated Circuit and Systems, Algorithms and Architectures. Springer, 2014, pp. 1-375.
- [2] G.J. Sullivan, J. Ohm, W.-J. Han, and T. Wiegand, "Overview of the High Efficiency Video Coding (Hvc) Standard", IEEE Transactions on Circuits and Systems for Video Technology. vol. 22, no. 12. 2012, pp. 1649-1668.
- [3] H. Schwarz, D. Marpe, and T. Wiegand, "Overview of the Scalable Video Coding Extension of the H. 264/Avc Standard", IEEE Transactions on Circuits and Systems for Video Technology. vol. 17, no. 9. 2007, pp. 1103-1120.
- [4] J.M. Boyce, Y. Ye, J. Chen, and A.K. Ramasubramanian, "Overview of Shvc: Scalable Extensions of the High Efficiency Video Coding Standard", IEEE Transactions on Circuits and Systems for Video Technology. vol. 26, no. 1. 2016, pp. 20-34.
- [5] G.J. Sullivan, J.M. Boyce, Y. Chen, J.-R. Ohm, C.A. Segall, and A. Vetro, "Standardized Extensions of High Efficiency Video Coding (Hvc)", IEEE Journal of Selected Topics in Signal Processing. vol. 7, no. 6. 2013, pp. 1001-1016.
- [6] M. Wien, High Efficiency Video Coding, Coding Tools and specification. 2015.
- [7] H. Sun, T. Chiang, and X. Chen, Digital Video Transcoding for Transmission and Storage, CRC press, 2004.
- [8] B. Li, H. Li, L. Li, and J. Zhang, " $\lambda$ -Domain Rate Control Algorithm for High Efficiency Video Coding", IEEE

- transactions on image processing. vol. 23, no. 9. 2014, pp. 3841-3854.
- [9] MARZUKI, I., AHN, Y.-J. & SIM, D. 2017. Tile-level rate control for tile-parallelization HEVC encoders. *Journal of Real-Time Image Processing*, 1-19.
- [10] S. Wang, S. Ma, S. Wang, D. Zhao, and W. Gao, "Rate-Gop Based Rate Control for High Efficiency Video Coding", *IEEE Journal of Selected Topics in Signal Processing*. vol. 7, no. 6. 2013, pp. 1101-1111.
- [11] H. Choi, J. Yoo, J. Nam, D. Sim, and I.V. Bajic, "Pixel-Wise Unified Rate-Quantization Model for Multi-Level Rate Control", *IEEE Journal of Selected Topics in Signal Processing*. vol. 7, no. 6. 2013, pp. 1112-1123.
- [12] C.-W. Seo, J.-H. Moon, and J.-K. Han, "Rate Control for Consistent Objective Quality in High Efficiency Video Coding", *IEEE transactions on image processing*. vol. 22, no. 6. 2013, pp. 2442-2454.
- [13] B. Lee, M. Kim, and T.Q. Nguyen, "A Frame-Level Rate Control Scheme Based on Texture and Nontexture Rate Models for High Efficiency Video Coding", *IEEE Transactions on Circuits and Systems for Video Technology*. vol. 24, no. 3. 2014, pp. 465-479.
- [14] M. Wang, K.N. Ngan, and H. Li, "An Efficient Frame-Content Based Intra Frame Rate Control for High Efficiency Video Coding", *IEEE Signal Processing Letters*. vol. 22, no. 7. 2015, pp. 896-900.
- [15] M. de-Frutos-López, J.L. González-de-Suso, S. Sanz-Rodríguez, C. Peláez-Moreno, and F. Díaz-de-María, "Two-Level Sliding-Window Vbr Control Algorithm for Video on Demand Streaming", *Signal processing: Image communication*. vol. 36. 2015, pp. 1-13.
- [16] M. Rezaei, M.M. Hannuksela, and M. Gabbouj, "Semi-Fuzzy Rate Controller for Variable Bit Rate Video", *IEEE Transactions on Circuits and Systems for Video Technology*. vol. 18, no. 5. 2008, pp. 633-645.
- [17] D. Fani and M. Rezaei, "A Gop-Level Fuzzy Rate Control Algorithm for High-Delay Applications of Hevc", *Signal, Image and Video Processing*. vol. 10, no. 7. 2016, pp. 1183-1191.
- [18] R. Kamran, M. Rezaei, and D. Fani, "A Frame Level Fuzzy Video Rate Controller for Variable Bit Rate Applications of Hevc", *Journal of Intelligent & Fuzzy Systems*. vol. 30, no. 3. 2016, pp. 1367-1375.
- [19] FANI, D. & REZAEI, M. 2017. Novel PID-Fuzzy Video Rate Controller for High-Delay Applications of the HEVC Standard. *IEEE Transactions on Circuits and Systems for Video Technology*, 28, 1379-1389.
- [20] L. Li, B. Li, D. Liu, and H. Li, " $\lambda$ -Domain Rate Control Algorithm for Hevc Scalable Extension", *IEEE Transactions on Multimedia*. vol. 18, no. 10. 2016, pp. 2023-2039.
- [21] L. Li, B. Li, and H. Li, Rate Control by R-Lambda Model for Shvc. vol. JCTVC-M0037. 2013.
- [22] T. Biatek, W. Hamidouche, J.-F. Travers, and O. Deforges, "Adaptive Rate Control Algorithm for Shvc: Application to Hd/Uhd", in *Acoustics, Speech and Signal Processing (ICASSP)*, 2016 IEEE International Conference on. 2016, IEEE, p. 1382-1386.
- [23] F. Raufmehrer and M. Rezaei, "Fuzzy Logic-Based Scalable Video Rate Control Algorithm for High-Delay Applications of Scalable High-Efficiency Video Coding", *Journal of Electronic Imaging*. vol. 27, no. 4. 2018, p. 043013.
- [24] S. Akramullah, *Digital Video Concepts, Methods, and Metrics: Quality, Compression, Performance, and Power Trade-Off Analysis*, Apress, 2014.
- [25] H. Zhao, W. Xie, Y. Zhang, L. Yu, and A. Men, "An Ssim-Motivated Lcu-Level Rate Control Algorithm for Hevc", in *Picture Coding Symposium (PCS)*, 2013. 2013, IEEE. p. 85-88.
- [26] H. Zeng, A. Yang, K.N. Ngan, and M. Wang, "Perceptual Sensitivity-Based Rate Control Method for High Efficiency Video Coding", *Multimedia tools and applications*. vol. 75, no. 17. 2016, pp. 10383-10396.
- [27] GAO, W., KWONG, S., ZHOU, Y. & YUAN, H. 2016. SSIM-based game theory approach for rate-distortion optimized intra frame CTU-level bit allocation. *IEEE Transactions on Multimedia*, 18, 988-999.
- [28] S. Wang, A. Rehman, K. Zeng, J. Wang, and Z. Wang, "Ssim-Motivated Two-Pass Vbr Coding for Hevc", *IEEE Transactions on Circuits and Systems for Video Technology*. vol. 27, no. 10. 2017, pp. 2189-2203.
- [29] ZUPANCIC, I., NACCARI, M., MRAK, M. & IZQUIERDO, E. 2016. Two-pass rate control for improved quality of experience in UHDTV delivery. *IEEE Journal of Selected Topics in Signal Processing*, 11, 167-179.
- [30] L.-X. Wang, *A Course in Fuzzy Systems*, Prentice-Hall press, USA, 1999.
- [31] Z. Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli, "Image Quality Assessment: From Error Visibility to Structural Similarity", *IEEE transactions on image processing*. vol. 13, no. 4. 2004, pp. 600-612.
- [32] Z. Wang, L. Lu, and A.C. Bovik, "Video Quality Assessment Based on Structural Distortion Measurement", *Signal processing: Image communication*. vol. 19, no. 2. 2004, pp. 121-132.
- [33] [Http://Hevc.Hhi.Fraunhofer.De/Svn/Svn\\_Shvcsoftware/Tags/Shm-11.0](http://Hevc.Hhi.Fraunhofer.De/Svn/Svn_Shvcsoftware/Tags/Shm-11.0)
- [34] V. Seregin and Y. He, *Common Shm Test Conditions and Software Reference Configurations*. vol. JCTVC-Q1009. 2014.

**Frahad Raufmehrer** received his BS degree in electronics engineering and his MSc degree in communications engineering from the University of Sistan and Baluchestan, Iran, in 2014 and 2016, respectively. He is currently pursuing his PhD in electronics engineering.

**Mehdi Rezaei (IEEE SM)** received the B.S. and M.Sc. degrees in electronics engineering from the Amir Kabir University of Technology (Polytechnic of Tehran) and Tarbiat Modares University, Tehran, Iran, in 1992, and 1996, respectively, and also Ph.D. degree in Signal Processing from the Tampere University of Technology (TUT), Finland, in 2008. His research interests include multimedia signal processing and communications. He has published several papers in these fields. During his Ph.D. program, he had a close collaboration with the Nokia Research Center and he filled several patents. He also received the Nokia Foundation Award in 2005 and 2006. Between 2015, he worked as Deputy Dean for Education and Postgraduate Studies, and as Deputy Dean for Research in Faculty of Electrical and Computer Engineering, and also Dean of Education at the University of Sistan and Baluchestan, Iran. Now he is Associate Professor and the head of Communications Engineering Department, University of Sistan and Baluchestan, Iran.

# DeepSumm: A Novel Deep Learning-Based Multi-Lingual Multi-Documents Summarization System

Shima Mehrabi

Computer Engineering Department, Faculty of Engineering, University of Guilan, Rasht, Iran  
shima.mehrabi85@gmail.com

Seyed Abolghasem Mirroshandel\*

Computer Engineering Department, Faculty of Engineering, University of Guilan, Rasht, Iran  
mirroshandel@guilan.ac.ir

Hamidreza Ahmadifar

Computer Engineering Department, Faculty of Engineering, University of Guilan, Rasht, Iran  
ahmadifar@guilan.ac.ir

Received: 26/Oct/2019

Revised: 20/Nov/2020

Accepted: 25/Dec/2020

## Abstract

With the increasing amount of accessible textual information via the internet, it seems necessary to have a summarization system that can generate a summary of information for user demands. Since a long time ago, summarization has been considered by natural language processing researchers. Today, with improvement in processing power and the development of computational tools, efforts to improve the performance of the summarization system is continued, especially with utilizing more powerful learning algorithms such as deep learning method. In this paper, a novel multi-lingual multi-document summarization system is proposed that works based on deep learning techniques, and it is amongst the first Persian summarization system by use of deep learning. The proposed system ranks the sentences based on some predefined features and by using a deep artificial neural network. A comprehensive study about the effect of different features was also done to achieve the best possible features combination. The performance of the proposed system is evaluated on the standard baseline datasets in Persian and English. The result of evaluations demonstrates the effectiveness and success of the proposed summarization system in both languages. It can be said that the proposed method has achieved the state of the art performance in Persian and English.

**Keywords:** Artificial Neural Networks; Deep Learning; Text Summarization; Multi-Documents; Natural Language Processing

## 1 Introduction

Nowadays, with the advances in science and technology, there is explosive growth in the amount of available data. As a result, it is useful to have desired information in smaller volumes but with maximum coverage of the original document. Text summarization by humans has some advantages such as accuracy, coverage, and cohesiveness, but it is a time consuming and expensive process. On the other hand, summarizing huge documents is really hard for a human. Mostly the internet provides people's information, and it contains a rich amount of textual data. As a result, automatic summarization systems could lead us to save our time and efforts, even if they could not perform as well as a human in generating a summary.

Generally, the goal of automatic text summarization is compressing a text into a shorter version with preserving its main aspects. Text summarization leads us to use more

resources in a faster and more efficient way. An ideal summary should contain important aspects of one or more documents with a minimize redundancy [1].

Text Summarization can be categorized in different ways; one way refers to how a summary is organized in terms of shapes and forms. So, in this case, a summary can be abstractive or extractive. In the extractive method, the significant sentences of the document are determined, and without any modification, they are placed in summary. In abstractive summarization, a conceptual summary of the document is produced and the original form of sentences may change. Abstractive summarization is similar to the human summarization technique [1]. Another way of classifying summarization methods is based on the number of documents involved in summary, so the summarization task is divided into single-document and multi-document summarization. In single document summarization, only one document is used to create a summary, but in multi-document summarization, several documents with the same topic area construct the input of the summarization system.

\* Corresponding Author:

One of the main challenges in summarization task is how to determine the most important sentences while summary covered all significant aspects of the document and of course, without redundancy. As a result, the document has to be preprocessed and the features which represent the importance of sentences have to be exploited. The preprocessing and feature extraction phases play an important role in achieving the best result. In extractive summarization, sentences form some vectors called feature vectors. Each vector contains some features that show the importance of a sentence based on various perspectives. A feature vector has  $N$  elements that each of them has a numerical value. The importance of sentences could be determined according to the values of the feature vectors.

One of the well-known multi-document summarization systems is called MEAD [2]. MEAD was developed in two versions at the University of Michigan in 2000 and 2001. It uses a clustering method for summarization. Gistsumm is an extractive summarizer which is composed of three parts: segmentation, sentence scoring, and extract function [3]. Gistsumm scores sentences based on keywords. Keywords are determined according to the frequency of words. Sentences with the highest scores describes the main point of context more efficiently. The other sentences are chosen based on relevance to the important sentences or entire content of the text.

The earliest work on Persian text summarization is a single document extractive summarizer called Farsisum [4], it is an online summarizer, and it is developed based on Swedish summarizing project called Swesum. FarsiSum summarizes the Persian news documents in Unicode format. In another study, a Persian single document summarizer was designed, which uses the graph-based method and lexical chains [5], as ranking metrics it uses sentences similarity, the similarity of the sentence with user query, title similarity, and existing of demonstrative pronouns in the sentence. As a summary, the sentences with higher rank are selected. In [6], a multi-document multi-lingual automatic summarization system is proposed, which is based on singular value decomposition (*SVD*) and hierarchical clustering. In another system, fuzzy logic was utilized to produce a summary. In this system, some textual features such as Mean-TF-ISF, sentence length, sentence position, similarity to the title, similarity to keywords are assumed as inputs of the fuzzy system [7].

Since 2006, deep learning has persuaded lots of machine learning researchers to study and work on different aspects of it. In recent years, deep learning has influenced a vast amount of researches on signal and information processing. Deep learning uses artificial neural networks. The upper layers of the network are defined based on the outputs of the lower layers. One of the most important researches in deep learning was

published in 2009 and declared that hierarchal learning and extracting features directly from raw input data are some of deep learning characteristics [8]. Hinton et al., provided an overview of recent successes in using deep neural networks for acoustic modeling in speech recognition [9]. It is shown that deep neural networks use data more efficiently; therefore, they do not require as much data to attain the same performance of other common methods.

The result of using deep learning in speech recognition and image processing were sounds promising, which convinced natural language processing researchers to apply deep learning in Natural Language Processing (NLP) tasks. In 2011, a unified deep learning-based architecture for NLP was introduced, which is able to solve different NLP tasks such as name entity recognition, part of speech tagging, semantic role labeling, and chunking [10,11], the architecture avoids task-specific engineering as much as possible and rely on great amount of unlabeled data sets to discover internal representations which are applicable for all mentioned tasks. In [12], deep learning was applied in language modeling, and it was shown that word error rate and perplexity were decreased compared with conventional  $n$ -gram Language Models (LMs). In another study, a multi-document summarization framework was proposed based on deep learning that its feature vector contains the frequency of predefined dictionary within the documents, the framework used the deep network for developing summarizer [13]. In the first layer, the network attempts to omit unnecessary words; then, keywords are distinguished among remained words; sentences that contain keywords are extracted as candidate sentences. Finally, a summary is generated from candidate sentences via dynamic programming.

In [14], some methods are presented for extractive query-oriented single-document summarization using a deep auto-encoder to measure a feature space from the term-frequency and provides extractive summaries, gained by sentence ranking. The advantage of their approach is that the auto-encoder produces a concept vectors for a sentence from a bag-of-words input. The obtained concept vectors are so affluent that cosine similarity is adequate as the means of query-oriented sentence ranking. In [15], an extensive summarization approach was presented, which works based on neural networks. The neural network was trained by extracting ten features, including word vector embedding from the training set. For summarization, the multi-layer perceptron is applied to predict the probability of each sentence belongs to a specific class. Sentences with higher probability have a higher chance of appearing in summary. In [16], an approach was introduced for extractive single-document summarization, which applies a combination of Restricted Boltzmann Machine (RBM) and fuzzy logic to choose important sentences from the document. The set of sentence position, sentence length,

numerical token, and Term Frequency/Inverse Sentence Frequency (TF-ISF) is their feature vector. It is shown that the results produced by their method give better evaluation parameters in comparison with the standard RBM method.

Considering the achievements of deep learning, in this paper, a new summarization system is introduced, which is a multi-lingual multi-document summarizer, and it was evaluated on Persian and English documents, which achieved the state of the art results. The task of sentence extracting is based on the scores that the network assigned to each sentence. The proposed deep neural network has nine layers. In the input layer, sentence features including Term Frequency/Inverse Document Frequency (TF/IDF), title similarity, sentence position, and Part Of Speech tagging (POS) are fed to the network. After the training phase, the network is able to score sentences based on feature vectors. In the end, sentences are sorted by their scores, and top sentences are chosen for a summary.

The remainder of this paper is organized as follows. Section 2 introduces deep learning briefly, and section 3 describes the proposed method by investigating the preprocessing phase, extracting features vector, the network topology, and scoring sentences by deep learning. Section 4 presents experiments and results on both Persian and English standard data sets. Finally, the paper is closed with a conclusion in section 5.

## 2 Deep Learning

Data processing mechanism by human-like hearing and sight somehow shows the need of deep architecture extracting complicated structures of input data. For example, the human sight system uses a hierarchal structure for comprehending picture; it takes features like color, position, and direction as inputs and makes a judgment about the picture [8].

Training deep networks are complicated and difficult. The methods which are used for training shallow artificial neural networks do not work efficiently in deep networks. This issue can be solved by using a method known as unsupervised layer-wise pre-training. More precisely, in a deep learning structure, each layer is assumed independent from the others, as soon as each layer is trained, the next layer starts training by obtained input data from the previous layer. In the end, there is a fine-tuning phase on the entire deep network [17].

RBM and autoencoders are two common models in deep learning. RBM is a model for representing data probability distribution. By providing a set of training data in order to train RBM, the network adjust its parameters to find out the best probability distribution of data. RBM can be stacked to form a network, that called Deep Belief Network (DBN). The idea of DBN is that the output of each RBM serves as the inputs of the next RBM.

