

Rate Distortion Optimized Inter-View Frame Level Bit Allocation Method for MV-HEVC

Hui Yuan, *Member, IEEE*, Sam Kwong, *Fellow, IEEE*, Xu Wang, Wei Gao, and Yun Zhang, *Member, IEEE*

Abstract—In multi-view video coding, since inter-view prediction has been adopted as an important coding tool which could improve coding efficiency greatly, inter-view dependency is inevitable, i.e., the distortion of the reference view (RV) picture could be propagated to the non-reference view (NRV) pictures. Therefore, in order to achieve higher coding efficiency, the inter-view dependency must be taken into account for inter-view bit allocation. In this paper, the inter-view dependency is analyzed in detail, and a rate-distortion (RD) model for NRVs is derived by taking the distortion of RV into account. Based on the derived RD model, the inter-view bit allocation is represented as a mathematical problem with an analytic form, and is solved by a convex optimization (Lagrangian Multiplier) method. Experimental results demonstrate that the RD performance and the inter-view quality consistency of the proposed method is better than existing methods, while the complexity of the proposed method is comparable with the existing methods.

Index Terms—Bit allocation, inter-view dependency, MV-HEVC, rate distortion optimization.

I. INTRODUCTION

Due to the rapid development of 3D multimedia technologies, e.g. IMAX movie theaters¹ and low cost personal 3D displays, 3D video (3DV) applications is becoming

Manuscript received December 28, 2014; revised May 24, 2015 and July 25, 2015; accepted September 02, 2015. Date of publication September 10, 2015; date of current version November 13, 2015. This work was supported in part by the National Natural Science Foundation of China under Grant 61201211, Grant 61571274, Grant 61272289, and Grant 61501299, in part by the Young Scholars Program of Shandong University under Grant 2015WLJH39, in part by the City University of Hong Kong Applied Research Grant 9667094, in part by the Ph.D. Programs Foundation, Ministry of Education of China under Grant 20120131120032, in part by the Excellent Youth Scientist Award Foundation of Shandong Province under Grant BS2012DX021, in part by the Key Laboratory of Wireless Sensor Network and Communication, Chinese Academy of Sciences under Grant 2013002, and in part by the City University of Hong Kong Shenzhen Research Institute. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Adrian Munteanu.

H. Yuan is with the School of Information Science and Engineering, Shandong University, Ji'nan 250100, China, with the Key Laboratory of Wireless Sensor Network & Communication, Shanghai Institute of Microsystem and Information Technology, Chinese Academy of Sciences, Shanghai 200050, China, and with the Department of Computer Science, City University of Hong Kong, Hong Kong (e-mail: yuanhui0325@gmail.com; huiyuan@sdu.edu.cn).

S. Kwong is with the Department of Computer Science, City University of Hong Kong, Hong Kong, and also with the City University of Hong Kong Shenzhen Research Institute, Shenzhen 5180057, China (e-mail: cssamk@cityu.edu.hk).

X. Wang is with the College of Computer Science and Software Engineering, Shenzhen University, Shenzhen 518060, China (e-mail: wangxu@szu.edu.cn).

W. Gao is with the Department of Computer Science, City University of Hong Kong, Hong Kong (e-mail: gaowei262@126.com).

Y. Zhang is with the Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen 518055, China (e-mail: yun.zhang@siat.ac.cn).

Digital Object Identifier 10.1109/TMM.2015.2477682

¹“The IMAX experience,” [Online]. Available: <http://www.imax.com/about/experience/3d/>

more and more popular for both movie industry and home entertainment. In 3DV applications, 2 pictures, i.e. the left and the right view pictures, which correspond to the left and the right eyes of an observer, are projected onto a special-coated screen simultaneously. The observer can use 3D glasses to separate the 2 pictures so as to let them be projected into his corresponding eyes. Then, the human brain will fuse the 2 pictures together and enjoy a depth perception. To support the free view function of 3D applications, a set of view pictures (at least 2 view pictures) should be captured/generated, compressed, and transmitted to the 3DV application terminal. Thus, compared with the traditional 2D videos, more storage space and network bandwidth are needed for 3DVs. Due to the limited storage capacity and network bandwidth, efficient compression of 3DVs becomes critical and important for 3DV applications.

In order to compress 3DVs efficiently, Multi-view Video Coding (MVC) [1] was developed as an important extension of H.264/AVC [2]. Unlike the traditional method, the MVC exploits the inter-view correlation [3] and encodes multiview pictures simultaneously. Recently, a Joint Collaborative Team for 3DV (JCT-3 V), which was established by Moving Pictures Experts Group (MPEG) of International Organization for Standardization (ISO) and Video Coding Experts Group (VCEG) of International Telecommunication Union (ITU), developed a new MVC encoder based on the state of the art High Efficiency Video Coding (HEVC) Standard [4] in order to further improve the coding efficiency. Accordingly, the new MVC standard is named as MV-HEVC [5]. In the current MV-HEVC encoder common test condition [6], only 2 and 3 view cases are considered. Based on the MV-HEVC codec, a full 3D-HEVC standard [7] was also developed to generate more views, in which depth maps have been taken into considerations.

In 3DV coding, bit allocation is very important to improve the rate distortion (RD) performance. The bit allocation problem of 3DVs can be divided as 2 steps, i.e., bit allocation for texture videos and depth maps (e.g. Wang's algorithm [8], Cheung's algorithm [9], and our previous work [10]), and bit allocation for multiview (only texture videos are considered) pictures. In this paper, the bit allocation for multiview pictures is discussed based on MV-HEVC.

In MV-HEVC, inter-view prediction [11] based on disparity compensation is adopted as an important coding tool. Recently, lots of researchers focused on developing novel inter-view prediction methods [12]–[14] so as to improve the coding efficiency. Besides, for the same rate constraint, different inter-view bit allocation methods may give different average qualities of reconstructed view pictures. Yao *et al.* [15] proposed a joint bit allocation scheme for multiple stereoscopic 3DVs. In this method, the reference view (RV) picture is encoded

independently which is the same with MV-HEVC encoder, but the non-reference view (NRV) pictures are encoded as an enhancement layer which is compatible with scalable video codec. Then, a fast recursive golden-section search method is employed to allocate bits for different view pictures. Although a sub-optimal result could be achieved, multi-pass coding which is time-consuming must be used so as to generate enough search points. *Chang et al.* [16] proposed a frame level joint rate control scheme for MVC. In this method, the bit allocation between the left and the right view is converted as an optimization problem. However, the inter-view dependency is not considered in the optimization model. Moreover, in [17]–[20], a predefined fixed bit ratio is used for inter-view bit allocation; while for different 3DVs, the predefined fixed bit ratios are different. In the current MV-HEVC codec, an empirically optimal bit ratio is also adopted for inter-view bit allocation. For 2 view case, the empirically optimal bit ratio is 0.8 for RV picture, and 0.2 for NRV picture; while for 3 view case, the empirically optimal bit ratio is 0.66 for RV picture, and 0.17 for each of the 2 NRV pictures [21]. Besides, it should be noted that only sequence level bit allocation method were proposed in [17]–[21]. Motivated by the above analysis, a frame level inter-view bit allocation method based on inter-view dependency analysis is proposed to improve the coding efficiency for MV-HEVC.