Therefore, by stacking RBMs, the network will be able to learn new features from previous features [18].

The Input layer of an autoencoder is the same as its output layer. This kind of network mostly is used to feature learning by encoding inputs data. Autoencoders provide a way to extract features without using tagged data. An autoencoder has an input layer that represents network input data (for example, pixels of a picture). Also, autoencoders have one or more hidden layers that indicate modified features, and it has an output layer, just like its input layer [19].

## 3 Proposed Method

For developing a text summarizer, some steps should be fulfilled to achieve a better result. First of all, the input text is preprocessed to gain a standard and less ambiguous form of the text. For showing the importance of the sentence, some metrics are described as features. Our proposed method uses deep learning for ranking sentences based on their features. To the best of our knowledge, it is the first time of utilizing deep learning in Persian text summarizer. Although the proposed summarizer is multi-lingual and it is evaluated in English as well.

In this section, the proposed summarization system (we call it DeepSumm) is explained in more detail. Preprocessing of the text, constructing feature vector, network topology, and sentence scoring task will be also covered.

### 3-1 preprocessing

Preprocessing input text is one of the basic steps in text summarization. First, the text should be normalized. Normalization refers to transforming the text into a canonical form. Sometimes a word has several dictations but the same meaning, so this sort of words should be normalized and transformed to a standard form that machine would be able to recognize them. For example, in Persian, one way to construct plural nouns is concatenation “hâ |هـ” at the end of the noun word. There are three different ways to use “hâ |هـ” based on blank space between the word and “hâ |هـ”, but all of them are correct and depend on the writer. In normalization, one of these three forms is determined as standard, and all the other forms are converted to standard form.

In the next step, the text should be segmented into sentences and words. The border of words and sentences are identified. For example, some symbols like “.” (if it is not surrounded with numbers) or newline character indicates the end of sentences. Blank space and comma indicate the borders of words. Also, in the preprocessing phase, words are stemmed, and the stop words are eliminated.

## 3-2 Constructing feature vector

In order to train DeepSumm, seven types of features were defined. In the most of the summarization tasks, these features are frequently used. The set of features includes frequency of words, title similarity, sentence position, part of speech tag, sentence stop words, sentence pronouns, and sentence length. Each sentence of the document has a feature vector that is constructed by the features mentioned earlier. Although after several experiments, it is shown that all of these seven features are not suitable for our summarization system and four of them lead us to the best result. The best four features are including TF/IDF frequency, title similarity, sentence position, and POS score. We will elaborate on the process of choosing the set of four features in more details in section 4-1.

### 3.2.1 Frequency feature

In this paper, TF/IDF is used to measure the frequency of each word of sentences. A weight is assigned to each word based on its frequency within the document. This system shows how important each word is. The frequency of a word in a document is shown by  $TF(t,d)$ , and the final weight is obtained by association of IDF. IDF means inverse document frequency, and it determines the frequency of the word in other documents. Does IDF indicate whether the word is common in all documents or not? Equation 1 shows how IDF is computed:

$$IDF(t, D) = \log\left(\frac{D}{d \in D : t \in d}\right) \quad (1)$$

$t$ ,  $D$ , and  $d$  refer to the word, all documents in the corpus, and the current document, respectively. " $d \in D : t \in d$ " is the number of documents that contain the word  $t$ .

In equation 2,  $TF(t,d)$  shows the frequency of the word  $t$  in document  $d$ . The TF/IDF of a word is obtained by multiplying TF and IDF of the word. For each sentence, the average of its word TF/IDF is assumed as the sentence TF/IDF.

$$TF/IDF_{(t.d.D)} = TF(t, d) \times IDF(t, D) \quad (2)$$

Equation 3 shows the sentence TF/IDF feature.  $S$  is the current sentence,  $w_i$  is the  $i^{\text{th}}$  word of the sentence  $S$ , and  $n$  is the sentence length (according to the number of words).

$$Sentence_{TF/IDF} = \frac{\sum_{i=1}^n TF/IDF(w_i, d, D)}{n} \quad (3)$$

### 3.2.2 Title similarity feature

The number of similar words between a sentence and the title of the document is normalized by the title length. The result is the value of the title similarity feature for the corresponding sentence. The title similarity feature is computed after preprocessing of the sentence and the title. Equation 4 shows the title similarity computation method. By normalizing the number of similarities with title length, the effect of title length is considerate on the result, because if the document has a long title, counting the number of similarities is not sufficient enough.

$$Sentence_{Title\ Similarity} = \frac{|S \cap T|}{|T|} \quad (4)$$

$|S \cap T|$  refers to the similarity between the sentence  $S$  and the document title  $T$ ,  $|T|$  refers to the title length.

### 3.2.3 Sentence position feature

Generally, the first sentences of a document (in some languages like Persian, the last sentence contains important information either [20]) are more informative than the other sentences. In the proposed summarizer, if the sentence is the first (or the last one for Persian), its corresponding value of position feature is one, and for the other sentences, the position feature is assumed zero.

### 3.2.4 Part of speech tag feature

Part of speech tagging is the process of notation a word in a text as corresponding to a specific part of speech like noun, verb, and adjective based on its description and its sense. Noun and adjective are two kinds of part of speech which can imply the most informative parts of sentences [21,22]. In this paper, the score of sentence POS is obtained by adding up the number of nouns and adjectives in the sentence, divided by the sentence length. Equation 5 shows how to calculate the POS score for a sentence  $S$ .

$$Sentence_{POS} = (S_{|N|} + S_{|Adj|}) / |S| \quad (5)$$

### 3.2.5 Sentence Stop words feature

Usually, the sentences that contain so many stop words have less important words; thus, these kinds of sentences do not imply significant information. The fraction of the sentence stop words can be considered as a metric for ranking sentences. Equation 6 shows the sentence stop words feature computation method.

$$Sentence_{Stop\ Words} = \frac{|S_{NS}^i|}{|S^i|} \quad (6)$$

The numbers of non-stop words in the sentence  $i$  shows by  $|S_{NS}^i|$  and  $|S^i|$  refers to sentence length [20].

### 3.2.6 Sentence pronouns feature

In general, when a sentence starts with a pronoun, the sentence contains some explanation about previous sentences, and it is associated with other sentences, including these sorts of sentences in summary without their related sentences may reduce the readability of the text, because it needs another sentence to complete the meaning that it is going to convey. Therefore, these types of sentences are not suitable for including in summary without their related sentence. The ratio of position of pronoun, the number of pronouns, and sentence length are represented as the value of the sentence pronoun feature. In fact, whatever the number of pronouns is more, the positive impact on sentence importance is less. Equation 7 shows how to calculate a sentence pronoun feature.

$$Sentence_{Pronoun} = \frac{BP_3^i + |S_{PR}^i|}{1 + |S^i|} \quad (7)$$

$|S_{PR}^i|$  is the number of the pronoun in a sentence  $i$ . if at least one of the first three words of the sentence  $i$  is a pronoun then  $BP_3^i$  is equal to 1 otherwise 0. Also  $|S^i|$  refers to sentence length [20].

### 3.2.7 Sentence length feature

Normally, very long or very short sentences are not suitable for including in summary text. The impact of sentence length on its importance is computed, as shown in equation 8. The ratio of  $i^{\text{th}}$  sentence length to the longest sentence in the document is shown as  $RS^i$ . Sentence length feature is obtained based on  $RS^i$  [20].

$$Sentence_{Length} = -RS^i \cdot \log(RS^i) - (1 - RS^i) \cdot \log(1 - RS^i) \quad (8)$$

$$RS^i = \frac{|S^i|}{\text{Max}_{j=1:n}(|S^j|)} \quad (9)$$

## 3-3 Scoring sentence by deep learning

After preprocessing and feature extracting phases, sentences should be ranked. In our proposed method, deep learning techniques and an autoencoder network are used for sentence ranking. The proposed autoencoder has nine layers, including the input layer. Autoencoders have an output layer equal to their input layer, and their goal is a reconstruction of input data at the output layer. So an autoencoder network is an unsupervised learning method that applies the back-propagation algorithm to achieve its goal.

An autoencoder always consists of two parts: encoder and decoder. In the encoding part, the network tries to construct new feature from input data, in decoding part the network tries to reconstruct input data from the new feature which are obtained at the end of encoding part.

Deciding the number of layers and hidden nodes of the network is an experiential task, and the best case will be determined after repetitious experiments. Different kinds of network topologies are investigated and the performances of the networks in recreating the input layer into the output layer are evaluated by measuring their errors. In this study, all of the networks have one characteristic in common, which is the number of neurons in the layer where the encoding phase is finished; this layer only has one neuron. This single-neuron layer plays an essential role in the network because it contains the score of the sentence based on its importance, which is assigned by the deep network. After all, the network with the lowest amount of error is chosen. So the network design is started by one layer as input, one hidden layer with one neuron, and one output layer as same as its input layer. Then the performance of the network is evaluated by means of calculating the error of recreating input in the output layer. Repeatedly more hidden layers are added to the networks, and in each step, the error is calculated and eventually, the network with minimum error is selected as a proposed network. The proposed network has an input layer contains neurons that are fed by a feature vector for each sentence (4 element feature vector in the best case that will discuss later in detail). The second, third, and fourth layer has 15, 10, and 5 neurons, respectively. At the

fifth layer, where the encoding phase is finished, the network has one neuron that contains the sentence score, which is assigned by the network. In fact, the network can rank a sentence and shows the importance of the sentence by the value of a single-neuron of layer five. Then the reconstruction or decoding phase is started. The network in sixth, seventh, and eighth layers has 5, 10, 15 neurons, respectively. Layer nine or output layer is the same as the input layer. Figure 1 shows the proposed neural network topology.

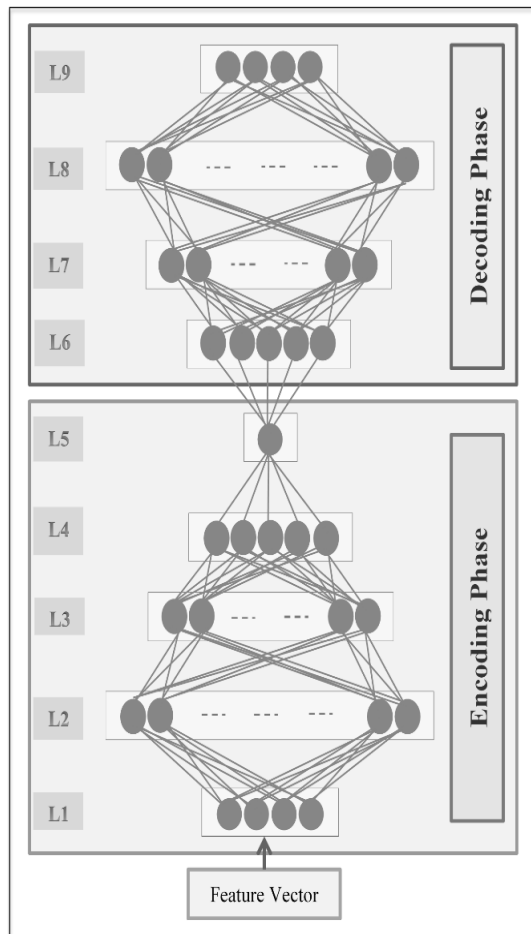


Fig. 1. Our proposed deep neural network structure.

In the training phase, the network uses a sigmoid function to predict the values of hidden neurons. The average of square error of the network error is measured, and by back-propagation error with gradient descent, the learning process is continued repeatedly till the error is minimized. After the training process, the ideal network weights are obtained; now, the trained network is ready to use for ranking the sentences.

At the end of the encoding phase, the network has one node, this node is a modified and compressed form of the

input features. Input features can be reconstructed by the value of this node. In fact, the value of this node is the score of the input sentence based on its feature vector. The ability of the network in scoring sentences based on feature vector and without interfering with the other methods or human is outstanding. One novelty of DeepSumm is the existence of a single neuron in layer five that contains a sentence score according to its importance in the document. The human assumption about the weight of each feature for assigning a score to a sentence is not considered; in fact, the network decides how important each feature is.

After training the network and adjusting its parameters, the trained network is used for scoring sentences. Sentences are sorted according to their scores. The sentences with the highest scores are selected for generating a summary, considering the compression rate.

Figure 2 illustrates the steps that the proposed method follows to generate a summary. Figure 3 is a pseudocode of the procedure of generating a summary which is used by DeepSumm.

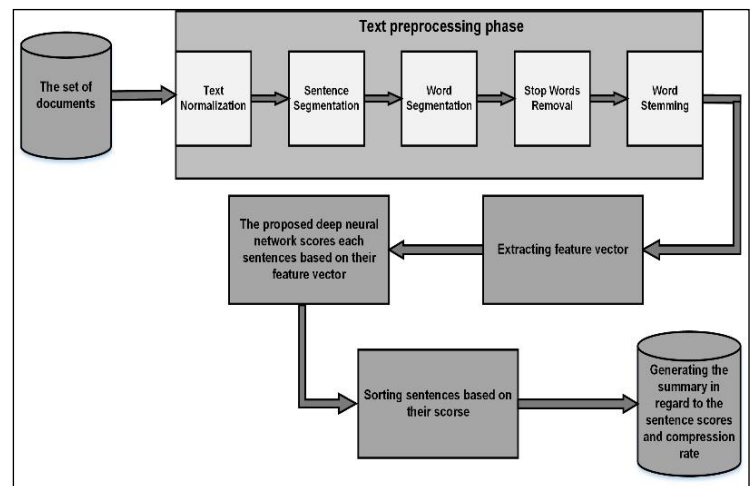


Fig. 2. Diagram of the proposed method for generating summary.



```

This procedure will generate summary by use of a deep learning method
to score the sentences.
function Preprocessing (argument document)
  normalDoc = Normalize (document)
  sentence = SentenceTokenize (normalDoc)
  foreach sentence do
    word = WordTokenize (sentence)
    if word is a Stopword then remove it from the sentence
    else
      word=stemming (word)
      preprocessedSentence=preprocessedSentence + word
  end of Preprocessing
end of Preprocessing

function FeatureExtraction (argument document)
  preprocessing (document)
  foreach sentence do
    Calculate features (
      TF/IDF,
      title similarity,
      sentence position,
      Part of Speech tagging)
    Construct feature vector
  end of FeatureExtraction
end of FeatureExtraction

function DeepLearning Scoring (featureVector)
  return the score of the sentence
end of DeepLearning Scoring

function DeepSumm (argument sentenceScores, argument CompersionRate)
  sort sentences by their scores
  regarding to the compersionRate select sentences with highest scores and generate the summery
end of DeepSumm

```

**Fig. 3.** Pseudocode of generating a summary.

## 4 Evaluation of the DeepSumm

In this section, DeepSumm is evaluated under multiple scenarios. As it was mentioned before, DeepSumm is a multi-lingual extractive summarizer, and it is tested in Persian and English. Persian can be regarded as a low resource language; therefore, the main focus of developing and testing this system is performed on Persian documents. The processing of the Persian texts is so complex and more difficult than English. Persian is among the languages with complicated preprocessing, because of different forms of writing, free word orderness, the symmetrical omission of words, and ambiguities on word segmentation [23]. Therefore, our experiments in Persian comes in various ways. Also, the DeepSumm summarizer is tested for English documents, and results sound promising.

In the following parts of section 4, the results of experiments in Persian and English are investigated thoroughly. Also, the Pasokh dataset for the Persian summarization task is introduced.

### 4-1 Experiments in Persian

The proposed summarization system used the Pasokh corpus for training and testing Persian summarizer [24]. The Pasokh is a standard corpus for evaluation and testing performance of Persian text summarization systems. Pasokh is a dataset including a variety of topics for Persian

news documents. Also, this corpus consists of gold summaries in forms of single-document, multi-document, extractive, and abstractive that is generated by a human. Pasokh has 50 topics in the multi-document section and each topic incorporating 20 documents. In total, Pasokh has 1,000 documents in the multi-document section, which 800 documents are used for training the proposed network and 200 documents for testing.

For evaluating DeepSumm, feature vectors are extracted for 2,493 sentences of test data, and the network scored each sentence according to their feature vector. In the end, one summary is generated for each topic. We used an evaluation method that considerate the exact similarity of sentences between human-generated summary and the summary of DeepSumm. It means the number of sentences in the system summery that have exact resemble sentences in the human-generated summary is considered for evaluation. The evaluation metrics (precision, recall, and f-score [25]) are calculated based on the exact similarity.

To figure out the impact of different features, DeepSumm is investigated with different kinds of features. Five different cases are assumed by combining seven features that we discussed earlier, and the network is trained and tested using them. Table 1 shows a combination of these features to form five cases of the feature vector for training the deep neural network.

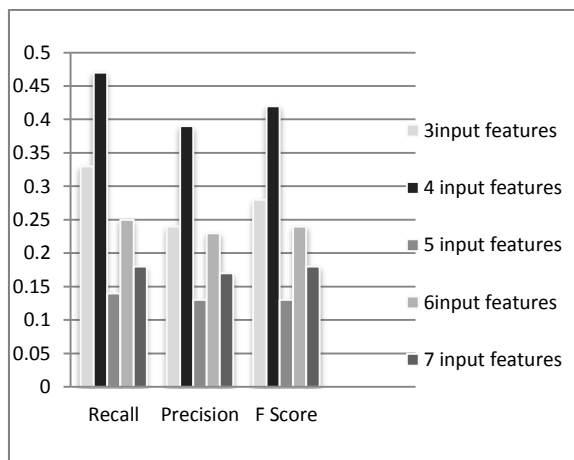
**Table 1** Five different cases according to the type and number of features.

Features	Numbers				
	3	4	5	6	7
Title similarity	*	*	*	*	*
Sentence position	*	*	*	*	*
TF/IDF	*	*	*	*	*
POS		*	*	*	*
Stop words			*	*	*
Sentence pronoun				*	*
Sentence length					*

At first, we considered a set of three features which is contained tittle similarity, sentence position, and TF/IDF. The network was trained by this set of features as its input to score sentences. The performance of DeepSumm is evaluated on test data by considering the exact similarity of human-generated summary and DeepSumm output. Based on the exact similarity of sentences, the evaluation metrics (precision, recall, and F-Score) are measured. Afterwards, by adding POS tagging to the set of features,

the four feature case is created. Respectively five, six, and seven feature cases are created by adding stop words, sentence pronoun, and sentence length feature to the former sets of the features. For all of these cases, the evaluation phase is executed, and precision, recall, and F-Score are measured. Figure 4 shows the results of the comparison of these five cases based on precision, recall, and f-score in sentence similarity evaluation. As it is shown in figure 4, by comparing recall, precision, and f-score metrics which are obtained in evaluation phase in a different set of the features cases, the best result was achieved using four features, i.e., TF/IDF, title similarity, sentence position, and POS.

One of the reasons for this result could be the data sparseness issue. In fact, increasing or decreasing the number of features may not be expressive enough to generalize on test data. The other reason could be that the network configuration was not adaptable to modifying the number of input neurons. Considering the fact of data sparseness and network configuration and based on the results of the evaluations, according to the value of all three metrics (recall, precision, and f-score), DeepSumm outperforms the other cases when it applies four features cases. As a result, this set of features is chosen for the proposed system. All further evaluation results are based on these four features.



**Fig.4.** The result of the evaluation DeepSumm with different kinds of features.

According to our studies, there is not any accessible Persian multi-document summarization system for comparing the result of the proposed summarizer. One of the prerequisites for comparing two different summarizer systems is the unity of test data. Because of the inaccessibility of other Persian multi-document summarization systems, which could be evaluated by Pasokh, our evaluation in Persian is limited to comparing the result of the proposed system with the human-generated summary that contained in the Pasokh corpus.

Table 2 shows the precise numerical result of sentence similarity evaluation on the test data when four features case is used.

**Table 2.** Sentence similarity evaluation result for Persian document from Pasokh corpus.

System	Recall	Precision	F-Score
DeepSumm	0.4667	0.3889	0.4243

Whereas DeepSumm is a multi-document summarizer, thus it is possible to have some sentences in output which are semantically similar to some sentences of the human-generated summary, but they did not use the same words. These sorts of sentences have not participated in sentence evaluation. According to table 2, it is comprehended that in sentence evaluation, the output of DeepSumm has about 50 percent similarity to Pasokh human-generated summaries.

Also, the system is evaluated by average recall scores of the Rouge toolkit. The performance of the system, in comparison with human-generated summaries, was evaluated by ROUGE [26]. ROUGE stands for Recall-Oriented Understudy for Gisting Evaluation. It contains a set of metrics for evaluating the automatic text summarization systems as well as machine translations. Its evaluation is based on comparing an automatically produced summary or translation against a set of reference summaries. ROUGE-N measures unigram, bigram, trigram, and higher-order n-gram overlap in system and reference summaries. DeepSumm is Evaluated based on ROUGE-1 (the overlap of unigrams between the system summary and reference summary) and ROUGE-2 (the overlap of bigrams between the system and reference summaries). The results are shown in Table 3.

According to table 2 and table 3 Considering the difficulty of the summarization task, especially in Persian documents because of the complexity of the Persian, the results of the evaluation are promising. It should be noted that, to the best of our knowledge, there is no study on Persian multi-document summarization task on Pasokh dataset. As a result, there do not exist any method which can be appropriately compared with our work.

**Table 3.** The result of System evaluation for Persian by Rouge-1 and Rouge-2.

System	ROUGE-1	ROUGE-2
DeepSumm	0.6850	0.5127

## 4-2 Experiments in English

In English Documents, the performance of DeepSumm was evaluated on the DUC 2005 dataset, which is the standard dataset on English summarization task. Based on Rouge scores, the performance of DeepSumm is compared with some of the most significant multi-document summarizer systems such as QODE [13], Manifold-Ranking [27], Ranking SVM [28], Regression Model [29], NIST baseline [30], MA-MultiSumm [31], MR&MR [32], and SRSum [33]. The results of the performance comparison are demonstrated in table 4.

**Table 4.** Comparison to other algorithms on DUC 2005.

System	ROUGE-1	ROUGE-2
MA-MultiSumm	<b>0.4001</b>	0.0868
SRSum	0.3983	0.0857
MR&MR	0.3932	0.0834
Manifold-ranking	0.3839	0.0676
<b>Proposed method (DeepSumm)</b>	0.3809	<b>0.1053</b>
Regression Model	0.3770	0.0761
QODE	0.3751	0.0775
Ranking SVM	0.3702	0.0711
NIST Baseline	—	0.0403

QODE is a query oriented multi-document summarizer that works by deep learning methods. It aims to extract significant concepts of documents layer by layer. Its proposed deep architecture can be divided into three distinct stages, concept extraction, reconstruction validation, and summary generation. Manifold-Ranking uses graph-based algorithms to rank sentences. The sentence relationships are divided into two categories, within the document, and cross-document relationships. Each Type of sentence relationship is considered as a separate graph with specific characteristics. In the learning phase, an extension of the basic manifold-ranking algorithm is used. Ranking SVM is a method based on Support Vector Machine (SVM) classification method. It uses a supervised learning method for ranking sentences

based on the SVM classifier. The Regression Model uses support vector regression (SVR) and some pre-defined features. It measures the importance of a sentence within a set of documents. By using different training data set, it is shown that the quality of the training data set has a significant roll in the learning process of the regression models. MA-MultiSumm is derived from CHC (Cross-generational elitist selection, Heterogeneous recombination, Cataclysmic mutation) algorithm and local search. MR&MR is an unsupervised text summarization, which can be applied to both single-document and multi-document summarization. This approach regards text summarization as a Boolean programming problem. For generating a summary, the optimization of text relevancy, redundancy, and the length of the summary are taken into account. SRSum is a deep neural network model that uses a multilayer perceptron for scoring the sentences. It works based on different kinds of sentence relations such as contextual sentence relation, title sentence relations, and query sentence relation.

The results show that the proposed system outperforms other algorithms on ROUGE-2. That means the summaries generated by DeepSumm have more bigrams overlap with reference summaries than the other systems mentioned in table 4. Based on Rouge-1, MA-MultiSumm has the best score and DeepSumm dedicates the fifth-best result to itself. As mentioned earlier QODE and SRSum use deep learning methods for generating a summary. according to Rouge-1 values, DeepSumm outperforms QODE but SRSum has 0.0174 improvements than our proposed system. In general, from the result of Rouge-2 in table 4, it can be concluded that DeepSumm achieves the best performance amongst the other representative algorithm.

## 5 Conclusion

6 In this paper, Deep Learning has been used to design and implement a multi-lingual multi-document extractive summarization system. DeepSumm is composed of two Phases, in the first phase, after preprocessing the texts, the deep network learns to rank sentences based on preset criteria and features and shows the importance of the sentence in the given document. In the second phase, according to the scores of sentences and compression rates, the system chooses the best sentences to form a summary. In the end, the result of DeepSumm has been evaluated under multiple scenarios. As our knowledge, DeepSumm is a first summarizer system based on deep learning for Persian, the result of experiment and compressions by Pasokh human-generated summary are magnificent. Also, DeepSumm is evaluated by DUC 2005, and the result is compared to some representative systems. Evaluations show that, even in English, the performance of the system

is very encouraging, and the system experiment results are successful. Based on the result of the Rouge-2, it is concluded that DeepSumm achieves state-of-the-art performance.

The main limitation of our study in Persian text summarization is the lack of any other accessible multi-document summarization system to evaluate the results. Therefore, our evaluation in the Persian document is bounded by compression of the result of DeepSumm to the human-generated summary. It is clear that having another summarization system for the assessment would give us a better view of the performance of the proposed system. In future work, we intend to design another deep network that used some other deep learning algorithms to see the results in comparison to DeepSumm.

## References

- [1] D. Das and A. Martins, "A Survey on Automatic Text Summarization," Literature Survey for the Language and Statistics II Course at Carnegie Mellon University, 2007, pp.1-31.
- [2] D. Timothy, T. Allison, S. Blair-goldensohn, J. Blitzer, A. Elebi, S. Dimitrov, E. Drabek, A. Hakim, W. Lam, D. Liu et al., "Mead A Platform For Multidocument Multilingual Text Summarization," in *International Conference on Language Resources and Evaluation*, 2004, pp. 699-702.
- [3] T. A. S. Pardo, L. H. M. Rino, and M. d. G. V. Nunes, "Gistsumm: A Summarization Tool Based On A New Extractive Method," in *International Workshop on Computational Processing of the Portuguese Language*. Springer, 2003, pp. 210–218.
- [4] M. Hassel and N. Mazdak, "Farsisum - A Persian Text Summarizer," in *Proceedings of the Workshop on Computational Approaches to Arabic Script-based Languages*, 2004, pp. 82–84.
- [5] Z. Karimi and M. Shamsfard, "Summarization of Persian Text," in *Proceedings of the 12th Computer Society of Iran*, 2007, pp. 1286-1294.
- [6] M. A. Honarpisheh, G. Ghassem-Sani, and G. Mirroshandel, "A Multidocument Multi-Lingual Automatic Summarization System," in *Proceedings of the Third International Joint Conference on Natural Language Processing: Volume-II*, 2008, pp. 733-738.
- [7] F. Kiyoumars and F. Rahimi Esfahani, "Optimizing Persian Text Summarization Based on Fuzzy Logic Approach," *Proceedings of the International Conference on Intelligent Building and Management*, 2011, pp. 264-269.
- [8] Y. Bengio, "Learning Deep Architectures for AI," *Foundations and Trends in Machine Learning*, 2009, vol. 2, no. 1, pp. 1–127.
- [9] G. Hinton, L. Deng, D. Yu, G. Dahl, A. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, B. Kingsbury et al., "Deep Neural Networks for Acoustic Modeling in Recognition," *IEEE Signal processing magazine*, 2012, vol. 29, no. 6, pp. 82-97.
- [10] R. Collobert and J. Weston, "A Unified Architecture for Natural Language Processing: Deep Neural Networks with Multitask Learning," in *Proceedings of the 25th international conference on Machine learning*. ACM, 2008, pp. 160–167.
- [11] R. Collobert, J. Weston, L. Bottou, M. Karlen, M. Kayukcuoglu, and P. Kuksa, "Natural Language Processing (almost) from Scratch," *Journal of Machine Learning Research*, 2011, vol. 12, no. Aug, pp. 2493-2537.
- [12] E. Arisoy, T. N. Sainath, B. Kingsbury, and B. Ramabhadran, "Deep Neural Network Language Models," in *Proceedings of the NAACL-HLT 2012 Workshop: Will We Ever Really Replace the N-gram Model? On the Future of Language Modeling for HLT*. Association for Computational Linguistics, 2012, pp. 20–28.
- [13] Y. Liu, S. Zhong, and W. Li, "Query-Oriented Multi-Document Summarization via Unsupervised Deep Learning," in *Proceedings of the 26th Conference on Artificial Intelligence*, 2012, pp. 1699-1705.
- [14] M. Yousefi-Azar and L. Hamey, "Text Summarization Using Unsupervised Deep Learning," *Expert System with Application*, 2017, vol. 68, pp. 93-105.
- [15] A. Jain, D. Bhatia, and M. K. Thakur, "Extractive Text Summarization using Word Vector Embedding," in *Proceedings of International Conference on Machine learning and Data Science*, 2017, pp. 51-55.
- [16] N. S. Shirwandkar and S. Kulkarni, "Extractive Text Summarization Using Deep Learning," in *Proceedings of 4<sup>th</sup> International Conference on Computing Communication Control and Automation*, 2018, pp. 1-5.
- [17] H. Geoffrey, O. Simon, and T. Yee-Whye, "A Fast Learning Algorithm for Deep Belief Nets," *Neural Computation*, 2008, vol. 18, pp. 1527-1554.
- [18] A. Fischer and C. Igel, "An Introduction to Restricted Boltzmann Machines," in *Proceedings of the 17th Iberoamerican Congress on Pattern Recognition*, 2012, pp. 14-36.
- [19] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P. Manzagol, "Stacked Denoising Autoencoders: Learning Useful Representation in a Deep Network with a Local Denoising Criterion," *Journal of Machine Learning Research*, 2010, vol. 11, pp. 3371-3408.
- [20] A. Pourmasoumi, M. Kahani, A. Toosi, A. Estiri, and H. Ghaemi, "Ijaz: A Single Document Summarization System for Persian News Text," *Signal and Data Processing*, 2014, vol. 21, no. 1, pp. 33-48.
- [21] M. Prabhakar and N. Chandra, "Automatic Text Summarization Based on Pragmatic Analysis," *International Journal of Scientific and Research Publications*, 2012, vol. 2, no. 5, pp. 1-4.
- [22] R. Mihalecea and P. Tarau, "TextRank: Bringing Order into Texts," in *Proceedings of the Empirical Methods in Natural Language Processing*, 2004, pp. 404-411.
- [23] M. Shamsfard, "Challenges and Open Problems in Persian Text Processing," in *Proceedings of the 5th Language and Technology Conference*, 2011, pp.65-69.
- [24] B. Behmadi Moghaddas, M. Kahani, S.A. Toosi, A. Pourmasoumi, and A. Estiri, "Pasokh: A Standard Corpus for the Evaluation of Persian Text Summarizers," in *Proceedings of the International Conference on Computer and Knowledge Engineering*, 2013, pp. 471-475.
- [25] D. M. Ward Powers, "Evaluation: From Precision, Recall, and F-Measure to Roc, Informedness, markedness &

- Correlation,” *Journal of Machine Learning Technologies*, 2011, vol. 2, no. 1, pp. 37-63.
- [26] C. Lin, “Rouge: A Package for Automatic Evaluation of Summaries,” in *Proceedings of the ACL Workshop on Text Summarization Branches out*, 2004, pp. 74-81.
- [27] X. Wan and J. Xiao, “Graph Based Multi-Modality Learning for Topic Focused Multi Document Summarization,” in *Proceedings of the 21st International joint conference on Artificial intelligence*, 2009, pp. 1586-1591.
- [28] T. Joachims, “Optimizing Search Engines Using Click Through Data,” in *Proceedings of the 8th International Conference on Knowledge Discovery and Data Mining*, 2002, pp. 133-142.
- [29] Y. Ouyang, W. J. Li, S. J. Li, and Q. Lu, “Applying Regression Models to Query Focused Multi Document Summarization,” *Information Processing and Management*, 2011, vol. 47, no. 2, pp. 227-237.
- [30] H. T. Dang, “Overview of DUC 2005,” in *Proceedings of the Document Understanding Conference*, 2005, pp. 1-12.
- [31] D. Mendoza, C. Cobos, E. Len, M. Lozano, F. Rodriguez, E. Herrera-Viedma, “A new memetic algorithm for multi-document summarization based on CHC algorithm and greedy search,” *Human-Inspired Computing and Its Applications*, Springer International Publishing, 2014, vol. 8856, pp. 125-138.
- [32] R.M. Alguliyev, R.M. Aliguliyev, N.R. Isazade, “An unsupervised approach to generating generic summaries of documents,” *Applied Soft Computing*, 2015, vol. 34, pp. 236-250.
- [33] P. Ren, Z. Chen, Z. Ren, F. Wei, L. Nie, J. Ma, M. de Rijke, “Sentence relations for extractive summarization with deep neural networks,” *ACM Transactions on Information Systems*, 2018, vol. 36, pp. 1-32.

2013. Now, he works as assistant professor in the Computer Engineering Department at University of Guilan. His research interests include Residue Number Systems, Computer Arithmetic and Distributed Systems.

**Shima Mehrabi** received the B.S. degree in Computer Software from Tabarestan University, Chalus, Iran in 2009, and M.S. degree in Computer Engineering from Guilan University, Rasht, Iran, in 2016. Her research interests include Information retrieval, Machine learning, Natural language processing and Data mining.

**Seyed Abolghasem Mirroshandel** received his B.Sc. degree from University of Tehran in 2005 and the M.Sc. and Ph.D. degree from Sharif University of Technology, Tehran, Iran in 2007 and 2012 respectively. Since 2012, he has been with Faculty of Engineering at University of Guilan in Rasht, Iran, where he is an Associate Professor of Computer Engineering. Dr. Mirroshandel has published more than 50 technical papers in peer-reviewed journals and conference proceedings. His current research interests focus on Natural Language Processing, Data Mining, and Machine Learning.

**HamidReza Ahmadifar** received the B.S. degree in Computer Engineering from Shahid Beheshti University, Tehran, Iran in 1997, and M.S. degree in Computer Systems Architecture from Amir Kabir University of Technology, Tehran, Iran, in 2001 and Ph.D. degree in Computer Systems Architecture from Shahid Beheshti University, Tehran, Iran in

# Social Groups Detection in Crowd by Using Automatic Fuzzy Clustering with PSO

Ali Akbari

Department of Electrical and Computer Engineering., University of Birjand, Birjand, Iran.  
ali.akbari@birjand.ac.ir

Hassan Farsi\*

Department of Electrical and Computer Engineering., University of Birjand, Birjand, Iran.  
hfarsi@birjand.ac.ir

Sajad Mohamadzadeh

Technical faculty of Ferdows, University of Birjand, Birjand, Iran.  
s.mohamadzadeh @birjand.ac.ir

Received: 04/Sep/2019

Revised: 24/Nov/2019

Accepted: 08/Dec/2019

## Abstract

Detecting social groups is one of the most important and complex problems which has been concerned recently. This process and relation between members in the groups are necessary for human-like robots shortly. Moving in a group means to be a subsystem in the group. In other words, a group containing two or more persons can be considered to be in the same direction of movement with the same speed of movement. All datasets contain some information about trajectories and labels of the members. The aim is to detect social groups containing two or more persons or detecting the individual motion of a person. For detecting social groups in the proposed method, automatic fuzzy clustering with Particle Swarm Optimization (PSO) is used. The automatic fuzzy clustering with the PSO introduced in the proposed method does not need to know the number of groups. At first, the locations of all people in frequent frames are detected and the average of locations is given to automatic fuzzy clustering with the PSO. The proposed method provides reliable results in valid datasets. The proposed method is compared with a method that provides better results while needs training data for the training step, but the proposed method does not require training at all. This characteristic of the proposed method increases the ability of its implementation for robots. The indexing results show that the proposed method can automatically find social groups without accessing the number of groups and requiring training data at all.

**Keywords:** Author Guide; Article; Camera-Ready Format; Paper Specifications; Paper Submission.

## 1- Introduction

Over the time and the extending use of a camera to maintain security and detect social anomalies, the importance of the processing video data has increased. Detecting social groups is concerned by governments for detecting dangerous situations and analyzing people's behavior [1]. Detection of theft and terrorist groups is of great importance. To identify groups with such aims, recognizing social groups is a prerequisite. In [2], social anomalies between two persons were analyzed.

According to the recent research in [3], people are interested in moving in social groups of crowds and the behavior of pedestrians in social groups of the crowd has been studied. According to the approach in [4], people are interested in moving in groups. Moving in the group means to be a subsystem in the group, in other words, a group containing two or more persons can be considered in

terms of the same direction of movement and speed of movement. According to the studies in [5], most social groups contain two persons and social groups including three and four persons are at top of the list in terms of people's willingness to appear in these groups and to move in crowds. According to the research in [6], a group is considered by the personal influence of someone or other people to move through crowds. The person who enters into a group is influenced by the group, and the speed and direction of that person are changed and grouped accordingly. Although many methods for the detection of social groups have been reported so far, all of them need training data. A method that does not need training data has not been reported yet.