In order to design a RD optimal bit allocation method for multiview pictures, accurate RD model must be obtained first. In the early studies, the RD model was derived based on the assumption that the transformed coefficients have Gaussian distributions [22]. Later, more accurate RD models were derived by considering that the transformed coefficients have Laplacian [23], [24], Generalized Gaussian [25], and Cauchy [26] distributions. Among them, the Cauchy distribution based RD model is with the best accuracy. Therefore Cauchy distribution based RD model is used in this paper. Exactly speaking, there are mainly 3 contributions, i.e.:

- a simplified Cauchy density [26] based RD model is proposed to describe the RD characteristic for RV picture;
- a novel RD model for NRV pictures is proposed based on inter-view dependency analysis; and
- an optimal analytic solution and the corresponding one-pass coding scheme for inter-view bit allocation.

The remainder of this paper is organized as follows. In Section II, a simplified Cauchy density based RD model is derived to describe the RD characteristic for RV pictures; besides, based on the inter-view dependency analysis, the RD model for NRV pictures is derived as well. In Section III, the inter-view bit allocation problem is modeled as a convex optimization problem by taking different temporal levels (TLs) into account, and is solved by Lagrangian Multiplier Method. Experimental results and conclusions are given in Sections IV and V, respectively.

II. RATE DISTORTION ANALYSIS

In H.264/MVC or MV-HEVC, the coding tools of RV pictures are the same with those of single view video coding; while for NRV pictures, inter-view prediction is used as an additional coding tool by taking the RV picture as a reference picture so as to achieve higher coding efficiency. In order to model the inter-view bit allocation problem as an analytical form, the rate

distortion characteristics of RV and NRV pictures are analyzed as follows.

A. Simplified Cauchy Model Derivation

The aim of video coding and transmission is to provide the best reconstruction quality under a certain network bandwidth. Therefore, it is no doubt that seeking the best tradeoff between rate and distortion is very important. In order to find the best tradeoff, the relationship between rate and distortion, i.e., rate-distortion (RD) model must be confirmed first.

In the Cauchy distribution based RD model, the relationship between distortion (D , measured by mean squared error, i.e., MSE) and rate (R , measured by average coding bits per pixel) could be written as

$$D(R) = \alpha \cdot R^{-\beta} \quad (1)$$

where $\alpha, \beta > 0$ are model parameters. However, the Cauchy based RD model cannot be applied easily because the power β is an uncertainty parameter. For easier application, the model (1) can be simplified by using the Taylor Series² expansion at $R = 0$. Therefore, we have

$$\begin{aligned} \frac{1}{D} &= \frac{1}{\alpha} R^\beta = \sum_{n=0}^{\infty} \frac{f^{(n)}(0)}{n!} (R - 0)^n \\ &= \frac{1}{\alpha} \cdot \left[\beta \cdot R + \frac{1}{2!} \beta(\beta - 1) R^2 + \dots \right] \\ &= \frac{\beta}{\alpha} R + \frac{1}{2!} \frac{\beta(\beta - 1)}{\alpha} R^2 + \dots \end{aligned} \quad (2)$$

Usually, the average coding bits of all pixels is far less than 1, therefore, (2) could be simplified as

$$\frac{1}{D} = \frac{\beta}{\alpha} R + O(R) \approx \frac{\beta}{\alpha} R + \varepsilon \quad (3)$$

where $O(R)$ represents the high order infinitesimals of R , which is approximated as a constant value ε . Thus, for easy illustration, (3) could be written as

$$D = (a \cdot R + b)^{-1}. \quad (4)$$

Where a (corresponds to β/α) and b (corresponds to ε) are model parameters.

Therefore, the RD model of RV pictures could be represented as

$$D_{RV} = (a_{RV} \cdot R_{RV} + b_{RV})^{-1} \quad (5)$$

in which D_{RV} and R_{RV} represent the distortion and average coding bits of RV pictures, a_{RV} and b_{RV} are parameters. Fig. 1 shows the relationship between $1/D_{RV}$ and R_{RV} , and the curves fitted by the proposed model and Laplacian based model [25]. Besides, Table I shows the detailed quantitative comparison between the proposed model and the Laplacian based model. From Fig. 1 and Table I, we can observe that the propose model is more accurate than the Laplacian based model.

B. Inter-View Distortion Dependency Analysis

In multi-view video coding, a NRV picture of the current frame could be predicted (called as disparity compensation prediction, DCP [27]) from the corresponding RV picture, as shown

²[Online]. Available: http://en.wikipedia.org/wiki/Taylor_Series

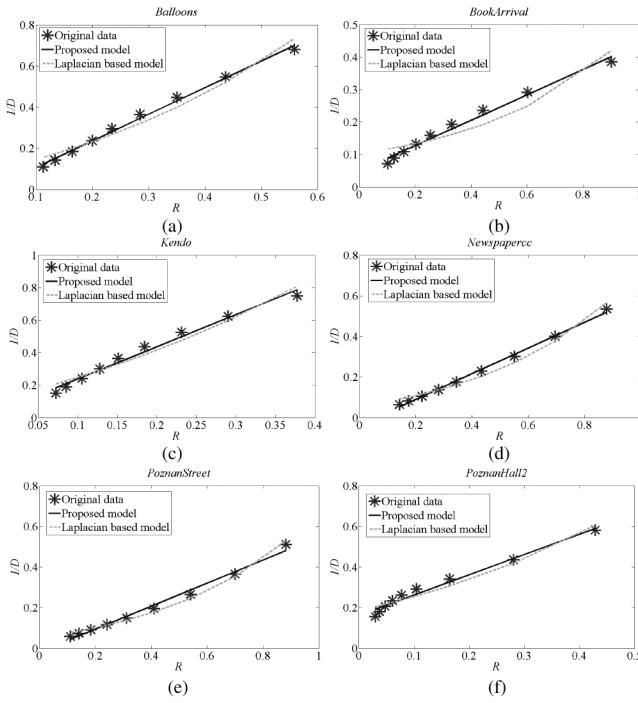


Fig. 1. Comparison between the proposed model and the Laplacian-based R-D model.

TABLE I

ROOT OF MEAN SQUARED ERROR (RMSE) AND CORRELATION COEFFICIENTS COMPARISONS (CC)

Sequence	Proposed Model		Laplacian Based Model	
	RMSE	CC	RMSE	CC
Balloons	0.0129	0.9975	0.0361	0.9803
BookArrival	0.0115	0.9930	0.0325	0.9430
Kendo	0.0245	0.9918	0.0411	0.9769
Newspapercc	0.0085	0.9984	0.0227	0.9886
PoznanStreet	0.0161	0.9937	0.0143	0.9950
PoznanHall2	0.0189	0.9892	0.0282	0.9757

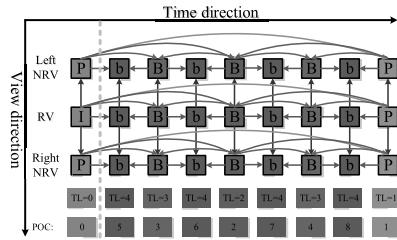


Fig. 2. Coding structure example of MV-HEVC.

in Fig. 2. During the DCP procedure, when *interview SKIP mode* is selected, the matching Prediction Units (PUs) in the RV picture will be directly copied to the current PU in the NRV picture. Thus, the distortion of the matching PU in the RV picture will also be propagated to the current PU in the NRV picture directly. Accordingly, for those PUs with interview SKIP mode, the distortion could be written as

$$D_{NRV,SKIP} = \eta(D_{RV} + \delta) \quad (6)$$

where D_{RV} is the distortion of the RV picture, η is the percentage of interview SKIP mode which is predicted from RV