This paper is organized as follows: In Section 2, a literature review is explained. In Section 3, the physical distance feature and automatic fuzzy clustering with the PSO are introduced and the proposed method is detailed for identifying the social groups. In Section 4,

\* Corresponding Author:

experimental results are demonstrated and in the final section, the conclusion is drawn.

## 2- Literature Review

For detecting social groups several methods have been reported so far. Some of the methods used many features and machine learning. These methods require training their networks. Some of the methods used people's movement and these methods are divided into three categories. Some new methods have been used deep learning, deep neural network and Long Short Term Memory (LSTM) in recent years.

In [7], the detection of social groups to help the robot's behavior in teamwork with humans has been reported, and an anticipation method by linear extrapolation of inter-event intervals has been used. In [8], social groups have been detected by skeletal data from participants, and the event anticipation method has been used to move robots among human groups. The method reported in [9] detects conversation and social groups for the robot by using the direction of people and the lower-body orientation. In [10], body, head orientation and body orientation from a distance social scene are used to detect social groups. As an example, the output of this method is shown in Figure 1.



Fig. 1. head orientation and body orientation used in [10].

In [11], social groups are detected by using physical distance, people's locations and head directions in definite spaces. In this method, fixed and unmoving groups are detected. Figure 2 shows the various models concerned by this method.

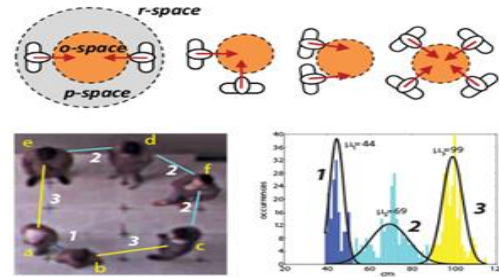


Fig. 2. Verity of direction in definite spaces [11].

Social group detection based on direction of people's movement is divided into three categories: group-based, individual-group joint and individual-based. The reported method in [12] applies the second-order derivative of the information matrix of people's movement path to detect social groups.

Individual-group joint approaches use more important information compared to group-based approaches, for instance, group tracking [13]. In individual-based approaches, the groups are detected by using single people's trajectories. As an example, in [14], the head is more accurately detected using crowd density estimation.



Fig. 3. Effect of crowd density estimation in head detection in a crowded area.

The reported method in [15] predicts two persons in one group using distance features, speed difference, and time overlap information with Support Vector Machines (SVM).

In [16], probabilistic similarity methods of two-persons variations in successive frames and soft areas are used to identify social groups.

An attractive method without training is introduced in [17], which uses group information in the previous frame. In this method, a potentially infinite mixture model of group probability with the mean value of distance in burst frames is examined. In this method, only the distance feature is used and, better results are achieved in the recognition of social groups by using more features.

In [18], proximity and velocity characteristics are used to identify social groups. In this method, in crowded conditions, Markov model, proximity and velocity features are assisted, and in some situations which are not very

crowded, the pursuit of individuals in successive frames is determined to recognize the social groups.

In Figure 4, the status of people is shown in sequential frames.

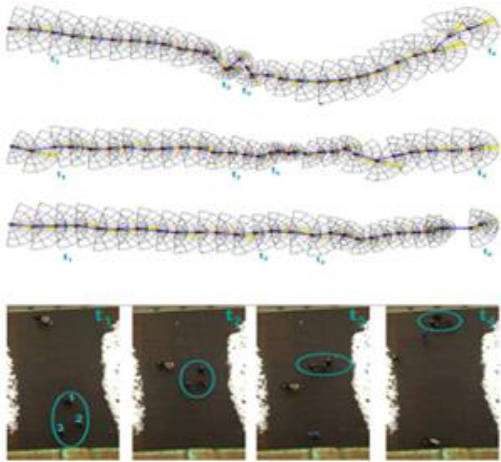


Fig. 4. Tracking people in consecutive frames to identify social groups [18].

The reported method in [18] has detected social groups by using the features of distance and attention to people or the direction of the people in sequences of frames and game theory tools.

In the method reported in [19], which is inspired by electric dipole shown in Figure 5, each person's eyesight is firstly examined, and if a relationship between the attention of the people's eyes is found, this group of people is placed in a social group.

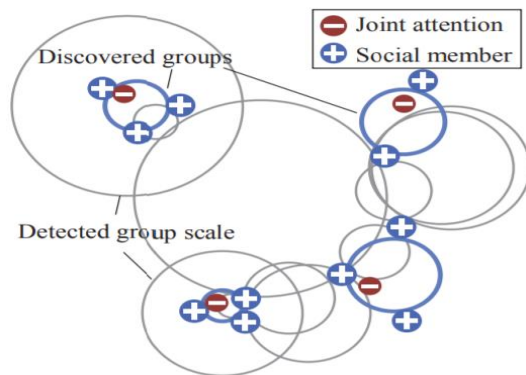


Fig. 5. Effect of investigating the relationship between visual attention in the detection of social groups [19].

In the reported method in [20], the relationship between people is identified by using graph-based clustering and further social group activity is detected by the SVM-based classification.

In [21], the track of salient points and adaptive clustering are used, and hierarchical social groups are achieved. In the method reported in [22], by weighing two features including distance and speed correlation between two persons, social groups are detected, and the group's

behavior is analyzed. In [23], instead of using the similarity feature between two persons, the clustering of pedestrians into different groups is based on using the start and endpoint to identify social groups. In further analysis of square matrices, the size of the number of people in the video is constructed and the probability of being grouped is calculated between all of them. The probabilities are calculated based on Euclidean distance between two individuals. Hierarchical clustering is then used to identify social groups. In [24], conjoint individuals and group tracking have been reported. RGB histogram, region covariance, and Histogram Of Gradient (HOG) similarity are used to detect social groups.

In [25], features of distance, motion causality, trajectory shape, and path convergence are used by optimized the SVM to detect social groups. The reported method in [26] has used deep learning algorithms, trajectory modeling approaches with LSTM, contextual information from the local neighborhood and Generative Adversarial Network (GAN) to detect social groups.

### 3- Proposed Method

In the proposed method, automatic fuzzy clustering with the PSO has been used to cluster all people in the video. After automatic fuzzy clustering with the PSO, groups are detected. Post-processing is used to remove scattered groups. After the post-processing, final social groups are maintained. In each iteration, the automatic algorithm merges two similar clusters and the fuzzy PSO changes 'w' to approach a global solution. In the following, the proposed method is discussed in four main parts as shown in Figures 6 and 7.

---

#### Algorithm:

---

- 1: Input = people's location in video
  - 2:  $k$  = number of people in video // only for the first iteration
  - 3: for iteration < max iteration
  - 4:     clustering by PSO for  $k$  clusters
  - 5:      $S$  // calculate scattering for all clusters (calculate using //equation (4))
  - 6:      $R$  // calculate similarity between two clusters (calculate //using equation (6))
  - 7:     if (max (similarity < 2) and similarity>0.4) between two // clusters
  - 8:         make one cluster with them and  $k = k-1$
  - 9:     end if
  - 10:      $w$  // fuzzy operation identified weight for PSO
  - 11:     end for
  - 12: post-processing // delete this group (maximum distance //between group members from the center of the group >0.4)
  - 13: output = social groups detection
- 

Fig. 6. Algorithm: the proposed method.



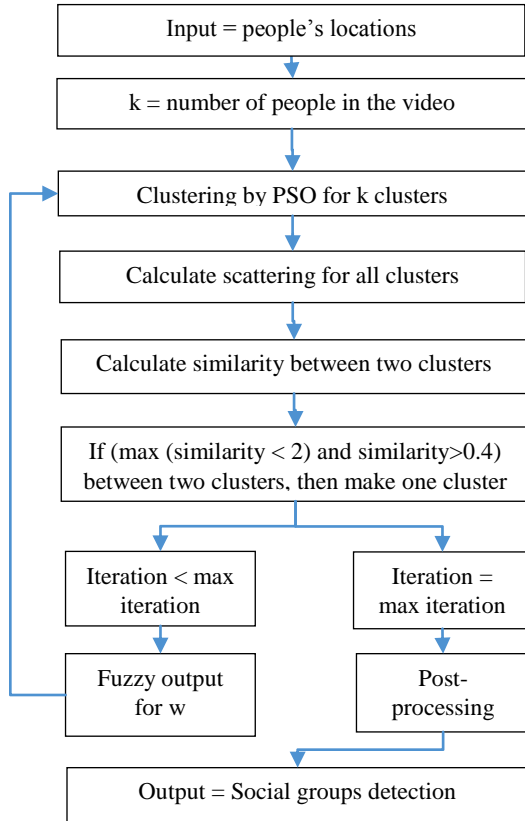


Fig. 7. Flowchart of the proposed method.

In the following, all parts of the flowchart depicted in Figure 7 are described. The physical distance feature is used for clustering by calculation of scattering and similarity.

### 3-1- Physical Distance Feature

One of the most important features in detecting social groups is physical distance between people in a scene. The proxemics theory is identified based on physical distance features [27]. This theory is shown in Table 1.

Table 1. The proxemics Theory [27]

space	Boundaries(m)	description
Intimate	0.0 – 0.5	Unmistakable involvement
Personal	0.5 – 1.2	Familiar interactions
Social	1.2 – 3.7	Formal relationships
public	3.7 – 7.6	Non-personal interactions

The proxemics theory includes 4 classes to identify the relation between pairs and their physical distance. Each class is separated from another class by the boundaries shown in Table 1.

### 3-2- Random Algorithm

In [28], a randomized algorithm, with the PSO has been reported based on collective intelligence. In this

algorithm, random particles are firstly generated. Next, all the particles are applied to a fitness function and the best solution among all solutions,  $P_g^t$ , is the leader of all particles. For all particles, the best position of the particles is saved as the leader of the  $P_i^t$  particle. In Equation 1,  $Y_i^t$  is the position of the  $i$ -th particle at the moment 't', and  $V_i^{t+1}$  is the velocity of the particle [28].

$$V_i^{t+1} = w * V_i^t + c_1 * r_1 * (P_g^t - Y_i^t) + c_2 * r_2 * (P_i^t - Y_i^t) \quad (1)$$

In this equation,  $r_1$  and  $r_2$  are random numbers between zero and one, and  $c_1$  and  $c_2$  are considered to be 2 for the PSO algorithm [28]. The speed of each particle is updated at each iteration. In Equation 1, 'w' is the setting parameter for exploiting or exploring the search space. If 'w' is large, the speed will be higher and the search step will increase. This search for the problem space is called exploration. If w is small, the speed will reduce and the search step will be short. This type of searching is called exploitation. Next, according to Equation 2, the positions of all particles are updated [28].

$$Y_i^{t+1} = Y_i^t + V_i^{t+1} \quad (2)$$

In Equation 3, 't' is the current repetition and 'maxt' is the maximum iteration to get the best solution in the search space [28].

$$w = \frac{-0.8 * t}{maxt} + 0.9 \quad (3)$$

### 3-3- Clustering

For a more precise clustering, the Davis-Bouldin (DB) index is used [29].

#### 3-3-1- Davis-Bouldin Index

Different measures are used for clustering. One of the most important measures is the Davis-Bouldin index [30, 31]. For a more precise clustering, the distance between cluster members and the center of their clusters should be minimized, and the cluster members' distance between the other clusters should be maximized. In the following, the formulation of this concept is discussed. Scattering of a cluster,  $S_i$ , is given by [30]:

$$S_i = \left( \frac{1}{T} \sum_{j=1}^{T_i} |X_j - A_i|^p \right)^{\frac{1}{p}} \quad (4)$$

where  $A_i$  is the cluster center and 'T' is the number of members located in the  $i$ -th cluster. If p is equal to 2, then Euclidean distance is calculated between the cluster members and the cluster centers.

The disparity between the two clusters is defined as follows [30]:

$$d_{ij} = \|A_i - A_j\|_p = \left( \sum_{k=1}^n |a_{k,i} - a_{k,j}|^p \right)^{\frac{1}{p}} \quad (5)$$

Where  $k$  is the dimension for the center in the cluster. For both equations,  $p$  is considered to be 2.

The concept of similarity between two clusters is defined as [30]:

$$R_{ij} = \frac{S_i + S_j}{d_{ij}} \quad (6)$$

The similarity between the two clusters,  $R_{ij}$ , has the following properties:

$$R_{ij} \geq 0$$

$$R_{ij} = R_{ji}$$

If  $S_i$  and  $S_j$  both are zero, then  $R_{ij}$  will be zero.

If  $S_j \geq S_b$  and  $d_{ij} = d_{ib}$  then  $R_{ij} > R_{ib}$

If  $S_j = S_b$  and  $d_{ij} \leq d_{ib}$  then  $R_{ij} > R_{ib}$

The DB index is now defined as [29]:

$$DB = \frac{1}{N} \sum_{i=1}^N D_i, \quad D_i = \max_{j \neq i} R_{ij} \quad (7)$$

where  $N$  is the number of clusters. By considering the maximization of the similarity between two clusters, since the distance between the two centers is in the denominator of the similarity between two clusters,  $R_{ij}$ , it is needed that the distance between the two centers of similar clusters to be minimized. This means that two similar groups are considered in the same group. Of course, by maximizing the similarity between two clusters, the, DB index will also be larger. As a result, with a larger DB index, the clustering can detect better clusters.

### 3-3-2 Automatic Clustering

Automatic clustering means that the number of clusters is unknown, as an unsupervised, and the algorithm must recognize the number of clusters and perform clustering [33, 34]. For this aim, in the process of finding the similarity between two clusters, the groups with the similarity between two clusters under the threshold of 0.4 are removed. Also, if the sum of the similarities between the two clusters is less than 2, then two groups with the highest similarity between two clusters are merged. This means that the number of groups is reduced and, of course, the condition of being less than 2 is to reduce the number of clusters at first when the number of groups is large. Then reducing the number of clusters will be stopped when approaching the correct numbers of the clusters.

### 3-3-3 Automatic Clustering with PSO

For all clusters, PSO randomly specifies cluster centers. In the following, with automatic clustering, the number of clusters is reduced, but 'w' and the particle size are not reduced. This means that in each iteration all centers are updated by the PSO, but the cluster centers that satisfy the automatic clustering conditions will be entered into the next processing stage. In the initial implementation, in the proposed method, the number of clusters per video is equal to the number of people existing in the same video. In the following, the cluster centers of each video are randomly identified by the PSO.

### 3-3-4 Automatic Fuzzy Clustering with PSO

PSO is a relatively recent heuristic search method which is based on the idea of collaborative behavior and swarming in biological populations. The PSO senses population-based search approaches and depending on information sharing among their population members enhances their search processes using a combination of deterministic and probabilistic rules. The weakness of the PSO is detecting local solutions and not approaching a global solution when it solves with fuzzy PSO.

In the PSO algorithm, the weight 'w' decreases when iteration increases and reaches to its minimum in final iterations. Decreasing 'w' means that the search is around the previous solution 's'. This also means that the PSO finds the approximate location of the solution  $s$  and decreases 'w' to find the exact location around the current location. However, the problem is that it may not be close to the proper solution. To solve this problem, fuzzy controllers are used. It is even possible to use fuzzy controllers for an individual and group leader's impact factor in the PSO algorithm. In the proposed method, fuzzy control is used only for 'w'. In each execution, the best scattering solution of a cluster  $S_i$  is examined. If the best dispersion solution among all clusters is very low, 'w' is considered small, which means we are close to the global solution, and we need to perform a thorough review around it. Now, if the best solution among all clusters is a large number, then we have a far-reaching global solution, and we need to search the global solution with a large 'w' within the entire search space.

### 3-4 Post-processing

The output includes clustering with clusters far from the center of the cluster and also a short distance from the center of the cluster, but long-distance clusters are not suitable for the detection of social groups. As a result, a threshold for maximum distance between group members from the center of the group is considered and the threshold in the proposed method is 0.4. The first video among twenty videos in student03 database contains 1475

frames. These frames are converted into 25 episodes and the average position of the people is recorded. Due to the difference in the number of frames and therefore the average positions of the people (for example in the twentieth video, there are five average positions of the people) in the proposed method we have converted all videos into five average positions of individuals in equal intervals. In each of the five average positions of the individuals per video, 400 replications are performed by the PSO for automatic fuzzy clustering.

Among all parts of each video the positions of all people are averaged as shown in Figure 8. This means that the positions of all people in groups that move from one part of an image to another part are placed near them.



Fig. 8. Few frames from the first set of the first video

As shown in Figure 9, the convergence is achieved by the PSO after 400 iterations for one part of the five episodes of a video. After convergence, the post-processing is started. As shown in Figure 10, the post-processing neglects the groups with high dispersion. In Figure 10, circles are drowned in the center of the groups.

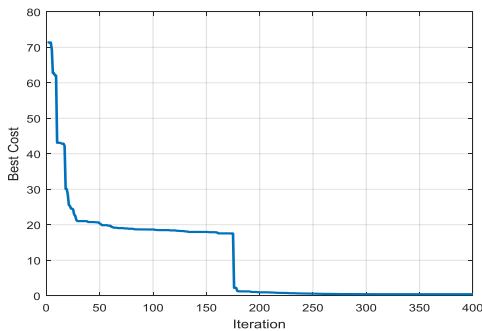


Fig. 9. Automatic fuzzy clustering convergence by the PSO for 400 iterations for one part of the five episodes of a video.

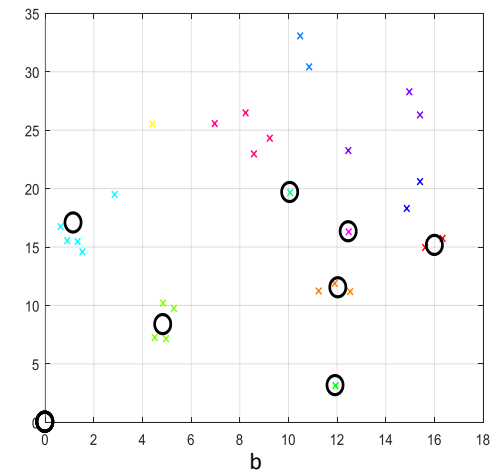
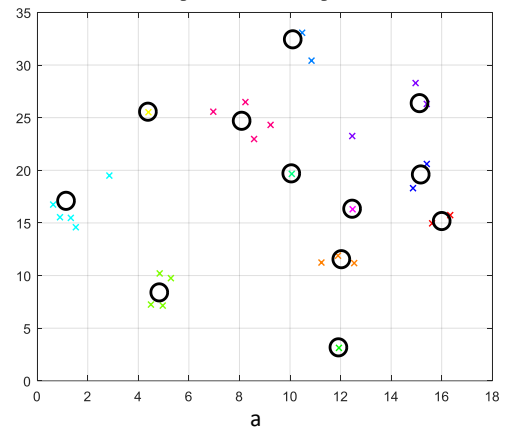


Fig. 10. (a) Automatic fuzzy clustering output by the PSO, (b) Automatic fuzzy clustering by PSO after post-processing and removal of high-dispersion clusters.

## 4- Experimental Results

### 4-1- Datasets

To evaluate the performance of the proposed method, ETH and Hotel [33], Student003 [25], GVEII [26, 35] and MPT-20x100 [25] datasets have been used. These datasets include the trajectory of people and the number of people. As an example in Figure 11, the trajectory and the number of people are shown in a sequence of frames.

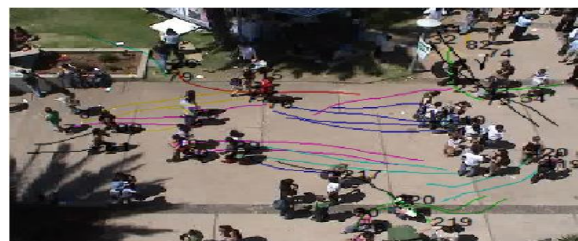


Fig. 11. The trajectory of people's movement and the number of people in the sequence of frames.

In Table 2, five datasets are compared. The parameters  $v$ ,  $p$ ,  $g$ ,  $d1$ , and  $d2$  represent the number of videos, the number of people, the number of groups, the minimum distance in the group ( $m$ ), and the maximum distance in the group ( $m$ ), respectively.

Table 2. Comparing four datasets.

	$v$	$p$	$g$	$d1$	$d2$
ETH	1	117	18	0.99	2.79
Hotel	1	107	11	0.75	2.00
Student003	20	406	108	0.41	0.71
GVEII	30	630	207	0.77	1.66
MPT-20x100	20	82	10	0.63	1.45

## 4-2- Evaluation Measures

The precision and recall parameters are used to compare the results. The precision is the ratio of the number of groups that are correctly identified to the number of all identified groups; the recall is the ratio of the number of groups that are correctly identified to the number of real groups of the database. The standard F-score,  $F1$ , is defined as follow [33]:

$$F1 = \frac{2 * \text{precision} * \text{recall}}{\text{precision} + \text{recall}} \quad (8)$$

## 4-3- Indexing Results

Since ETH, Hotel, student003, GVEII and MPT-20x100 datasets have different properties, they are used to evaluate real scenarios. In this section, we present the obtained results of the proposed method and compare it with other methods. The precision and recall parameters are used to compare the proposed method with social constrained structural learning (SCSL) for group detection in the crowd [25], vision-based analysis of small groups (VASG) in pedestrian crowds [18], who are you with and where are you going? (WWG) [15], online Bayesian non-parametric (OBNP) for group detection [17], scene-independent group profiling in the crowd (SIGP) [12], conjoint individual and group tracking (CIGT) [24] and Generative Adversarial Networks for Trajectory Prediction and Group Detection in Crowds (GD-GAN) [26].

Table 3. The obtained results in ETH and Hotel datasets.

	ETH		Hotel	
	precision	recall	precision	recall
proposed method	90.56	<b>86.79</b>	<b>89.23</b>	<b>91.81</b>
	$\pm 0.34$	$\pm 0.42$	$\pm 0.67$	$\pm 0.26$
CIGT [24]	<b>96.24</b>	56.16	-	-
SCSL [25]	91.14	83.43	89.12	91.32
VASG [18]	80.72	80.78	88.93	89.35
WWG [15]	60.64	76.42	84.06	51.27
OBNP [17]	79.12	80.76	88.12	73.25
SIGP [12]	69.33	68.26	67.38	64.11

The precision and recall are calculated and presented in Table 3 for the ETH and the Hotel datasets. Each measure is reported in terms of mean and standard deviation over 5 runs to account for the stochastic nature of the clustering of our algorithm. Although the proposed method is not needed to be trained, the recall of the proposed method in the ETH dataset is 86.79 and the recall of the proposed method in the Hotel dataset is 91.81, which is the best between all methods. The recall of the SCSL method in the ETH dataset is 83.43 and the recall of SCSL in the Hotel dataset is 91.32. The precision of the CIGT and the SCSL methods in the ETH dataset are 96.24 and 91.14, respectively, where the precision of the proposed method is 90.56.

The precision of the proposed method in the Hotel dataset is 89.23, which is the best performance among other methods. The precision of the SCSL method in the Hotel dataset is 89.12.

To find the precision and recall for all 20 videos, 20 videos of the Student003 dataset and 30 videos of the GVEII dataset are calculated separately in five sections, and the averages of these five measures are compared. Tables 4 and 5 present the average of all sections for the student003 and the GVEII datasets, respectively. Each measure is reported in terms of mean and standard deviation over 5 runs to account for the stochastic nature of the clustering of the proposed algorithm for the student003 and the GVEII datasets.

Table 4. Obtained precision and recall for the student003 dataset

	Precision	recall
proposed method	80.78	<b>91.81</b>
	$\pm 0.36$	$\pm 0.21$
SCSL	81.72	82.51
GD-GAN	82.14	63.47
VASG	77.28	73.69
WWG	56.76	76.02
OBNP	71.12	78.76
SIGP	40.48	48.63

The results of the same experiment for student003 dataset are summarized in Table 4. The recall of the proposed method in the student003 dataset is 91.81, which has the best performance. The results of the same experiment for the GVEII dataset are summarized in Table 5. The recall of the proposed method in GVEII dataset is 96.31, which has the best performance.

Table 5. Obtained precision and recall for the GVEII database

	precision	recall
proposed method	<b>85.53</b>	<b>96.31</b>
	$\pm 0.24$	$\pm 0.14$
SCSL	84.12	84.11
GD-GAN	83.16	79.54
VASG	80.14	79.45
WWG	57.84	75.51
OBNP	70.15	76.65

SIGP	44.85	49.92
------	-------	-------

As shown in Tables 4 and 5, the recall of the proposed method in the student003 and the GVEII datasets is the best among all methods. Note that the proposed method is not needed to be trained.

Some of the videos in the MPT-20x100 database are selected and F-score is calculated for the proposed method. Table 6 presents the F-score of the proposed method in each video, and compare it with the SCSL and the VASL methods. The results of the same experiment for the MPT-20x100 dataset are summarized in Table 6. As shown in Table 6, F-scores of the proposed method are 79.38, 97.86, 98.47, and 96.19 for the airport1, chinacross4, grand1, and thu10 videos, respectively; these are the highest scores compared to other methods.

Table 6 Results of comparison in MPT-20x100 database

F-score	airport1	china cross4	grand1	thu10
proposed method	<b>79.38</b>	<b>97.86</b>	<b>98.47</b>	<b>96.19</b>
SCSL	78	96	97	90
VASG	58	92	83	82

Effect of the threshold in the post-processing is examined, and the precision and recall are determined with the threshold of 0.3, 0.4 and 0.5 for the student003 and the GEVII datasets. The obtained results are presented in Tables 8 and 9. Note that in all previous Tables the threshold is 0.4.

Table 8. Precision and recall for different thresholds in student003 dataset.

Threshold=0.3		Threshold=0.4		Threshold=0.5	
precision	recall	precision	recall	precision	recall
64.25	92.35	64.38	91.81	62.43	90.13

Table 9. Comparison of thresholds in the GEVII dataset

Threshold=0.3		Threshold=0.4		Threshold=0.5	
precision	recall	precision	recall	precision	recall
85.32	97.11	84.48	96.31	80.97	92.69

As observed, the best performance is provided by a threshold of 0.3. This means that removing larger groups results in better precision and recall.

As an example, Figures 12 and 13 show detecting social groups by the proposed method.

## 5- Conclusion

Detecting social groups is one of the most important problems that has been concerned recently to analysis interpersonal relations in groups. In this study, an

automatic fuzzy clustering with the PSO was used, and acceptable results were achieved.

It is a matter of great importance, without supervision and training, of the automatic fuzzy clustering with the PSO. This method does not require to be trained and is easier to be calculated and implemented for human robots compared to the methods in which the parts of the search space are considered as the training data. The skipping of two persons' movement in opposite directions causes a mistake in identifying the group. It is even possible that two persons who are in the same group will not be together from the beginning and continue on the same path for sometimes, in this case it would be a mistake to distinguish the two persons from the proposed method at the beginning of the route.



a



b

Fig. 12. (a) Input video, (b) video after detecting social groups with the proposed method.



a



Fig. 13. (a) Input video, (b) video after detecting social group with the proposed method.

## References

- [1] Mehran, R., Oyama, A., Shah, M., : 'Abnormal crowd behavior detection using social force model,' in Computer Vision and Pattern Recognition, Conference on IEEE., 2009, pp. 935-942.
- [2] Manzi, A., Fiorini, L., Limosani, R., Dario, P., Cavallo, F.: 'Two-person activity recognition using skeleton data', IET Computer Vision, vol. 12, no.1, pp. 27-35, 2017.
- [3] Moussaïd, M., Perozo, N., Garnier, S., Helbing, D., Theraulaz G., : 'The walking behaviour of pedestrian social groups and its impact on crowd dynamics,' PloS one, vol.5, no.4, pp. 10047-10054, 2010.
- [4] Bandini, S., Gorrini, A., Manenti, L., Vizzari G.,: 'Crowd and pedestrian dynamics: Empirical investigation and simulation,' in Proceedings of Measuring Behavior, pp. 308-311, 2012.
- [5] Van Stekelenburg, J., & Klandermans, B. Individuals in movements: A social psychology of contention. In Handbook of social movements across disciplines (pp. 103-139). Springer, Cham, 2017.
- [6] Luber, M., Stork, J. A., Tipaldi, G. D. Arras, K. O.: 'People tracking with human motion predictions from social forces,' in Robotics and Automation (ICRA), 2010 IEEE International Conference on IEEE, pp. 464-469, 2010.
- [7] Iqbal, T., Moosaei, M., Riek, L. D.,: 'Tempo adaptation and anticipation methods for human-robot teams,' in RSS, Planning HRI: Shared Autonomy Collab. Robot. Workshop, 2016.
- [8] Iqbal, T., Rack, S., Riek, L. D.,: 'Movement coordination in human-robot teams: a dynamical systems approach,' IEEE Transactions on Robotics, vol.32, no.4, pp. 909-919, 2016.
- [9] Vázquez, M., Steinfeld, A., Hudson, S. E.,: 'Parallel detection of conversational groups of free-standing people and tracking of their lower-body orientation,' in Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on IEEE, pp. 3010-3017, 2015.
- [10] Ricci, E., Varadarajan, J., Subramanian, R., Rota Bulò, S., Ahuja, N., Lanz, O.,: 'Uncovering interactions and interactors: Joint estimation of head, body orientation and f-formations from surveillance videos,' in Proceedings of the IEEE International Conference on Computer Vision, pp. 4660-4668, 2015.
- [11] Cristani M., et al.,: 'Social interaction discovery by statistical analysis of F-formations,' in BMVC, pp 4-11, 2011.
- [12] Shao, J., Change Loy, C., Wang, X.: 'Scene-independent group profiling in crowd,' in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2219-2226, 2014.
- [13] Pang, S. K., Li, J., Godsill, S. J.,: 'Detection and tracking of coordinated groups,' IEEE Transactions on Aerospace and Electronic Systems, vol. 47, no.1, pp. 472-502, 2011.
- [14] Rodriguez, M., Laptev, I., Sivic, J., Audibert, J.-Y.,: 'Density-aware person detection and tracking in crowds,' in Computer Vision (ICCV), 2011 IEEE International Conference on IEEE, pp. 2423-2430, 2011.
- [15] Yamaguchi, K., Berg, A. C., Ortiz, L. E., Berg, T. L.,: 'Who are you with and where are you going?,' in Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on IEEE, pp. 1345-1352, 2011.
- [16] Chang, M.-C., Krahnstoeber, N., Ge, W.,: 'Probabilistic group-level motion analysis and scenario recognition,' in Computer Vision (ICCV), 2011 IEEE International Conference on IEEE, pp. 747-754, 2011.
- [17] Zanotto, M., Bazzani, L., Cristani, M., Murino, V.,: 'Online bayesian nonparametrics for group detection,' in Proc. of BMVC, 2012.
- [18] Ge, W., Collins, R. T., Ruback, R. B.,: 'Vision-based analysis of small groups in pedestrian crowds,' IEEE transactions on pattern analysis and machine intelligence, vol.34, no.5, pp. 1003-1016, 2012.
- [19] Park H. S., Shi, J.,: 'Social saliency prediction,' in Computer Vision and Pattern Recognition (CVPR), Conference on IEEE, pp. 4777-4785, 2015.
- [20] Tran, K. N., Bedagkar-Gala, A., Kakadiaris, I. A., Shah, S. K.,: 'Social Cues in Group Formation and Local Interactions for Collective Activity Analysis,' in VISAPP, vol. 1, pp. 539-548, 2013.
- [21] Shao, J., Dong, N., Zhao, Q.,: 'An adaptive clustering approach for group detection in the crowd,' in Systems, Signals and Image Processing (IWSSIP), 2015 International Conference on IEEE, pp. 77-80, 2015.
- [22] Voon, W. P., Mustapha, N., Affendey, L. S., Khalid, F.,: 'A new clustering approach for group detection in scene-independent dense crowds,' in Computer and Information Sciences (ICCOINS), 2016 3rd International Conference on IEEE, pp. 414-417, 2016.
- [23] Khan, S. D., Vizzari, G., Bandini, S., Basalamah, S.,: 'Detection of social groups in pedestrian crowds using computer vision,' in International Conference on Advanced Concepts for Intelligent Vision Systems, Springer, pp. 249-260, 2015.
- [24] Yight, A., Temizel, A.,: 'Individual and group tracking with the evaluation of social interactions,' IET Computer Vision, vol. 11, no.3, pp. 255-263, 2016.
- [25] Solera, F., Calderara, S., Cucchiara, R.,: 'Social constrained structural learning for groups detection in crowd,' IEEE transactions on pattern analysis and machine intelligence, vol.38, no.5, pp. 995-1008, 2016.
- [26] Fernando, T., Denman, S., Sridharan, S., & Fookes, C., 'GD-GAN: Generative Adversarial Networks for Trajectory Prediction and Group Detection in Crowds,' In Asian Conference on Computer Vision, Springer, Cham, pp. 314-330, 2018.

- [27] Reicher, S. D., Spears, R., Postmes, T.,: 'A social identity model of deindividuation phenomena,' *European review of social psychology*, vol.6, no.1, pp. 161-198, 1995.
- [28] Eberhart R., Kennedy, J.: 'A new optimizer using particle swarm theory,' in *Micro Machine and Human Science*, IEEE, pp. 39-43, 1995.
- [29] Peng, Hong, et al. "An automatic clustering algorithm inspired by membrane computing." *Pattern Recognition Letters* 68 pp. 34-40, 2015.
- [30] Farsi, H., Mozaffarian MA., Rahmani, H.,: 'Improving voice activity detection used in ITU-T G.729.B,' *Proc. of 3rd WSEAS International Conference on Circuits, Systems, and Telecommunications*, pp. 11-15, 2009.
- [31] Farsi, H.,: 'Improvement of minimum tracking in minimum statistics noise estimation method,' *Signal Processing An International Journal (SPIJ)*, vol. 4, no.1, pp. 1-17, 2010.
- [32] Khosravi, H., Moradi, E., Darabi, H.,: 'Identification Of Homogeneous Groundwater Quality Regions Using Factor And Cluster Analysis; A Case Study Ghir Plain Of Fars Province,' 2015.
- [33] Hosseini, S. M., Farsi, H., Yazdi, H. S.,: 'Best clustering around the color images,' *International Journal of Computer and Electrical Engineering*, vol.1, no.1, pp. 20-29, 2009.
- [34] Hosseini, S. M., Nasrabadi, A., Nouri, P., Farsi, H.,: 'A novel human gait recognition system,' *International Journal of Computer and Electrical Engineering*, vol.2, no.6, pp. 1043-1055, 2010.
- [35] Sadeghian, A., Kosaraju, V., Sadeghian, A., Hirose, N., Rezaatofghi, H., & Savarese, S. Sophie.,: 'An attentive gan for predicting paths compliant to social and physical constraints,' In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1349-1358, 2019.
- [36] Wang, Q., Chen, M., Nie, F., & Li, X., 'Detecting coherent groups in crowd scenes by multiview clustering,' *IEEE transactions on pattern analysis and machine intelligence*, vol.40, no.7, pp. 910-919, 2018.

**Ali Akbari** received the B.Sc. and M.Sc degrees in Electronics engineering from University of Birjand, Birjand, Iran, in 2013 and 2015, respectively. Since 2015 he was accepted as a Ph.D candidate student in communications engineering, university of Birjand, Birjand, Iran. His area research interests include Image and Video Processing, Retrieval, Pattern recognition, Digital Signal Processing.

**Hassan Farsi** received the B.Sc. and M.Sc. degrees from Sharif University of Technology, Tehran, Iran, in 1992 and 1995, respectively. Since 2000, he started his Ph.D in the Centre of Communications Systems Research (CCSR), University of Surrey, Guildford, UK, and received the Ph.D degree in 2004. He is interested in speech, image and video processing on wireless communications. Now, he works as professor in communication engineering in department of Electrical and Computer Eng., University of Birjand, Birjand, IRAN.

**Sajad Mohamadzadeh** received the B.Sc. degree in communication engineering from Sistan & Balochestan, University of Zahedan, Iran, in 2010. He received the M.Sc. and Ph.D. degree in communication engineering from South of Khorasan, University of Birjand, Birjand, Iran, in 2012 and 2016, respectively. Now, he works as assistant professor in Faculty of Technical and Engineering of Ferdows, University of Birjand, Birjand, Iran. His area research interests include Image and Video Processing, Retrieval, Pattern recognition, Digital Signal Processing, Sparse Representation, and Deep Learning.