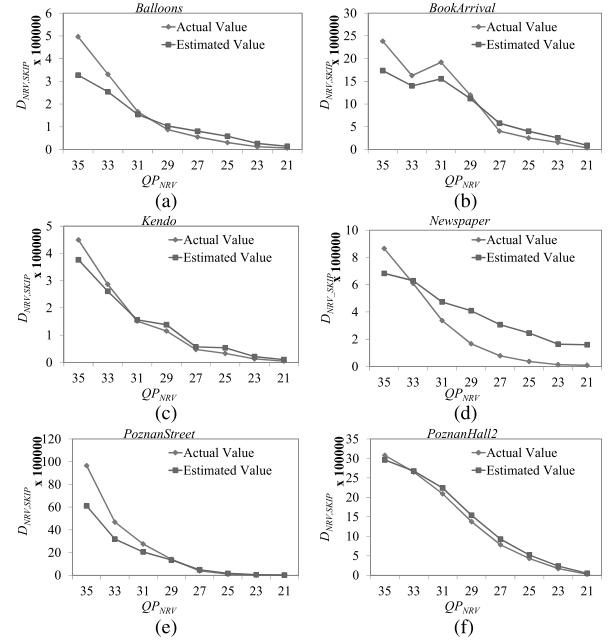


Fig. 3. Relationship between the actual $D_{NRV,SKIP}$ and the estimated one.

TABLE II
CC AND RMSE BETWEEN THE ACTUAL $D_{NRV,SKIP}$
AND THE ESTIMATED ONE

Sequence	CC	RMSE	Sequence	CC	RMSE
Balloons	0.986	0.086	Newspaper	0.9773	0.226
BookArrival	0.985	0.292	PoznanStreet	0.995	0.662
Kendo	0.995	0.038	PoznanHall2	0.996	0.053

picture, δ is determined by wrong estimated disparity, non-verified block-based correlation assumption, and the inherent luminance discrepancy between the RV and NRV pictures, etc.

The relationships between the actual $D_{NRV,SKIP}$ and the estimated $D_{NRV,SKIP}$ by using (6) are shown in Fig. 3, in which D_{RV} is obtained by encoding the RV picture with a quantization parameter (QP) of 28 (denoted as QP_{RV}); $D_{NRV,SKIP}$ is obtained by encoding the NRV picture with variable QPs (i.e. 21, 23, 25, 27, 29, 31, 33, and 35, denoted as QP_{NRV}); both D_{RV} and $D_{NRV,SKIP}$ are normalized by the total number of pixels in a picture (this is why the distortions are small). For easy statistics, only the first frame was encoded. From Fig. 3, we can observe that the actual $D_{NRV,SKIP}$ and the estimated one are similar for all the test video sequences. The estimation errors are caused by δ which is usually non-zero and could also be affected by the compression procedure. For example, since the interview discrepancy of the sequence *Newspapercc* is larger than other sequences, the estimation of $D_{NRV,SKIP}$ is not accurate enough compared with other sequences. Besides, Table II shows the correlation coefficients (CC) and the root of mean squared estimation errors (RMSE) between the actual $D_{NRV,SKIP}$ and the estimated one, from which it could be concluded that (6) is accurate to some extent.

When other (non-SKIP) prediction modes are selected, the matching PU in the RV picture is only used as a prediction of the current PU, and the residual signal is then obtained by subtracting it from the original signal of the current PU. Let

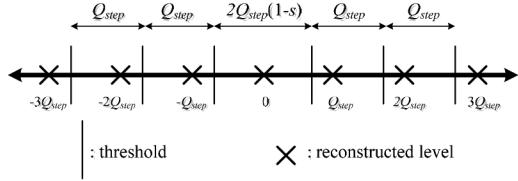


Fig. 4. Quantization example.

$\mathbf{L}_{\text{NRV},\text{orig}}$ denotes a transform unit (TU) of the current PU, $\mathbf{L}_{\text{NRV},\text{pred}}$ denotes the prediction (with a non-SKIP mode) of the TU. Then, the transform coefficients could be written as

$$\mathbf{W} = \mathbf{C} [\mathbf{L}_{\text{NRV},\text{orig}} - \mathbf{L}_{\text{NRV},\text{pred}}] \mathbf{C}^T \quad (7)$$

where \mathbf{W} is the transform coefficients matrix, \mathbf{C} is the transform basis matrix. During the quantization procedure, a transformed coefficient w (an element of \mathbf{W}) could be quantized as

$$z = \lfloor |w| / Q_{\text{step}} + s \rfloor \cdot \text{sgn}(w) \quad (8)$$

where z is the quantization level of w , Q_{step} is the quantization step size (determined by QP), s denotes rounding offset [28], $\lfloor \cdot \rfloor$ is the floor operation, and $\text{sgn}(\cdot)$ is a sign function which returns the sign of the input signal. From (7) and (8), we can conclude that when the distortion of the RV picture is changed, e.g., the QP of the RV picture is changed, $\mathbf{L}_{\text{NRV},\text{pred}}$ will also be changed; and then, the quantization level z will be changed as well. At the decoder, the quantization level z could be reconstructed to w_{rec} by the inverse quantization operation as shown in (9) and Fig. 4

$$w_{\text{rec}} = z \cdot Q_{\text{step}}. \quad (9)$$

It can be concluded from (9) and Fig. 4 that the w_{rec} may be changed when z is changed. However, no matter what the change is, the difference between w and w_{rec} will fall within the range of $[-(1-s)Q_{\text{step}}, (1-s)Q_{\text{step}}]$ which depends on Q_{step} . Furthermore, consider N level uniform quantizer, the quantized distortion could be calculated as

$$D = \int_{-(1-s)Q_{\text{step}}}^{(1-s)Q_{\text{step}}} (w)^2 f(w) dw + \sum_{n=1}^N \int_{(n-s)Q_{\text{step}}}^{(n-s+1)Q_{\text{step}}} (w - nQ_{\text{step}})^2 f(w) dw \quad (10)$$

from which we can observe that the quantized distortion depends on Q_{step} and the distribution function of w , i.e. $f(w)$. For a uniform quantizer, it is well known that the quantized distortion could be approximated as [29]

$$D = Q_{\text{step}}^2 / 12 \quad (11)$$

no matter what $f(w)$ is. This means that the distortion of RV picture will not affect the distortion of blocks with non-SKIP modes in NRV pictures. In order to verify this point, Fig. 5 shows the relationship between D_{RV} and the distortions of non-SKIP PUs in the NRV pictures (denoted as $D_{\text{NRV},\text{non-SKIP}}$). In Fig. 5, D_{RV} is obtained by encoding the RV picture with $QP = 20, 22, 24, 26, 28, 30, 32, 34$; while $D_{\text{NRV},\text{non-SKIP}}$ is obtained by encoding the NRV picture with