# Facial Images Quality Assessment based on ISO/ICAO Standard Compliance Estimation by HMAX Model

Azamossadat Nourbakhsh

Science and Research Branch, Islamic Azad University, Tehran, Iran  
a.nourbakhsh@srbiau.ac.ir

Mohammad-Shahram Moin\*

IT Research Faculty, ICT Research Institute, Tehran, Iran  
moin@itrc.ac.ir

Arash Sharifi

Science and Research Branch, Islamic Azad University, Tehran, Iran  
a.sharifi@srbiau.ac.ir

Received: 20/Aug/2019

Revised: 24/Oct/2019

Accepted: 20/Nov/2019:

## Abstract

Facial images are the most popular biometrics in automated identification systems. Different methods have been introduced to evaluate the quality of these images. FICV is a common benchmark to evaluate facial images quality using ISO / ICAO compliancy assessment algorithms. In this work, a new model has been introduced based on brain functionality for Facial Image Quality Assessment, using Face Image ISO Compliance Verification (FICV) benchmark. We have used the Hierarchical Max-pooling (HMAX) model for brain functionality simulation and evaluated its performance. Based on the accuracy of compliancy verification, Equal Error Rate of ICAO requirements, has been classified and from those with higher error rate in the past researches, nine ICAO requirements have been used to assess the compliancy of the face images quality to the standard. To evaluate the quality of facial images, first, image patches were generated for key and non-key face components by using Viola-Jones algorithm. For simulating the brain function, HMAX method has been applied to these patches. In the HMAX model, a multi-resolution spatial pooling has been used, which encodes local and public spatial information for generating image discriminative signatures. In the proposed model, the way of storing and fetching information is similar to the function of the brain. For training and testing the model, AR and PUT databases were used. The results has been evaluated by FICV assessment factors, showing lower Equal Error Rate and rejection rate, compared to the existing methods.

**Keywords:** Facial Images Quality; ISO/IEC19794 Standard; ICAO; FICV; HMAX Model

## 1- Introduction

Facial recognition quality assessment is one of the most important factors in the automatic face recognition accuracy. Face recognition has been announced by the International Civil Aviation Organization (ICAO), as a biometric feature of machine-verified. The International Institute for Standardization (ISO) has proposed ISO / IEC19794-5, which includes face image information requirements, environmental conditions and shooting features [1]. Since there are many testing requirements, it is difficult to determine the compliancy of a face image with ISO / ICAO standards. Fully automating of a face image compliancy detection with ISO / ICAO standards has many benefits such as no need to human experts and accelerating the document production process. Researches

about the quality assessment of commercial systems have shown that their performance needs to be improve for standard compliancy verifying and has not still reached the human accuracy level [2], [3]. Therefore, automated facial images quality assessment, is one of the most challenging issues in automated document production process. To evaluate the quality of produced algorithms, and comparing their performance, the FICV's Biolab benchmark is provided by the University of Bologna Biometrics Research Group. The Face Image ISO Compliance Verification (FICV) test, which includes assessing the requirements introduced in the face recognition standard, performs face evaluation and recognition. This benchmark includes a ground truth database, a well-defined testing protocol, and baseline algorithms for all ISO / ICAO requirements. Some of the 24 requirements (looking away, unnatural skin tone, hair across eyes, head rotation (roll/pitch/yaw



Greater than 8°), red eyes, shadow across face, frame too heavy, frame across eyes and mouth open) were used as quantitative variables of the problem. The compliancy of each of these requirements in the image is returned with a 1, 0, 1 score by the proposed model which shows three logically compliance, noncompliance, and dummy modes [4].

The Hierarchical Maximum (HMAX) pooling model is used to encode the properties of the noticeable face components. HMAX acts like a MAX operator and extracts location and scale-independent features for detection. This model expresses a hierarchy of brain regions through which object recognition is performed in the cerebral cortex. The purpose of this model is the cognitive phenomena describing in terms of simple and complex computational processes in an acceptable physiological model. To perform computational processes, two layers are embedded in this model:

- Simple "S" layers are derived from local filters convolution to compute higher order features from the different types of units in the previous layer.
  - Complex "C" layers are stabilized by fetching units and the number of units is reduced by sub-sampling and all position and scale information are deleted simultaneously.
- Object detection in the cerebral cortex extends through the feedforward ventral visual pathway. It travels through the primary visual cortex (V1) and reaches the InferoTemporal cortex (IT) by passing other visual areas V2 and V4. Different layers of HMAX model used for feedforward brain ventral visual pathway simulation in this research, are shown in Table 2 (which are driven from [5]).

Table 2: different HMAX layers for brain Visual pathway simulating (driven from [5])

Layers of the HMAX model	Ventral Visual Pathway of the Cerebral Cortex
S <sub>1</sub> (The first Simple layer)	V1 (The primary Visual Cortex)
C <sub>1</sub> (The first Complicated layer)	V2 (The second Visual cortex)
S <sub>2</sub> (The Second Simple layer)	V4 (The fourth Visual cortex)
C <sub>2</sub> ((The Second Complicated layer)	IT (The InferoTemporal cortex)

In this work, we have proposed, for the first time, a new integrated system for ISO/ICAO face image compliancy, which is based on HMAX method, inspired from the human brain functionality. The main contribution of our work is using an integrated method for different requirement compliancy assessment, compared to the existing methods, which use their specific features for each requirement compliancy assessment, separately. We have also approved the superiority of our approach over the existing methods, using the experimental results.

## 2- Related Works

Many researches have performed on facial images quality assessment based on ISO / ICAO standard requirements [6-12]. From the results of these studies, it can be concluded that this area of research still needs further development.

The Hierarchical Max-pooling model (HMAX) is a feedforward model mimicking the structures and functions of V1 to posterior inferotemporal (PIT) layer of the primate visual cortex, which could generate a series of position and scale- invariant features.

In [13] to mimic the attention modulation mechanism of V1 layer, a bottom-up saliency map is computed in S1 layer of the HMAX model, which can support the initial feature extraction for memory processing. Also to mimic the short-term memory to long-term memory conversion abilities of V2 and IT, an unsupervised iterative clustering method is used for clusters learning with multiscale middle level patches. Simulation results show that the enhanced model with a smaller memory size, exhibits higher accuracy than the original HMAX model and other unsupervised feature learning methods in multiclass categorization task.

An object recognition model by extracting features temporally and utilizing an accumulation to bound decision-making model is introduced in [14]. This model accounted recognition time and accuracy. In face recognition, for temporally extracting informative features, a hierarchical spiking neural network, called spiking HMAX is modified. In the decision making part of the model, the extracted information accumulates over time using accumulator units. The input category is determined if any of the accumulators reaches a threshold, called decision bound. Testing Results showed that the model follows human accuracy in a psychophysics task better than the classic spiking HMAX model.

For Image classification, a method based on ontology and HMAX features performed by integrating clusters [15]. This method relied on training visual-feature classifiers according to the taxonomic relationships between image categories. Using the HMAX model, the visual features and the concepts were extracted from the image categories. The taxonomic relationship between visual features and concepts were created to make an ontology that represents the semantic information associated with the training images. Using ontology-based HMAX and Bag-of-Visual-Words (BoVW) models, superior performance achieved over baseline methods. To evaluate this method, the Inception-v3 deep learning network was used, and the classification results performed better in some image classes than Inception-v3.

Bottom-up attention is crucial to primary vision and helps reducing computational complexity. In [16], a bottom-up attention model was presented based on the C1 features of

HMAX model. Attention modeling in layer C1 of the HMAX model showed better results than Graph-Based Visual Saliency (GBVS).

In [17], a face recognition model was presented that used the visual attention model using skin color features to find saliency maps of the face candidate areas and the C2 texture features in the visual cortex of the HMAX model for face recognition. After finding candidate face areas, C2 texture features were extracted for face or non-face areas classification using a support vector machine classifier. Experimental results on the Caltech Face Database with background, showed that the proposed model was reliable against variations in face brightness, expression and cluttered backgrounds.

In [18] Binary HMAX model (B-HMAX) was introduced. In the C1 layer of this model, using image patches selection instead of random usage in the standard model at the training phase increased the accuracy and decreased the calculating costs. Also using Hamming distance instead of Euclidean distance for calculating the distance between patches, increased the speed.

In [19], the original HMAX model [5] was used and the end-of-network filters, which integrated local filters, were modified for producing complex filters to cover larger and more complex areas of the image. To better discriminate the image content, they trained the coefficients of each filter in the last layer. This increased the discrimination and also the invariance.

A flexible multilayer radial method for the outputs' pooling of the filters in the image was presented. Neurons in the inferior temporal visual cortex (IT) are known for regions by varying sizes accepting [20]. This is called local areas of various sizes pooling in the visual field, causing slight variability with respect to spatial location. The multi-resolution pooling introduced in the study was equivalent to applying a specific filter in a spatial neighborhood with different radial pooling that caused different levels of invariance. The optimal level of invariance with a single classifier was obtained by training at higher levels of the network [21]. The classification in this research showed better results than previous architectures. Since this method achieved very good results with increasing discrimination and invariance, it has been used in the present study.

As it is mentioned above, there are many researches that payed attention to the face recognition subject using HMAX model. However, none of them has worked on facial image quality assessment by this model. Thus, in this study, we evaluated the suitability and effectiveness of using HMAX model for the facial image quality assessment.

In the following sections, first, ICAO requirements selection process has been described and after creating key and non-key patches from face components by Viola-Jones algorithm, HMAX model is introduced and employed in

FICV process. The results of executing the proposed model on AR and PUT databases have been evaluated using standard Facial image quality verification factors.

### 3- Requirements Selection

Considering the scope and content of the assessment factors operations; which are introduced in ISO / IEC19794-5, first, some requirements from 24 FICV requirements should be selected. According to Ferrara et al. [9], the error rate obtained for each requirement has been divided into three categories. Table 1 can be used for identifying the need for further research on requirements and selecting new research areas for the requirements assessment results improvement.

Table 1: Biolab's ICAO Requirements' classification based on their Accuracy Rates (Driven from [9])

ICAO Requirements Difficulty Diagnosing	Accuracy Rate	Name of Requirement
Easy Diagnosing Requirements	$EER < 3\%$	ICAO08(pixelation),ICAO10 (Eye Closed),ICAO13(Flash Reflection on Skin), ICAO15 (shadow behind Head), ICAO17(Dark Tinted lenses) , ICAO18( Flash Reflection on Lenses), ICAO22(Veil over Face)
Middle rate Diagnosing Requirements	$3\% \leq EER \leq 7\%$	ICAO02(Blurred) , ICAO04 (Ink Marked/Creased), ICAO05 (Unnatural Skin Tone), ICAO06 (Too Dark/Light), ICAO11 (Varied Background) , ICAO14( Red Eyes) , ICAO19 (Frames too Heavy) , ICAO20( Frame Covering Eyes) , ICAO23( Mouth Open)
Hard Diagnosing Requirements	$EER > 7\%$	ICAO01(Eye Location), ICAO03 (Looking Away) , ICAO07(Washed Out) , ICAO09(Hair Across Eyes) , ICAO12 (roll/pitch/yaw Greater than 8°),ICAO16 (Shadow Across Face), ICAO21 (Hat/CAP), ICAO24 (Presence of other Faces or Toys too Close to Face)

Requirements selection in this study, are based on the results of Table 1. Requirements with less than 3% error rates, which are not challenging, are ignored. Image and background features have also been excluded, and we have focused on nine following facial requirements:

Look Away (ICAO03), Unnatural Skin Tone (ICAO05), Hair Across Eyes (ICAO09), Head Rotation

(roll/pitch/yaw Greater than  $8^\circ$ ) (ICA012), Red Eyes (ICA014), Shadow Across Face (ICA016), Frame Too Heavy (ICA019), Frame Across Eyes (ICA020), Mouth Open (ICA023).

#### 4- Proposed HMAX model for Face Image Quality Assessment

The proposed model for face image quality assessment using the HMAX method is shown in Figure 1. In the following, its different parts has been described:

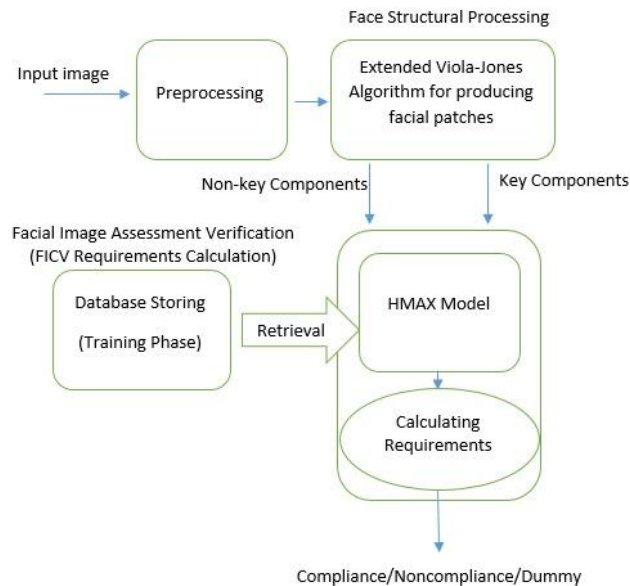


Fig. 1 Proposed method for Facial Image Quality Assessment using HMAX Model

##### 4-1- Pre-Processing

In this section, basic image processing is performed to produce a suitable image. First the image is examined for possibility of being tokenizable (without padding). So the database images should be corrected as follows:

- Distance between eyes ( $E_{Dist}$ ) be at least 60 pixels.
- Rectangular area be  $W * H$  (with  $W = 4 * E_{Dist}$  and  $H = W * 4 / 3$ ) size. Eyes aligned horizontally, centered at  $C_E = (W * 1 / 2, W * 3 / 5)$  which is generally enclosed in the original image (Figure 2).

Basic tasks of this part are face detection (ROI), face alignment and normalization (which can be used for reducing illumination effects).

##### 4-2- Face Structural Processing

In this part, face semantic patches of the key and non-key components are obtained using elemental images based on a component-based approach.

Based on biological evidences, the diagnosing operation is performed in two operation of location estimation and semantic division. In the estimation section, the center of four face key components locations (left eye, right eye, nose, and mouth) are estimated. To implement this section, a hybrid method using Viola-Jones and skin color pixel detection is used; that causes more accurate detection of the facial components location and increases detection speed [22]. Different semantic patches are formed by segmentation based on the location estimation results and by component-based method. In this way, the primary image is divided into 8 patches, 4 of which consist of key components of the face (left eye, right eye, nose, and mouth) and 4 patches containing non-key facial components (left cheek, right cheek, forehead and chin). For being almost every patch component-based, the size of divisions must be specific to each individual. Hence, the size of each patch is obtained based on a constant rate of the distance between the two eyes centers.

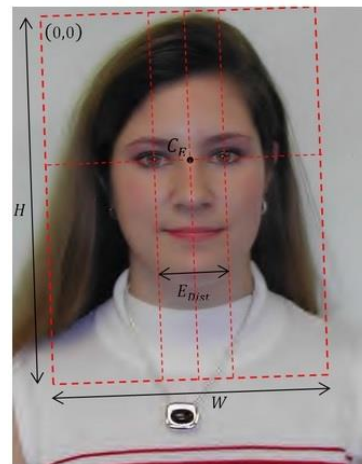


Fig. 2 Geometric properties of the obtained image format [1]

##### 4-2-1- Producing Facial Patches using Viola-Jones algorithm

Michael Jones and Paul Viola [23] developed the famous Viola-Jones Algorithm for face detection in 2003. In this algorithm, learning is performed by measuring the similarity of two sample images. A set of computationally efficient rectangular features (Haar features), are described and operate on a pair of input images. The features compare inside the input images areas in different

locations, scales and directions. To quickly evaluating these features from integral images and performing the training of facial similarity by features, the Adaboost algorithm is used. Finally, a hierarchical classifier is considered for rejecting windows that have not been recognized as a face. As this method is very accurate but time consuming, a very fast detection algorithm based on skin tone pixel detection has been merged to it in [22]. In the case of eyes and mouth detection, physical location approximation is made in detected face to locate the eyes and mouth. This method increased the accuracy of system and decreased its consumed time.

Since the introducing of the Viola-Jones algorithm, many researchers have used this method in their researches, which [24], [25] and [26] are amongst them.

#### 4-2-2- Feature Selection

Selecting Features for Multiclass Classification is an essential step in pattern recognition and machine learning applications. Specially, a big challenge is an optimal subset selecting from high-dimensional data, which has much more variables than observed and contains noise or outliers. We used the feature selection method presented in [27]. In that research, a feature selector named Fisher-Markov is presented to identify the features; which are more important in describing the essential differences around possible groups.

It is a systematic method of factors optimizing for the best feature subset selecting, to identify factors for sparsity and separability in the high dimensional scenarios. Since the introduced method is linear in number of features and quadratic in number of observations, it operates very quickly. In pattern recognition and model selection view, in the proposed model, it is easily possible to select the most discriminable subset of variables by solving an objective function without constraint.

In supervised classification, with the training data  $\{(x_k, y_k)\}_{k=1}^n$ , where  $x_k \in R^p$  are the  $p$  dimensional feature vectors and  $y_k \in \{w_1, \dots, w_m\}$  are classes' tags, the most important features should be selected for the most separable representation of the multi-class classification with  $m$   $C_i$  class, where  $i = 1, 2, \dots, m$ . Each  $C_i$  class has  $n_i$  observation. With a new test observation, the selected features are used for predicting an unknown class tag for each observation. In order to global optimization and

effective feature selection by Fisher-Markov method, in the feature subset, for large  $p$ , some specific kernels including polynomial kernels  $k$  have been considered [28] [29].

$$k(x_1, x_2) = (1 + \langle x_1, x_2 \rangle)^d \quad (1)$$

where  $d$  is the parameter degree, alternatively:

$$k'(x_1, x_2) = (\langle x_1, x_2 \rangle)^d \quad (2)$$

### 4-3- Face Image Assessment Confirmation Module (FICV Attribute Detection)

Inspired by biological evidences, facial image compliancy verification of the proposed model, for simulating memory structure such as brain functionality, includes code generation, storage, retrieval and final decision. The inputs of this section are the face key and non-key components patches (and thus it is a component-based model) and the output of this section is the result of ICAO requirements recognition.

#### 4-3-1- HMAX Model

As illustrated in [5] and shown in Figure 3, a general HMAX model is designed of frequency of pooling and convolution layers. Each convolutional layer has a series of feature maps and each pooling phase produces changing resistance against these feature maps. In the following, different layers of HMAX model are described.  $S_1$ ,  $C_1$ ,  $S_2$  and  $C_2$  layers of HMAX model are named L1, L2, L3 and L4, respectively.

- First layer

Each feature map  $Ll_{\delta, \theta}$  is produced by the input image convolution against a set of Gabor filters  $g_{\delta, \theta}$  (Eq. 3), with orientation  $\theta$  and scale  $\delta$ . These filters are used to simple cell activation in the V1 region of the visual cortex modeling [5].

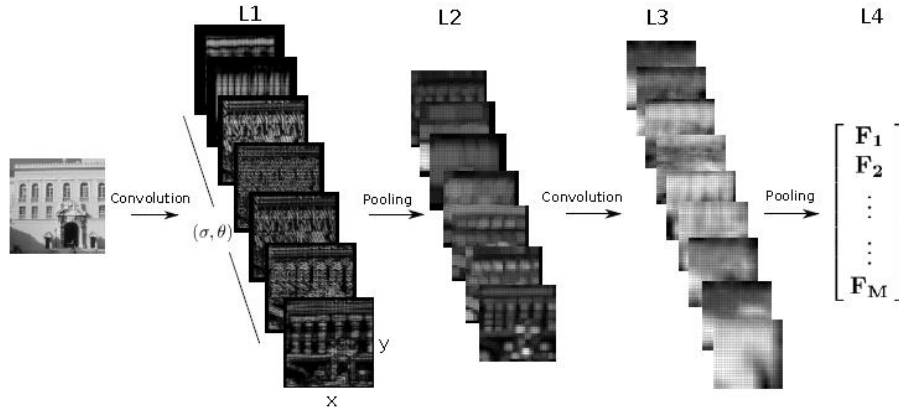


Fig. 3 A general convolution network: this network alternates feature mapping layers (convolution) and feature pooling layers alternately. The convolution layers produce information of a particular feature, and the pooling layers create invariance by relaxing the configuration of these features [5].

$$g_{\sigma,\theta}(x,y) = \exp\left(\frac{x_0^2 + \gamma y_0^2}{2\sigma^2}\right) \cdot \cos\left(\frac{2\pi}{\lambda} x_0\right) \quad (3)$$

where  $x_0 = x \cos\theta + y \sin\theta$  and  $y_0 = -x \sin\theta + y \cos\theta$ . Parameter  $\gamma$  shows the aspect ratio of the filter and  $\lambda$  is its wavelength.

With image  $I$ , the first layer in orientation  $\theta$  and the scale  $\delta$ , can be expressed as an absolute value convolution product, as follows:

$$L1_{\sigma,\theta} = |g_{\sigma,\theta} * I| \quad (4)$$

- Second layer

Each feature map of  $L2_{\delta,\theta}$  is a dimension reduction of  $L1_{\delta,\theta}$ , that is obtained by the maximum number of local neighborhoods selecting. Maximum pooling impact on local neighborhoods is the invariance of local conversions and global transformations [30].

The second layer divides each  $L1_{\delta,\theta}$  map into small neighborhoods  $u_{i,j}$  and finds the maximum value inside each  $u_{i,j}$  such that:

$$L2_{\sigma,\theta}(i,j) = \max_{u_{i,j} \in L1_{\sigma,\theta}} u_{i,j} \quad (5)$$

By keeping only the maximum output at two scales adjacent to each point (i, j), scale invariance can be achieved to some extent.

- Third layer

The  $L3$  layer at the  $\delta$  scale is obtained by the  $\alpha^m$  filters convolution against the  $L2_{\delta,\theta}$  layer, which are called HL filters.

$$L3_{\sigma}^m = \alpha^m * L2_{\sigma} \quad (6)$$

HL filters are visual descriptors of mid-level regions in the image that combine low-level Gabor filters with multiple orientation in one scale.

- Fourth layer

To achieve general invariance, the final step (last signature) is calculated by the maximum  $L3_{\delta}^m$  output selecting in all location conditions and scales. Thus the last layer is a vector of  $M \sim 1000$  dimension, which determines each coefficient of each HL filter maximum output on the scale  $\delta$  and location (x, y).

$$L4 = \begin{bmatrix} \max_{(x,y),\sigma} L3_{\sigma}^1(x,y) \\ \vdots \\ \max_{(x,y),\sigma} L3_{\sigma}^M(x,y) \end{bmatrix} \quad (7)$$

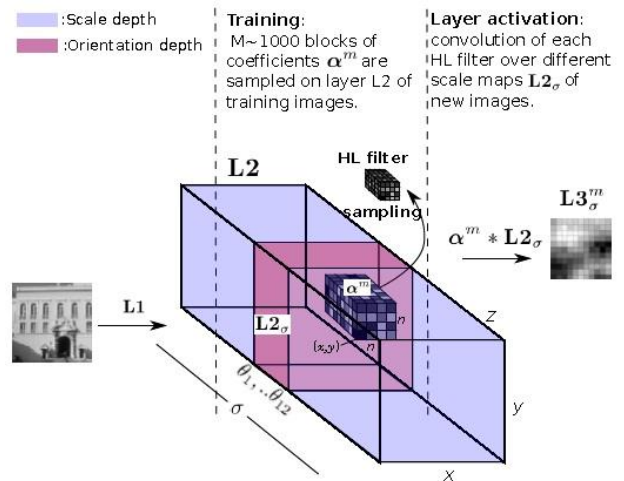


Fig. 4 Third level of HMAX: in training Phase,  $M \sim 1000$  HL filters is defined by L2 coefficients of sampling blocks. Layer activation in each image is obtained by convolving each HL filter on all positions of each  $L2_{\sigma}$  scale map [5].

HMAX model has been used in several researches including [19], [20], [24] and [25]. HMAX method in the model proposed in this work is based on [19], where the first level local filters are merged with more sophisticated filters at the previous level, producing a flexible descriptor of the object regions and combining local information across multiple scales and orientation. These filters are invariant and discriminative, making them more suitable for visual classification. It also introduces a multi-resolution spatial pooling that encodes local and public spatial information for generating image discriminative signatures. In Figure 5, each HL filter is convolved simultaneously on several scales that focus on the scale  $\delta$ . In training phase, the coefficients associated with weak scales and orientations, receive zero values; which makes the filter more discriminative and ignores weaker scales and orientations during the test phase.

**4-3-2- Code Generation**

Real code creating in the human brain, means information sensing and receiving from the environment in the form of physical and chemical stimuli. Especially when looking at a particular face, the brain encodes various facial features with many patterns. It is assumed that there are 9 coding patterns (for each face, 9 ICAO requirements would be stored in long-term memory). To mimic this fact functionally, the HMAX descriptor can be used to encode these requirements.

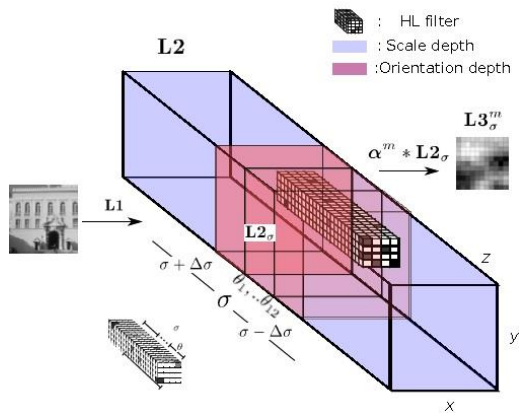


Fig. 5 Third level operations - Each HL filter measures simultaneously on several scales [19].

To prepare and extract the *CI* level from the *HMAX* model, split face patches  $cpath_{ij}$  from known individuals can be used where  $i = 1, 2, \dots, 8$  (for 8 Key and non-Key face components patches) And  $j = 1, 2, \dots, N$ , where  $N$  represents the number of known people (known people refers to those whose ICAO attribute recognition tags exist in the database). *CI* patches of an ICAO requirement, produce a patch cluster  $C_i$  where  $C_i \in cpath_{ij}$ . For each  $C_i$ , at the *C2* level of the *HMAX* model for each face image, the  $f_i \in R^N$  feature is extracted.  $f_i$  is the

final consideration feature that introduces a face component for ICAO requirement recognition (Figure 6).

**4-3-3- Storing**

In the mammals' long-term memory, different features of a known object are regularly stored in distributed areas, and common features of various known objects are stored through aggregating. Labeled faces are database images for which ICAO requirements assessment are labeled. As shown in Figure 7, the same strategy was chosen to store the ICAO requirements of labeled faces. For example, the first studied ICAO requirement of different individuals  $f_{ij}$  ( $j = 1, 2, \dots, N$ ) are stored together and constitute a subspace storage of an ICAO requirement. The first person's attributes  $f_{i1}$  (where  $i$  is one of the 9 case study requirements) are stored separately in the distributed subspace. The storing phase is similar to a training procedure in a general requirements estimation method, and does not include unlabeled faces.

**4-3-4- Decision Making**

This step identifies and assesses the ICAO requirements of a new face image from labeled faces images. To identify a person's requirements, it is needed to retrieve requirements for all labeled faces before decision making, which is called retrieval.  $R_{ij}$  is used for the requirement assessment of a new face image based on the  $j^{\text{th}}$  labeled faces requirements by using a notable face feature  $f_i$ .  $R_{ij}$  can be estimated using a support vector machine. There are 9 Binary SVM classifier for assessment of each studied requirements. Final classification results shows 1 for compliance, 0 for non-compliance and -1 for dummy classes. For 1 and 0 classes, the compliancy of each requirement in the image is returned with a score in range of zero to 100 by regression.

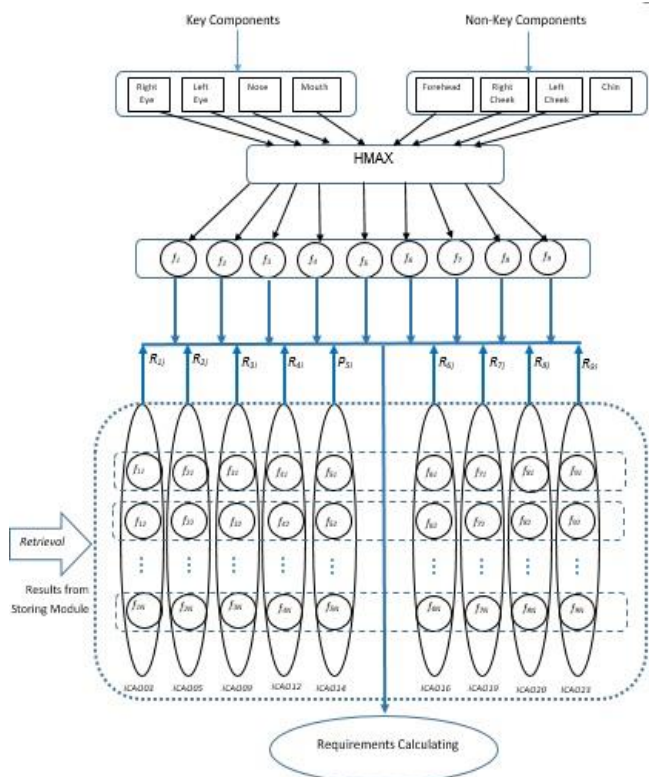


Fig. 6 Proposed Face Image Assessment Confirmation Module containing HMAX model for detecting compliancy with ICAO requirements.

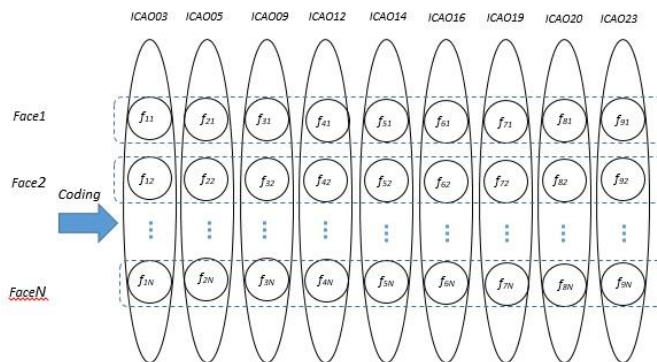


Fig. 7 Storage module structure.

### 5- Simulation Results

For simulating the proposed model, a system with 16GB of RAM, Core i7 processor, 2TB hard drive has been used.

In this study, 1741 images of the AR database with size of 576 \* 768 [32] and 291 images of the PUT database

with size of 1536 \* 2048 [33] were used to train and test the proposed method, which include 310 fully compatible images (compatible with all requirement) and 1722 incompatible images (incomparable at least with one requirement). Database division to test and train set, has been done by K-fold cross-validation algorithm, where the best result was obtained with K=10. For objective performance evaluation in the database, ground-truth data has been employed. Some of these images are manually labeled, each image containing eye corners information and position and shooting features based on three logically compliance, noncompliance, and dummy modes. Dummy value is used for uncertainty situations (for example, when one uses sunglasses, it is almost impossible to detect open or closed eyes). Two types of errors can occur during the compliancy assessment of face images:

- 1) Declaring compliancy for an image which is not compliant (False Match Rate) (FMR):
- 2) Declaring incompliance for an image which is compliant (False Non Match Rate) (FNMR).

$$FMR = \frac{FP}{FP+FN} \tag{8}$$

$$FNMR = \frac{FN}{TP+TN} \tag{9}$$

A good biometric system should illustrate a small amount of FMR and FNMR. High FMR indicates high system error and low FNMR indicates low system functionality for all studied cases' acceptance. ROC, DET charts and EER are used to analyze a biometric system.

**EER:** Equal Error Rate (EER) is indicated by interaction between FMR and FNMR. EER represents the error rate at a  $t$  threshold such that the False match rate and the false non-match rate are equal ( $FMR(t) = FNMR(t)$ ). EERs are calculated for compliancy rate checking and used for each feature performance evaluating.

**Rej:** Rejection Rate refers to the percentage of the face images, which cannot be processed by the proposed method. This rejection can be due to the pixellation, hair across eyes and shadow across the face. For comparing the results of simulation, three SDKs (two commercial SDKs and the BiolabSDK [9]) were used. The name of commercial SDKs cannot be disclosed, because of the specific agreement with their providers in [9]. Table 3 shows the EER and Rejection rates for comparing three SDKs results with the proposed method.

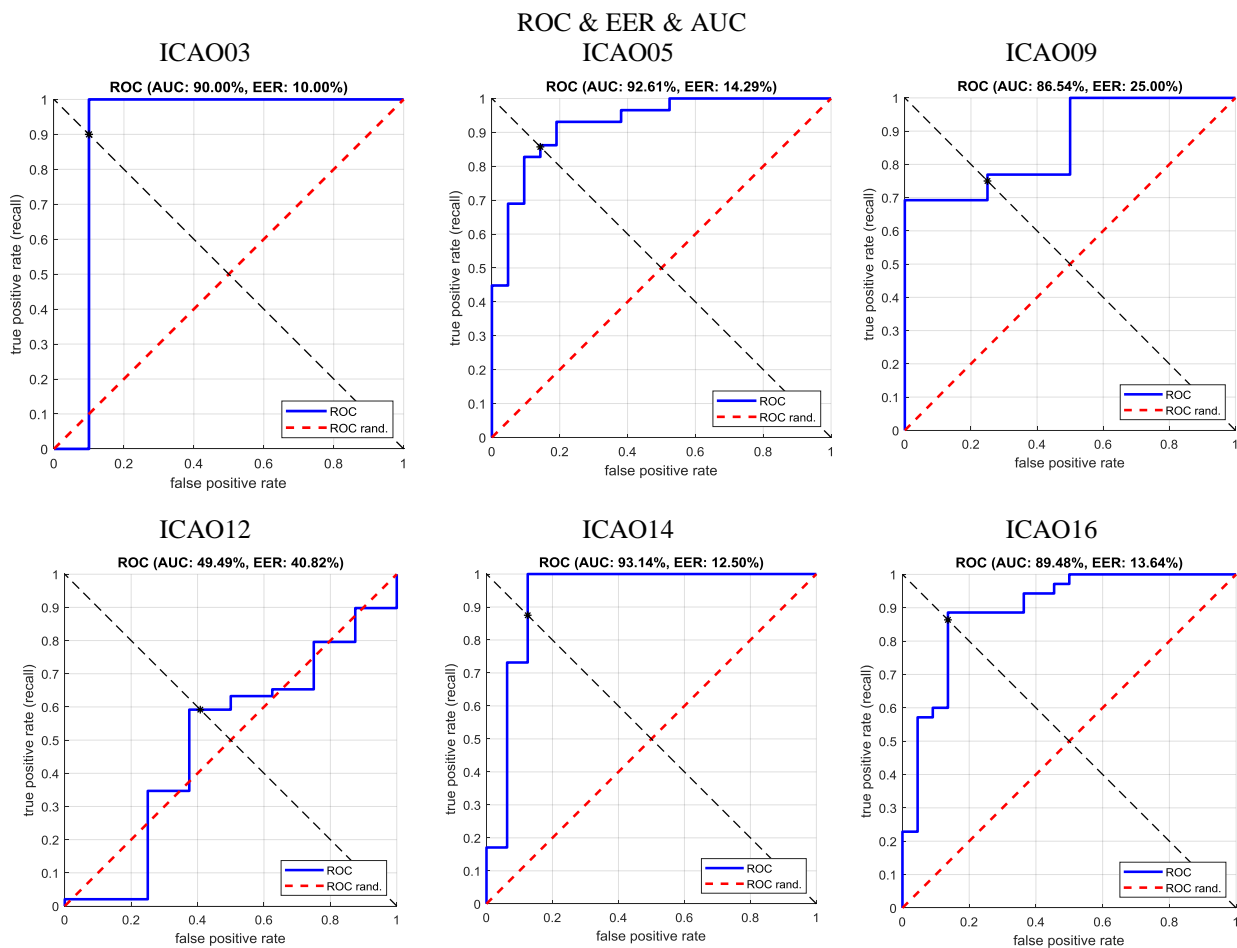
Table 3: EER and Rejection rates for the three SDKs and proposed method

Row	Characteristic	SDK1		SDK2		BioLabSDK		HMAX Method	
		EER	Rej	EER	Rej	EER	Rej	EER	Rej
1	ICAO03- Looking Away	27.5%	7.1%	-	-	20.6%	0.0%	10.00%	0.16%
2	ICAO05- Unnatural Skin Tone	18.7%	4.8%	50.0%	0.8%	4.0%	0.2%	14.29%	0.3%
3	ICAO09- Hair Across Eyes	50.0%	81.9%	-	-	12.8%	0.0%	25.00%	0.0%
4	ICAO12-Roll/Pitch/Yaw>8°	-	-	26.0%	2.9%	12.7%	0.2%	40.82%	0.0%
5	ICAO14- Red Eyes	5.2%	4.5%	34.2%	0.0%	7.4%	0.0%	12.5%	0.2%
6	ICAO16- Shadows Across Face	36.4%	8.1%	-	-	13.1%	0.4%	13.64%	0.4%
7	ICAO19- Frames Too Heavy	-	-	-	-	5.8%	0.0%	0.0%	0.0%
8	ICAO20- Frame Covering Eyes	50.0%	62.3%	-	-	6.3%	0.0%	0.0%	0.1%
9	ICAO23- Mouth Open	3.3%	52.1%	-	-	6.2%	0.0%	10.71%	0.4%

- Shows that the SDK does not support the test for this Characteristic

**ROC Curve:** The System Performance Characteristic Curve or Receiver Operating Characteristic (ROC) is an objective evaluation method that is a two-dimensional diagram, where the x-axis corresponds to False Positive Rate or FMR and the y-axis corresponds to True Positive Rate or 1- FNMR (often replaces by FNMR). The

area under the ROC diagram represents the Area Under Curve (AUC). The high AUC value indicates higher accuracy of the model. For each of the ICAO requirements investigated in this research, the ROC diagrams are calculated and are shown in Figure 8.