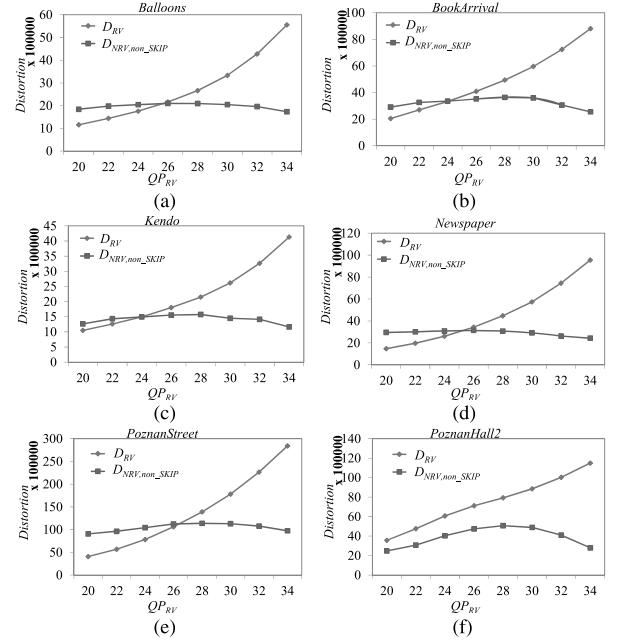
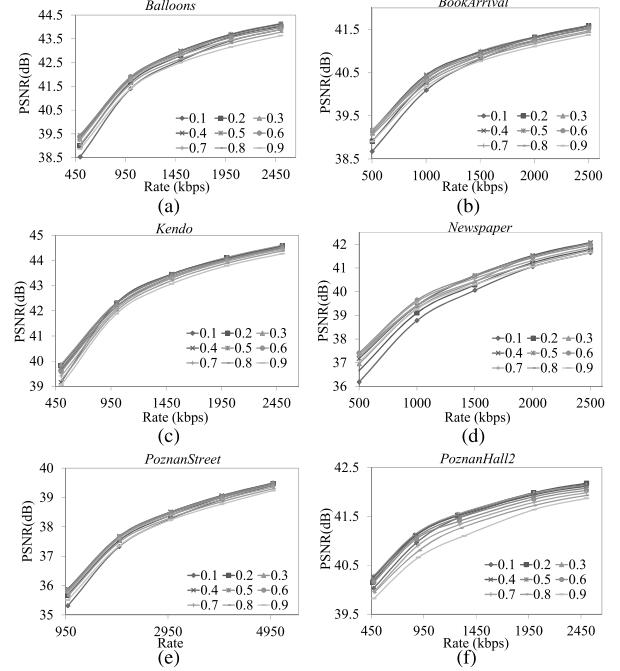
Fig. 5. Relationship between D_{RV} and $D_{\text{NRV},\text{non-SKIP}}$.

Fig. 6. RD performances under different initial bit ratio for 2 view case.

$QP=27$, from which we can observe that $D_{\text{NRV},\text{non-SKIP}}$ remain stable when D_{RV} is changed. Actually, the standard deviations of $D_{\text{NRV},\text{non-SKIP}}$ are only 0.163, 0.480, 0.179, 0.317, 0.425, and 0.484 for Fig. 5(a), (b), (c), (d), (e), and (f) respectively. Therefore, it could be concluded that the influence of D_{RV} on $D_{\text{NRV},\text{non-SKIP}}$ is small.

Since rate and distortion are one to one correspondence, the coding bit rate of the PU with non-SKIP modes is also mainly determined by the quantization process. Furthermore, since the residual coding bits of SKIP mode is 0, the coding bits of the NRV picture depends only on those non-SKIP PUs. Based on the simplified Cauchy rate distortion model shown in (5), the

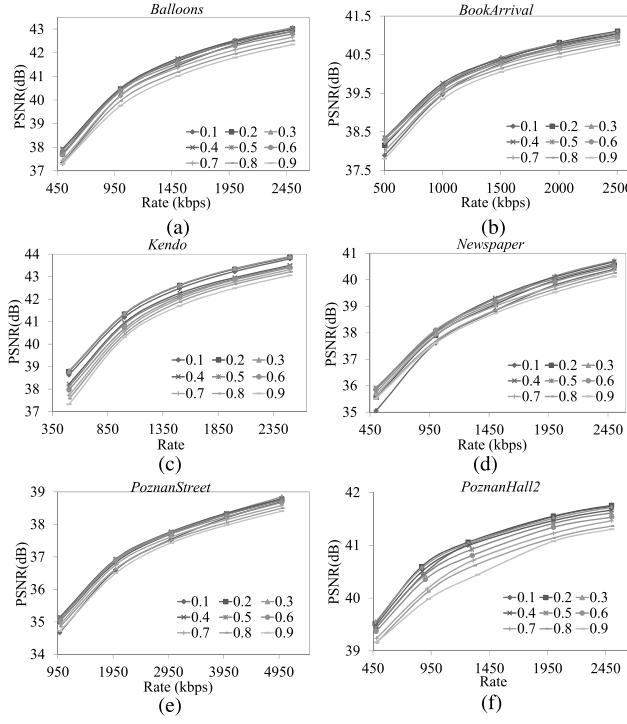


Fig. 7. RD performances under different initial bit ratio for 3 view case.

TABLE III
CODING CONFIGURATIONS

Encoder Parameters	Values	Encoder Parameters	Values
Unit Definition			
MaxCUWidth	64	FastSearch	1
MaxCUHeight	64	SearchRange	64
MaxPartitionDepth	4	BipredSearchRange	1
QuadtreeTULog2MaxSize	5	HadamardME	1
QuadtreeTULog2MinSize	3	FEN	1
QuadtreeTUMaxDepthInter	3	FDM	1
QuadtreeTUMaxDepthIntra	3		
Coding Structure			
IntraPeriod	-1	WaveFrontSyncro	0
DecodingRefreshType	1	UniformSpacingIdc	0
GOPSsize	8	PCMEnabledFlag	0
		SliceMode	0
Quantization			
MaxDeltaQP	0	DeblockingFilterControlPresent	1
MaxCuDQPDepth	0	LoopFilterOffsetInPPS	0
DeltaQpRD	0	LoopFilterDisable	0
RDOQ	1	LoopFilterBetaOffset_div2	0
RDOQTS	1	LoopFilterTcOffset_div2	0
ScalingList	0	DeblockingFilterMetric	0
Coding Tools			
SAO	1	IvMvPred	1
AMP	1	AdvMultiviewResPred	1
TransformSkip	1	IlliCompEnable	1
TransformSkipFast	1	ViewSynthesisPred	1
SAOLcuBoundary	0	DepthRefinement	1
Rate Control			
RateControl	1		
KeepHierarchicalBit	1		
LCULevelRateControl	1		
RCLCUSeparateModel	1		
InitialQP	0		
RCForceIntraQP	0		
Multiview Coding Tools			
		NA	

relationship between the distortion of non-SKIP PUs and the coding bits could be written as

$$D_{NRV,non_SKIP} = (\alpha_{NRV} R_{NRV} + \beta_{NRV})^{-1} \quad (12)$$

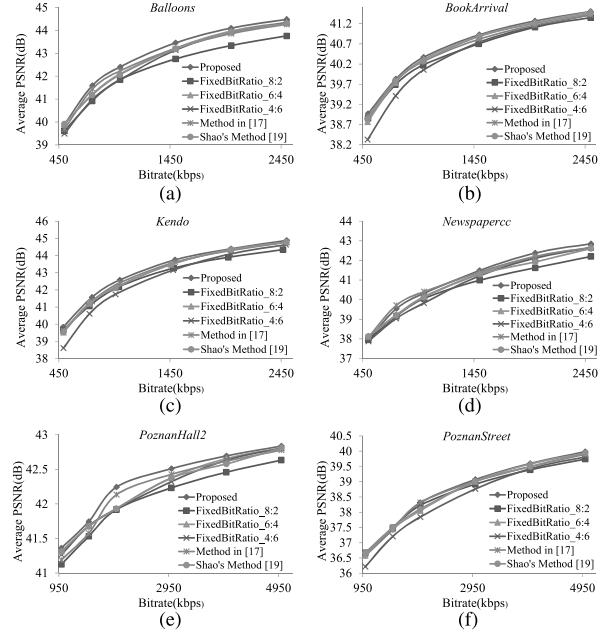


Fig. 8. RD curves comparison for 2 view case.