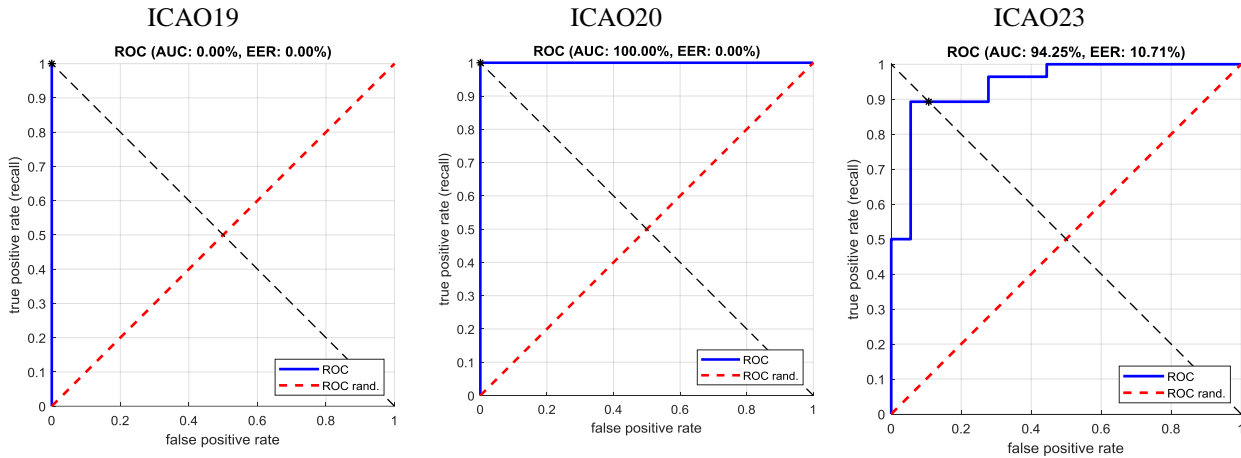


Fig. 8 ROC diagrams for each of the nine ICAO features investigated in this work.

In the ROC diagram, the greater the accuracy of the test, the closer the curve to the left boundary and then to the upper boundary of the ROC space. The closer the bend to the 45-degree diameter of the ROC space, the less accurate the test. The tangent slope at a cutting point shows the Likelihood Rate (LR) for the test.

**Recall:** A positive diagnosis probability being true when the actual results are positive. This parameter is also called the true positive rate.

$$Recall = TPR = \frac{TP}{TP+FN} \tag{10}$$

**Precision:** A positive diagnosis probability being true when the experimental results are positive. This parameter is also called the false negative rate.

$$Precision = FNR = \frac{TP}{TP+FP} \tag{11}$$

**Average Precision:** For each positive recall sample, the sum of precision values in positive recalls to all positive diagnoses, determines the average Precision parameter.

$$AP = \frac{\sum_{k=1}^n (P(k) \times rel(k))}{\text{number of compliant requirements}} \tag{12}$$

where  $rel(k)$  is an index that is equal to 1 if the requirement is compliant, otherwise it will be zero [34]. The average contains all the associated requirements and those associated requirements that have not been compliant, have a zero value for precision.

**AP11:** This is an index calculated by averaging the precision over a set of evenly spaced recall levels  $\{0, 0.1, 0.2, \dots, 1.0\}$ . This factor is used for reducing the impact of wiggles in the curve.

$$AP11 = \frac{1}{11} \sum_{r \in \{0, 0.1, \dots, 1.0\}} p_{int}(r) \tag{13}$$

where  $p_{int}(r)$  is the interpolated precision, shows the maximum precision over all recalls greater than  $r$  (in 11 points) [35].

Figure 9 contains the accuracy-recall diagrams, the area under the curve (AUC), the average precision (AP) and AP11.

PRECISION & RECALL & AUC & AVERAGE PRECISION (AP)

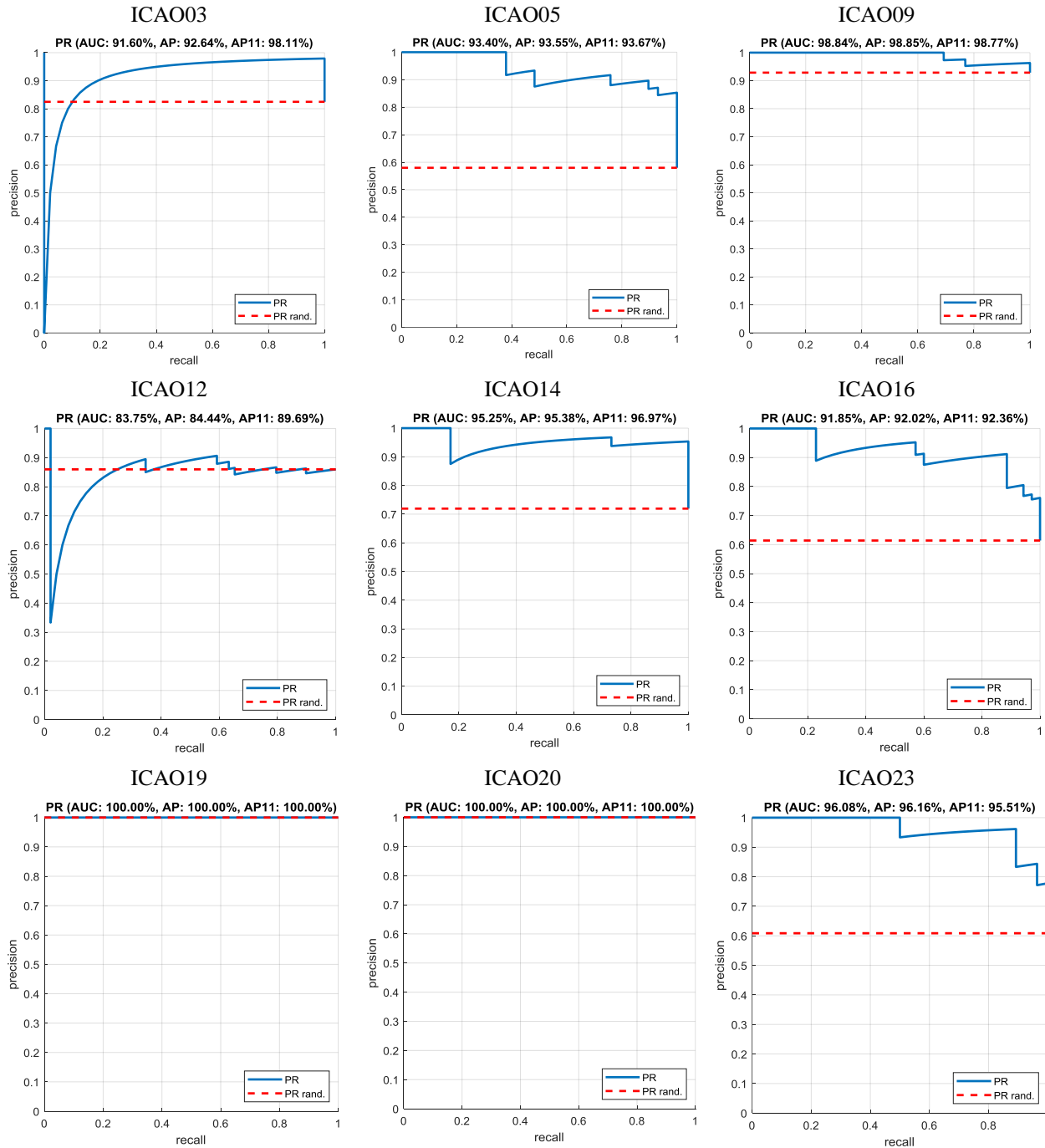


Fig.9 Accuracy-recall diagrams, Area under the Curve (AUC) and Average Precision (AP) for the selected ICAO requirements

## 6- Conclusions

One of the most important factors affecting the accuracy of automatic face recognition is the face images quality assessment. Accordingly, a benchmark was provided by the University of Bologna called BIOLAB-ICAO, which Facial

image Compliancy part is called FICV. In this research we proposed a new approach for facial images quality assessment using HMAX model (as the perceptual brain modeling).

The way of information storing and fetching it for training, is like the way of storing information in the brain. Nine

ICAO requirements are used to assess quality. The AR and PUT databases were used to train and test the model. The assessment factors introduced in the FICV benchmark were used to evaluate the modeling results. The results showed improvement in the detection of some requirements, particularly Frame Too Heavy (ICAO19), Frame Across Eyes (ICAO20). So, it is recommended to use HMAX model for these requirements detecting in the SDKs. As a follow-up, a model based on brain decision-making paths approaches to assess the quality of facial images, can be suggested.

## References

- [1] ISO/IEC 19794-5, Information technology - Biometric data interchange formats - Part 5: Face image data, 2011.
- [2] Y. Wong, Sh. Chen, S. Mau, C. Sanderson, B. C. Lovell, "Patch-based Probabilistic Image Quality Assessment for Face Selection and Improved Video-based Face Recognition", Computer Society Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), 2011, IEEE.
- [3] P. Ji, Y. Fang, Zh. Zhou, and J. Zhu, "Fusion of mSSIM and SVM for Reduced-Reference Facial Image Quality Assessment", (Eds.): CCBR 2012, LNCS 7701, pp. 75–82, Springer-Verlag Berlin Heidelberg.
- [4] M. Ferrara, A. Franco, D. Maio, "BIOLAB-ICAO: A New Benchmark to Evaluate Applications Assessing Face Image Compliance to ISO/IEC 19794-5 Standard", 978-1-4244-5654-3/09/IEEE ICIP 2009, pp. 41-44.
- [5] Th. Serre, L. Wolf, S. Bileschi, M. Riesenhuber, and T. Poggio, "Robust object recognition with cortex-like mechanisms," IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 29, 2007, pp. 411–426.
- [6] R. Hsu; M. Abdel-Mottaleb, A.K. Jain, "Face detection in color images," IEEE Transactions on Pattern Analysis and Machine Intelligence, May 2002, Vol.24, No.5, pp.696-706.
- [7] C. Bouvier, A. Benoit, A. Caplier, P. Y. Coulon "Open Or Closed Mouth State Detection: Static Supervised Classification Based On Log-Polar Signature", Springer Berlin Heidelberg on Advanced Concepts for Intelligent Vision Systems, 2008, pp. 1093-1102.
- [8] W. Qi, Y. Sheng, L. Xianwei, "A Fast Mouth Detection Algorithm Based On Face Organs" IEEE International Conference on Power Electronics and Intelligent Transportation System, December 2009, pp. 250-252.
- [9] M. Ferrara, A. Franco, D. Maio, D. Maltoni, "Face Image Conformance to ISO/ICAO Standards in Machine Readable Travel Documents", IEEE Transactions on Information Forensics and Security, AUGUST 2012, Vol 7, No. 4, pp. 1204-1213.
- [10] T.H.B Nguyen, V.H. Nguyen, and H. Kim, "Automated conformance testing for ISO/IEC 19794-5 Standard on facial photo specifications", Int. J. Biometrics, Vol. 5, No. 1, 2013, pp.73–98.
- [11] S. Coronel Castellanos, I. Solis Moreno, J. A. Cantoral Ceballos, R. Alvarez Vargas, P. L. Martinez Quintal, "An Approach to Improve Mouth-State Detection to Support the ICAO Biometric Standard for Face Image Validation", International Conference on Mechatronics, Electronics and Automotive Engineering, 978-1-4673-8329-5/15, 2015 IEEE, DOI 10.1109/ICMEAE.2015.12
- [12] R. L. Parente, L. V. Batista, Igor L. P. Andrezza, Erick V. C. L. Borges, Rajiv A. T. Mota, "Assessing Facial Image Accordance to ISO/ICAO Requirements", 29th SIBGRAPI Conference on Graphics, Patterns and Images (SIBGRAPI), Sao Paulo, Brazil October 2016, IEEE, DOI: 10.1109/SIBGRAPI.2016.033
- [13] Y. Li, W. Wu, B. Zhang and F. Li, "Enhanced HMAX model with feedforward feature learning for multiclass categorization". Front. Comput. Neurosci, 2015, Vol.9, pp. 123. DOI: 10.3389/fncom.2015.00123
- [14] HH. Gorji, S. Zabbah, R. Ebrahimpour, "A temporal neural network model for object recognition using a biologically plausible decision making layer", 2018, arXiv preprint arXiv:1806.09334.
- [15] J. Filali, H. Zghal, J. Martinet. "Ontology and HMAX Features-based Image Classification using Merged Classifiers". International Conference on Computer Vision Theory and Applications 2019 (VISAPP'19), Feb 2019, Prague, Czech Republic. hal-02057494.
- [16] H. Yu, Z. Xu, Ch. Fu and Y. Wang, "Bottom-up attention based on C1 features of HMAX model", Proc. SPIE 8558, Optoelectronic Imaging and Multimedia Technology II, 85580W (21 November 2012), <https://doi.org/10.1117/12.999263>.
- [17] S. Saraf Esmaili, K. Maghooli & A. Motie Nasrabadi, "A new model for face detection in cluttered backgrounds using saliency map and C2 texture features", International Journal of Computers and Applications, Vol. 40, No. 4, 2018, pp. 214-222.
- [18] H.Zh. Zhang, Y.F. Lu, T. K. Kang and M.T. Lim, "B-HMAX: A Fast Binary Biologically Inspired Model for Object Recognition", Neurocomputing, Vol. 218, 19 December 2016, pp. 242-250. <http://dx.doi.org/10.1016/j.neucom.2016.08.051>
- [19] C. Th'eriault, N. Thome, and M. Cord, "Extended coding and pooling in the HMAX model", IEEE

- Transactions on Image Processing , Vol. 22 , No. 2 , Feb. 2013 , pp. 764 – 777.
- [20] E.T. Rolls and G. Deco, "Computational neuroscience of vision", Press: Oxford, 1st edition, 2006.
- [21] J. Mutch and D.G. Lowe, "Object class recognition and localization using sparse features with limited receptive fields," *International Journal of Computer Vision*, vol. 80, October 2008, pp. 45–57.
- [22] I. Khan, H. Abdullah and M. Shamian Bin Zainal, "Efficient eyes and mouth detection algorithm using combination of viola jones and skin color pixel detection", *International Journal of Engineering and Applied Sciences*, Vol. 3, No. 4, June 2013, pp. 51-60.
- [23] M. Jones and P. Viola., "Face Recognition Using Boosted Local Features". Mitsubishi Electric Research Laboratories Technical Report Number: TR2003-25. Date: April, 2003.
- [24] NH. Barnouti, S. Sameer, "Face Detection and Recognition Using Viola-Jones with PCA-LDA and Square Euclidean Distance", (*IJACSA*) *International Journal of Advanced Computer Science and Applications*, Vol. 7, No. 5, 2016, pp.371-377.
- [25] S. Kaur, A. Chadha, "Supervised Descent Method Viola-Jones and Skin Color Based Face Detection and Tracking ", *International Journal of Engineering Development and Research (IJEDR)*, Volume 5, Issue 2, ISSN: 2321-9939, 2017, pp. 1689-1695.
- [26] I. Dagher and H. Al-Bazzaz, "Improving the Component-Based Face Recognition Using Enhanced Viola-Jones and Weighted Voting Technique", *Hindawi Modelling and Simulation in Engineering*, Volume 2019, Article ID 8234124, 2019, 9 pages. <https://doi.org/10.1155/2019/8234124>
- [27] Q. Cheng, H. Zhou, and J. Cheng, "The Fisher-Markov Selector: Fast Selecting Maximally Separable Feature Subset for Multiclass Classification with Applications to High-Dimensional Data", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 33, No. 6, June 2011, pp. 1217-1233.
- [28] V.N. Vapnik, "Statistical Learning Theory", Wiley, 1998.
- [29] B. Scholkopf and A.J. Smola, "Learning with Kernels", MIT Press, 2002.
- [30] M. Riesenhuber and T. Poggio, "Hierarchical models of object recognition in cortex," *Nature Neuroscience*, vol. 2, 1999, pp. 1019–1025.
- [31] A. Mittal, A.K. Moorthy, A.C. Bovik, "No-reference image quality assessment in the spatial domain". *IEEE Trans. Image Process.* Vol. 21, No. 12, 2012, pp. 4695–4708.
- [32] A. M. Martinez and R. Benavente, "The AR Face Database", CVC Technical Report No.24, June 1998.
- [33] A. Kasinski, A. Florek, A. Schmidt, "The PUT Face Database", *Image Processing & Communication*. Vol. 13. No. 3, 2008, pp. 59-64.
- [34] A. Turpin, F. Scholer, "User performance versus precision measures for simple search tasks", *Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (Seattle, WA, August 06–11, 2006)*. New York, NY: ACM. pp. 11–18.
- [35] K.H. Brodersen, C.S. Ong, K.E. Stephan, J.M. Buhmann, "The binormal assumption on precision-recall curves ", at the Wayback Machine. December 8, 2012, *Proceedings of the 20th International Conference on Pattern Recognition*, pp. 4263-4266.

**Azamossadat Nourbakhsh** received the B.S. degree in Computer Engineering from Azad University, Lahijan Branch, Iran in 1998, and M.S. degree in Artificial Intelligence & Robotics from Azad University, Science & Research Branch, Iran, in 2007. She is Ph.D. Candidate in Azad University, Science & Research Branch, Tehran, Iran. Her research interests include Image Processing, Machine Vision, Biometrics, Cognitive Science and Machine Intelligence.

**Mohammad. Shahram Moin** received his B.Sc. degree from Amir Kabir University of Technology, Tehran, Iran, in 1988; M.Sc. degree from University of Tehran, Iran, in 1991; and Ph.D. degree from École Polytechnique de Montréal, Montréal, Canada, in 2000, all degrees in electrical engineering. Dr. Moin is associate professor and head of IT Research Faculty in ICT Research Institute (ITRC). His research interests are Pattern Recognition, Image Processing, Biometrics, Data Mining and Big data Analytics.