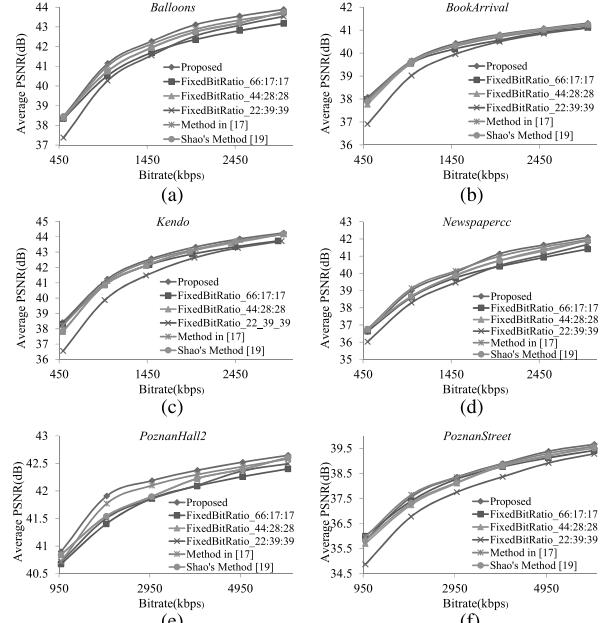


Fig. 9. RD curves comparison for 3 view case.

where R_{NRV} denotes the coding bits of the NRV picture, α_{NRV} and β_{NRV} are model parameters. Accordingly, the distortion of NRV picture could be written as

$$\begin{aligned} D_{NRV} &= D_{NRV,SKIP} + D_{NRV,non_SKIP} \\ &= \eta(D_{RV} + \delta) + (a_{NRV} R_{NRV} + b_{NRV})^{-1} \\ &= \eta(a_{RV} R_{RV} + b_{RV})^{-1} \\ &\quad + (a_{NRV} R_{NRV} + b_{NRV})^{-1} + \eta\delta. \end{aligned} \quad (13)$$

Moreover, for 3 view case of MV-HEVC, since the middle view picture is coded as the RV picture, both the left and the right view pictures are coded as NRV pictures, the distortion dependency between RV and NRV pictures is much larger than

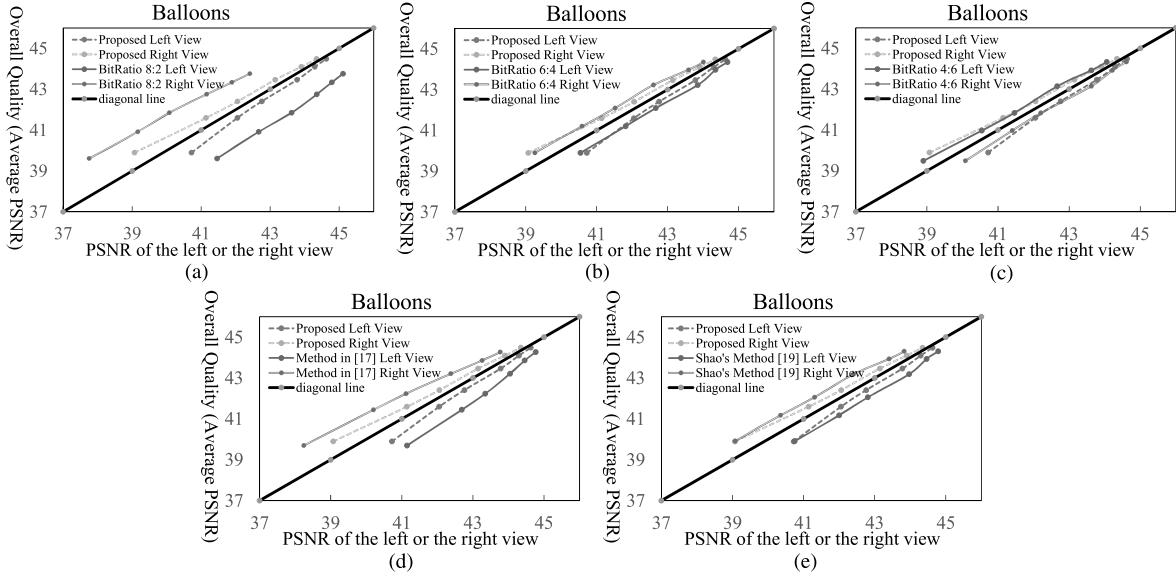


Fig. 10. Relationship between the overall quality and the distortion of each view of the *Balloons* sequence, 2 view case.

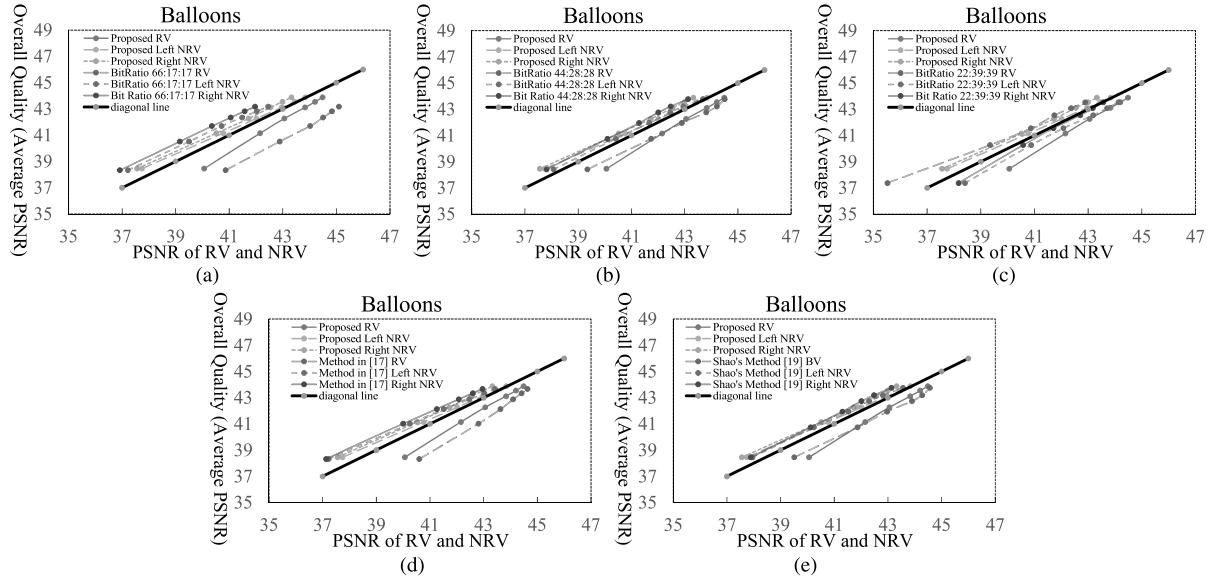


Fig. 11. Relationship between the overall quality and the distortion of each view of the *Balloons* sequence, 3 view case.

In the experiment, we vary initial bit ratio (i.e. 0.1:0.9, 0.2:0.8, ..., and 0.9:0.1 for 2 view case, and 0.1:0.45:0.45, 0.2:0.4:0.4, ..., and 0.9:0.05:0.05 for 3 view case) for the first 16 frames; while fixing the bit ratio for the remaining frames as 0.5:0.5 and 0.4:0.3:0.3 for 2 view and 3 view cases respectively. The total rate constraints are set as 500 kbps, 1000 kbps, 1500 kbps, 2000 kbps, and 2500 kbps. Then, the RD performances of different initial bit ratios are compared as shown in Fig. 6 and Fig. 7 (only the bit ratios of RV picture are labeled in the figures), from which we can observe that 0.4:0.6 and 30:35:35 are the quasi optimal initial bit ratios for 2 view and 3 view case respectively.

IV. EXPERIMENTAL RESULTS

To evaluate the performance of the proposed method, experiments were divided into 2 parts, i.e. bit allocation performance

comparison, and complexity comparison. 3DV sequences [31] which are adopted by JCT-3 V, i.e., *Balloons* (view 3 was encoded as RV, view 1 and 5 were encoded as NRV), *BookArrival* (view 8 was encoded as RV, view 10 and 6 were encoded as NRV), *Kendo* (view 3 was encoded as RV, view 1 and 5 were encoded as NRV), *Newspapercc* (view 4 was encoded as RV, view 2 and 6 were encoded as NRV), *PoznanHall2* (view 6 was encoded as RV, view 5 and 7 were encoded as NRV), and *PoznanStreet* (view 4 was encoded as RV, view 3 and 5 were encoded as NRV), were used. The MV-HEVC platform, 3D-HTM version 9.2 [32] was employed to encode those videos. The encoder configuration files were set as the same with that of common test condition [6] of MV-HEVC. Detailed encoder parameters are shown in Table III. Besides, the experiments were implemented on a computer with Intel Core i5-2300 CPU (2.8 GHz), and 4 GB memory size.

TABLE VIII
COMPLEXITY COMPARISONS FOR 3 VIEW CASE

Sequences	Encoding Time (s)										
	Proposed Method	Bit Ratio 66:17:17	CQ	Bit Ratio 44:28:28	CQ	Bit Ratio 22:39:39	CQ	Method in [17]	CQ	Shao's Method [19]	CQ
	2960.99	2937.18	1.01	2978.36	0.99	3014.11	0.98	2899.39	1.02	3423.23	0.86
Balloons	3242.51	3225.29	1.01	3282.86	0.99	3316.48	0.98	3249.39	1.00	3526.43	0.92
	3486.46	3478.37	1.00	3478.06	1.00	3527.29	0.99	3496.93	1.00	3571.04	0.98
	3615.41	3679.60	0.98	3638.00	0.99	3766.02	0.96	3649.70	0.99	3716.55	0.97
	3790.51	3722.70	1.02	3794.04	1.00	3910.52	0.97	3829.01	0.99	3868.10	0.98
	3925.15	3803.59	1.03	4016.96	0.98	4038.81	0.97	3966.88	0.99	4154.64	0.94
BookArrival	2982.29	2962.66	1.01	3000.00	0.99	3029.98	0.98	2976.67	1.00	3094.96	0.96
	3307.70	3297.45	1.00	3321.97	1.00	3379.31	0.98	3333.67	0.99	3440.63	0.96
	3564.14	3497.44	1.02	3572.97	1.00	3567.55	1.00	3604.74	0.99	3670.71	0.97
	3744.84	3691.48	1.01	3749.61	1.00	3765.82	0.99	3787.88	0.99	4005.09	0.94
	3977.11	3861.67	1.03	3909.29	1.02	3910.16	1.02	3942.97	1.01	3969.75	1.00
Kendo	4153.39	4001.59	1.04	4042.81	1.03	4064.75	1.02	4127.13	1.01	4124.39	1.01
	3145.23	3055.11	1.03	3087.13	1.02	3106.91	1.01	3044.25	1.03	3109.69	1.01
	3538.59	3390.33	1.04	3451.88	1.03	3482.94	1.02	3398.54	1.04	3515.06	1.01
	3702.11	3630.04	1.02	3722.44	0.99	3701.59	1.00	3656.51	1.01	3729.06	0.99
	3865.90	3830.18	1.01	3938.89	0.98	3887.26	0.99	3842.90	1.01	3874.24	1.00
Newspapercc	4026.69	4005.01	1.01	4070.27	0.99	4065.63	0.99	4039.52	1.00	4075.76	0.99
	4155.57	4078.25	1.02	4264.13	0.97	4166.55	1.00	4179.04	0.99	4237.81	0.98
	2773.30	2740.35	1.01	2858.87	0.97	2806.37	0.99	2761.01	1.00	2814.42	0.99
	3041.44	2951.95	1.03	3476.01	0.87	3014.70	1.01	2998.38	1.01	3015.12	1.01
	3216.71	3092.22	1.04	3486.27	0.92	3143.74	1.02	3179.10	1.01	3157.85	1.02
PoznanHall2	3306.37	3248.03	1.02	3325.50	0.99	3373.11	0.98	3320.58	1.00	3337.56	0.99
	3462.96	3365.83	1.03	3460.12	1.00	3445.20	1.01	3485.40	0.99	3543.66	0.98
	3571.21	3453.64	1.03	3667.31	0.97	3549.51	1.01	3551.76	1.01	3716.47	0.96
	8084.65	7972.16	1.01	8093.15	1.00	8244.51	0.98	8025.61	1.01	7946.21	1.02
	8749.43	8666.61	1.01	8931.09	0.98	8883.37	0.98	8687.21	1.01	8796.39	0.99
PoznanStreet	9232.94	9122.84	1.01	9740.91	0.95	9270.78	1.00	9119.29	1.01	9782.02	0.94
	9652.64	9382.23	1.03	9976.00	0.97	9695.82	1.00	9436.95	1.02	10088.62	0.96
	10079.77	9734.40	1.04	9853.29	1.02	9961.59	1.01	9744.54	1.03	10633.97	0.95
	10548.35	10055.86	1.05	10229.98	1.03	10186.47	1.04	10114.58	1.04	10923.57	0.97
	7381.88	7515.01	0.98	7732.35	0.95	7240.38	1.02	7453.42	0.99	7897.35	0.93
	7899.89	7768.00	1.02	8254.27	0.96	7833.64	1.01	7936.18	1.00	8379.31	0.94
	8316.25	8109.44	1.03	8556.62	0.97	8038.63	1.03	8374.50	0.99	8194.28	1.01
	8587.80	8380.97	1.02	8517.89	1.01	8619.30	1.00	8657.27	0.99	8442.26	1.02
	8946.48	8660.15	1.03	8734.04	1.02	8908.47	1.00	8889.92	1.01	8571.22	1.04
	9097.68	8869.35	1.03	9019.89	1.01	9069.82	1.00	9215.40	0.99	8719.20	1.04
	Average		1.02	0.99	1.00	1.00	1.00	1.00	0.98		

For 3 view case, the inter-view quality consistency ($IVQC_{3view}$) is defined as shown in (29), at the bottom of the page.

The inter-view consistency comparisons for 2 and 3 view cases are shown in Table VI. We can observe that, the average $IVQC_{2view}$ and $IVQC_{3view}$ of the proposed method are only 0.41 dB and 0.82 dB which are considerably smaller than other methods.

Furthermore, we have also given the relationship between the overall quality and the distortion of each view of *Balloon* sequence so as to analyze the interview quality consistency, as shown in Fig. 10 and Fig. 11. From the two figures, it can be observed that both the qualities of the left and the right view for 2 view case (or RV and NRVs for 3 view case) are closer to the overall quality (closer to the diagonal line) compared with other methods in most cases.

B. Complexity Comparison

In the experiments, the encoding time (T) was used to evaluate complexity. The encoding time of the proposed method was

used as benchmark. The ratio between the encoding times of other methods and that of the proposed method is employed to represent the complexity quotient (CQ)

$$CQ = T_{proposed} / T_{other_method} \quad (30)$$

where, T_{other_method} denotes the encoding time of one of the reference methods, $T_{proposed}$ denotes the encoding time of the proposed method. If CQ is larger than 1, the complexity of the proposed method (benchmark) is larger than the compared method; if CQ is smaller than 1, the complexity of the proposed method is smaller than the compared method. The complexity comparisons are shown in Table VII and Table VIII, from which we can observe that the CQs are close to 1, which means that the complexity of the proposed method is comparable with the existing methods. The reason why the complexity of the proposed method is a little higher or lower than the other methods is because of the random fluctuation of CPUs. Despite of the random fluctuation of CPUs, from Table VII and Table VIII, it can also be observed that the complexity of the proposed method is only a little larger than that of the Bit Ratio 8:2 and Bit Ratio 66:17:17

$$IVQC_{3view} = \frac{(|PSNR_{RV} - PSNR_{NRV_L}| + |PSNR_{RV} - PSNR_{NRV_R}|)}{2} \quad (29)$$

method for most of the test sequences and coding bit rates. Compared with other method, i.e., *Bit Ratio 6:4*, *Bit Ratio 44:28:28*, *Bit Ratio 4:6*, *Bit Ratio 22:39:39*, and *Shao's method*[18], the complexity of the proposed method is the least for almost all the test sequences and coding bit rates.

V. CONCLUSION

Inter-view dependency is analyzed in detail in this paper. Based on the analysis, the distortion propagation between RV picture and NRV pictures is mainly caused by the percentage of PUs with SKIP mode. Then, RD models for NRV pictures are derived. Subsequently, by employing the RD models, the inter-view bit allocation problem is formulated as a convex optimization problem, and was solved by Lagrangian Multiplier Method. Experimental results demonstrated that the RD performance of the proposed method is better than the existing method. Besides, the inter-view quality consistency of the proposed method is also better than the existing methods, and the complexity of the proposed method is comparable with other existing methods.

ACKNOWLEDGMENT

The authors would like to thank the editors and anonymous reviewers for their valuable comments. The authors would also like to thank the JCT-3 V Group for providing their 3D video sequences and their valuable work on 3DV.

REFERENCES

- [1] Y. Chen, Y.-K. Wang, K. Ugur, M. M. Hannuksela, J. Lainema, and M. Gabbouj, "The emerging MVC standard for 3D video services," *EURASIP J. Appl. Signal Process.*, vol. 2009, pp. 8:1–8:13, Jan. 2008.
- [2] A. Puri, X. Chen, and A. Luthra, "Video coding using the H.264/MPEG-4 AVC compression standard," *Signal Process. Image Commun.*, vol. 19, no. 9, pp. 793–849, Jun. 2004.
- [3] P. Merkle, A. Smolic, K. Müller, and T. Wiegand, "Efficient prediction structures for multiview video coding," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 17, no. 11, pp. 1461–1473, Nov. 2007.
- [4] G. J. Sullivan, J. R. Ohm, W. J. Han, and T. Wiegand, "Overview of the high efficiency video coding (HEVC) standard," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 22, no. 12, pp. 1649–1668, Dec. 2012.
- [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 (HEVC)," *IEEE J. Sel. Topics Signal Process.*, vol. 7, no. 6, pp. 1001–1016, Dec. 2013.
- [6] D. Rusanovskyy, K. Müller, and A. Vetro, "Common test conditions of 3DV core experiments," in *Proc. 3rd Meeting ITU-T/ISO/IEC Joint Collaborative Team 3D Video Coding (JCT-3 V)*, Jan. 2013, Doc. JCT3 V-C1100.
- [7] K. Müller, H. Schwarz, D. Marpe, C. Bartnik, S. Bosse, H. Brust, T. Hinz, H. Lakshman, P. Merkle, F. H. Rhee, G. Tech, M. Winken, and T. Wiegand, "3D high-efficiency video coding for multi-view video and depth data," *IEEE Trans. Image Process.*, vol. 22, no. 9, pp. 3366–3378, Sep. 2013.
- [8] Q. Wang, X. Ji, Q. Dai, and N. Zhang, "Free viewpoint video coding with rate-distortion analysis," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 22, no. 6, pp. 875–889, Jun. 2012.
- [9] G. Cheung, V. Velisavljević, and A. Ortega, "On dependent bit allocation for multiview image coding with depth-image-based rendering," *IEEE Trans. Image Process.*, vol. 20, no. 11, pp. 3179–3194, Nov. 2011.
- [10] H. Yuan, Y. Chang, J. Huo, F. Yang, and Z. Lu, "Model-based joint bit allocation between texture videos and depth maps for 3-D video coding," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 21, no. 4, pp. 485–497, Nov. 2011.
- [11] G. Tech, K. Wegner, Y. Chen, and M. Hannuksela, "MV-HEVC Draft Text 5," in *Proc. 5th Meeting ITU-T/ISO/IEC Joint Collaborative Team 3D Video Coding (JCT-3 V)*, Jul. 2013, Doc. JCT3 V-E1004.
- [12] Y. Chen, L. Zhang, V. Serigin, and Y.-K. Wang, "Motion hooks for the multiview extension of HEVC," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 24, no. 12, pp. 2090–2098, Dec. 2014.
- [13] D. B. Sansli, K. Ugurt, M. M. Hannukselat, and M. Gabbouj, "Inter view motion vector prediction in multiview HEVC," in *Proc. 3DTV-Conf.: True Vis.—Capture, Transmiss. Display 3D Video (3DTV-CON)*, Jul. 2014, pp. 1–4.
- [14] H. Schwarz and T. Wiegand, "Inter view prediction of motion data in multiview video coding," in *Proc. Picture Coding Symp.*, May 2012, pp. 101–104.
- [15] W. Yao, L. P. Chau, and S. Rahardja, "Joint rate allocation of stereoscopic 3D videos in next-generation broadcast applications," *IEEE Trans. Broadcast.*, vol. 59, no. 3, pp. 445–454, Sep. 2013.
- [16] Y. Chang and M. Kim, "A joint rate control scheme in a hybrid stereoscopic video codec system for 3DTV broadcasting," *IEEE Trans. Broadcast.*, vol. 59, no. 2, pp. 265–280, Jun. 2013.
- [17] Y. Liu, Q. Huang, S. Ma, D. Zhao, W. Gao, S. Ci, and H. Tang, "A novel rate control technique for multiview video plus depth based 3D video coding," *IEEE Trans. Broadcast.*, vol. 57, no. 2, pp. 562–572, Jun. 2011.
- [18] F. Shao, G. Jiang, W. Lin, M. Yu, and Q. Dai, "Joint bit allocation and rate control for coding multi-view video plus depth based 3D video," *IEEE Trans. Multimedia*, vol. 15, no. 8, pp. 1843–1854, Dec. 2013.
- [19] F. Shao, G. Jiang, and M. Yu *et al.*, "Asymmetric coding of multi-view video plus depth based 3D video for view rendering," *IEEE Trans. Multimedia*, vol. 14, no. 1, Feb. 2012.
- [20] F. Shao *et al.*, "Depth map coding for view synthesis based on distortion analyses," *IEEE J. Emerging Sel. Topics Circuits Syst.*, vol. 4, no. 1, pp. 106–117, Mar. 2014.
- [21] W. Lim, H. Jo, J. Yoo, D. Sim, and I. V. Bajić, "Inter-view MAD prediction for rate control of 3D multi-view video coding," in *Proc. 5th Meeting ITU-T/ISO/IEC Joint Collaborative Team 3D Video Coding (JCT-3 V)*, Aug. 2013, Doc. JCT3 V-E0227.
- [22] T. Wiegand and B. Girod, "Lagrange multiplier selection in hybrid video coder control," in *Proc. IEEE Int. Conf. Image Process.*, Oct. 2001, vol. 3, pp. 542–545.
- [23] E. Y. Lam and J. W. Goodman, "A mathematical analysis of the DCT coefficient distributions for images," *IEEE Trans. Image Process.*, vol. 9, no. 10, pp. 1661–1666, Oct. 2000.
- [24] H. J. Lee, T. Chiang, and Y. Q. Zhang, "Scalable rate control for MPEG-4 video," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 10, no. 6, pp. 878–894, Sep. 2000.
- [25] R. C. Reininger and J. D. Gibson, "Distributions of the two-dimensional DCT coefficients for images," *IEEE Trans. Commun.*, vol. COM-31, no. 6, pp. 835–839, Jun. 1983.
- [26] N. Kamaci, Y. Altunbasak, and R. M. Mersereau, "Frame bit allocation for the H.264/AVC video coder via cauchy-density-based rate and distortion models," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 15, no. 8, pp. 994–1006, Aug. 2005.
- [27] M. Flierl, A. Mavlnkar, and B. Girod, "Motion and disparity compensated coding for multiview video," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 17, no. 11, pp. 1474–1484, Nov. 2007.
- [28] Q. Xu, X. Lu, Y. Liu, and C. Gomila, "A fine rate control algorithm with adaptive rounding offsets (ARO)," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 19, no. 10, pp. 1424–1435, Nov. 2009.
- [29] R. M. Gray and D. L. Neuhoff, "Quantization," *IEEE Trans. Inf. Theory*, vol. 44, no. 6, pp. 2325–2383, Oct. 1998.
- [30] S. Boyd and L. Vandenberghe, *Convex Optimization*, 6th ed. Cambridge, U.K.: Cambridge Univ. Press, 2008, pp. 215–227.
- [31] Joint Collaborative Team for 3DV, 3D Video Test Sequences Oct. 2014 [Online]. Available: <ftp.hhi.fraunhofer.de>, Accessed on: Oct. 30, 2014
- [32] Joint Collaborative Team for 3DV, , Oct. 2014, 3D-HTM Software Platform [Online]. Available: https://hevc.hhi.fraunhofer.de/svn/svn_3DVCSoftware/tags/ Accessed on: Oct. 30, 2014
- [33] G. Bjontegaard and Tandberg, "Improvements of the BD-PSNR model," in *Proc. 35th Meeting ITU-T Video Coding Experts Group*, Jul. 2008, Doc. AI11.
- [34] Calculation of Average PSNR Differences Between RD Curves, ITU-T SG16/Q6 (VCEG), doc. VCEG-M33, Apr. 2001.



Hui Yuan (S'08–M'12) received the B.E. and Ph.D. degree in telecommunication engineering from Xidian University, Xi'an, China, in 2006 and 2011, respectively.

He was a Post-Doctoral Fellow with the Department of Computer Science, City University of Hong Kong, Hong Kong, from 2013 to 2014. He is currently an Associate Professor with the School of Information Science and Engineering, Shandong University, Jinan, China. His current research interests include video coding and multi-

media communication.



Sam Kwong (M'93–SM'04–F'13) received the B.S. degree from the State University of New York at Buffalo, Buffalo, NY, USA, in 1983, the M.S. degree in electrical engineering from the University of Waterloo, Waterloo, ON, Canada, in 1985, and the Ph.D. degree from the University of Hagen, Hagen, Germany, in 1996.

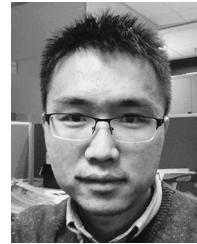
From 1985 to 1987, he was a Diagnostic Engineer with Control Data Canada, Mississauga, ON, Canada. He then joined Bell Northern Research Canada, Ottawa, ON, Canada, as a Member of Scientific Staff. In 1990, he became a Lecturer with the Department of Electronic Engineering, City University of Hong Kong, Hong Kong, where he is currently a Professor with the Department of Computer Science. His research interests are video and image coding and evolutionary algorithms.

Prof. Kwong serves as an Associate Editor of the IEEE TRANSACTIONS ON INDUSTRIAL ELECTRONICS and the IEEE TRANSACTIONS ON INDUSTRIAL INFORMATICS.



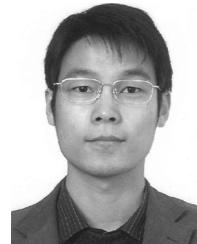
Xu Wang received the B.S. degree from South China Normal University, Guangzhou, China, in 2007, the M.S. degree from Ningbo University, Ningbo, China, in 2010, and the Ph.D. degree in computer science from The City University of Hong Kong, Hong Kong, in 2014.

He joined the College of Computer Science and Software Engineering, Shenzhen University, Shenzhen, China, in 2015, as an Assistant Professor. His research interests include video coding and stereoscopic image/video quality assessment.



Wei Gao received the M.S. degree in pattern recognition and intelligent systems from Huazhong University of Science and Technology, Wuhan, China, in 2012, and is currently working toward the Ph.D. degree in computer science at The City University of Hong Kong, Hong Kong.

His research interests include image and video processing, video coding and transmission, rate control, multimedia communication, visual perception, very-large-scale integration signal processing, and system-on-a-chip systems.



Yun Zhang (M'12) received the B.S. and M.S. degrees in electrical engineering from Ningbo University, Ningbo, China, in 2004 and 2007, respectively, and the Ph.D. degree in computer science from the Institute of Computing Technology, Chinese Academy of Sciences (CAS), Beijing, China, in 2010.

From 2009 to 2014, he was a Post-Doctoral Researcher with the Department of Computer Science, The City University of Hong Kong, Hong Kong. In 2010, he became an Assistant Professor with the Shenzhen Institutes of Advanced Technology, CAS,

where he has been an Associate Professor since 2012. His research interests include video compression, 3D video processing, and visual perception.