

Bruno Zatt · Muhammad Shafique
Sergio Bampi · Jörg Henkel

3D Video Coding for Embedded Devices

Energy Efficient Algorithms and
Architectures

 Springer

3D Video Coding for Embedded Devices

Bruno Zatt • Muhammad Shafique
Sergio Bampi • Jörg Henkel

3D Video Coding for Embedded Devices

Energy Efficient Algorithms
and Architectures

 Springer

Bruno Zatt
Department of Computer Science
Karlsruhe Institute of Technology
Karlsruhe, Germany

Institute of Informatics
Federal University of Rio
Grande do Sul (UFRGS)
Porto Alegre, Brazil

Sergio Bampi
Institute of Informatics
Federal University of Rio
Grande do Sul (UFRGS)
Porto Alegre, Brazil

Muhammad Shafique
Karlsruhe Institute of Technology
Karlsruhe, Germany

Jörg Henkel
Department of Computer Science
Karlsruhe Institute of Technology
Karlsruhe, Germany

ISBN 978-1-4614-6758-8 ISBN 978-1-4614-6759-5 (eBook)
DOI 10.1007/978-1-4614-6759-5
Springer New York Heidelberg Dordrecht London

Library of Congress Control Number: 2013934705

© Springer Science+Business Media New York 2013

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

Acknowledgements

This work was partly supported by the German Research Foundation (DFG) as part of the Transregional Collaborative Research Centre “Invasive Computing” (SFB/TR 89); <http://invasic.de>.

This work was also partly funded by the Brazilian Coordination for the Improvement of Higher Level Personnel (CAPES—“*Coordenação de Aperfeiçoamento de Pessoal de Nível Superior*”) and the Graduate Program on Microelectronics (PGMICRO), Institute of Informatics, Federal University of Rio Grande do Sul (UFRGS).

Contents

1	Introduction	1
1.1	3D-Video Applications.....	2
1.2	Requirements and Trends of 3D Multimedia	3
1.3	Overview on Multimedia Embedded Systems	5
1.4	Issues and Challenges.....	6
1.5	Monograph Contribution.....	7
1.5.1	3D-Neighborhood Correlation Analysis	7
1.5.2	Energy-Efficient MVC Algorithms	8
1.5.3	Energy-Efficient Hardware Architectures	9
1.6	Monograph Outline	9
2	Background and Related Works	11
2.1	2D/3D Digital Videos.....	11
2.2	Multiview Correlation Domains.....	14
2.2.1	Spatial Domain Correlation.....	14
2.2.2	Temporal Domain Correlation.....	15
2.2.3	Disparity Domain Correlation	16
2.3	Multiview Video Coding	16
2.3.1	MVC Encoding Process	18
2.3.2	Motion and Disparity Estimation	22
2.3.3	MVC Mode Decision	27
2.3.4	MVC Rate Control	28
2.4	3D-Video Systems.....	29
2.5	Multimedia Architectures Overview	30
2.5.1	Multimedia Processors/DSPs	30
2.5.2	Reconfigurable Processors for Video Processing.....	31
2.5.3	Application-Specific Integrated Circuits	32
2.5.4	Heterogeneous Multicore SoCs.....	33

2.6	Energy-Efficient Architectures for Multimedia Processing	33
2.6.1	Video Memories	34
2.6.2	SRAM Dynamic Voltage-Scaling Infrastructure.....	34
2.6.3	Dynamic Power Management for Memories.....	35
2.6.4	Energy Management for Multimedia Systems	36
2.6.5	Energy-Efficient Video Architectures	37
2.7	Energy/Power Consumption Background	38
2.8	Energy-Efficient Algorithms for Multiview Video Coding.....	40
2.8.1	Energy-Efficient Mode Decision	40
2.8.2	Energy-Efficient Motion and Disparity Estimation.....	42
2.9	Video Quality on Energy-Efficient Multiview Video Coding	45
2.9.1	Control Techniques Background	46
2.10	Summary of Background and Related Works	50
3	Multiview Video Coding Analysis for Energy and Quality	53
3.1	Energy Requirements for Multiview Video Coding.....	53
3.1.1	MVC Computational Effort.....	57
3.1.2	MVC Memory Access.....	59
3.1.3	Adaptivity in MVC Video Encoder	60
3.2	Energy-Related Challenges in Multiview Video Coding	62
3.3	Objective Quality Analysis for Multiview Video Coding	63
3.4	Quality-Related Challenges in Multiview Video Coding.....	65
3.5	Overview of Proposed Energy-Efficient Algorithms and Architectures for Multiview Video Coding	66
3.5.1	3D-Neighborhood.....	67
3.5.2	Energy-Efficient Algorithms	68
3.5.3	Energy-Efficient Architectures.....	69
3.6	Summary of Application Analysis for Energy and Quality	71
4	Energy-Efficient Algorithms for Multiview Video Coding.....	73
4.1	3D-Neighborhood Correlation Analysis	74
4.1.1	Coding Mode Correlation Analysis.....	74
4.1.2	Motion Correlation Analysis.....	82
4.1.3	Bitrate Correlation Analysis.....	84
4.2	Thresholds	87
4.3	Multilevel Mode Decision-based Complexity Adaptation	90
4.3.1	Multilevel Fast Mode Decision	90
4.3.2	Energy-Aware Complexity Adaptation	95
4.3.3	Multilevel Fast Mode Results.....	100
4.3.4	Energy-Aware Complexity Adaptation Results	105
4.4	Fast Motion and Disparity Estimation.....	107
4.4.1	Fast Motion and Disparity Estimation Algorithm.....	107
4.4.2	Fast ME/DE Algorithm Results	109

- 4.5 Video-Quality Management for Energy-Efficient Algorithms..... 111
 - 4.5.1 Hierarchical Rate Control for MVC..... 111
 - 4.5.2 Frame-Level Rate Control..... 113
 - 4.5.3 Basic Unit-Level Rate Control..... 119
 - 4.5.4 Hierarchical Rate Control Results..... 121
- 4.6 Summary of Energy-Efficient Algorithms for Multiview Video Coding 126
- 5 Energy-Efficient Architectures for Multiview Video Coding..... 127**
 - 5.1 Motion and Disparity Estimation Hardware Architecture 127
 - 5.1.1 SAD Calculator 130
 - 5.1.2 Programmable Search Control Unit 131
 - 5.1.3 On-Chip Video Memory..... 133
 - 5.1.4 Address Generation Unit..... 134
 - 5.2 Parallelism in the MVC Encoder and ME/DE Scheduling 136
 - 5.2.1 Parallelism in the MVC Encoder..... 136
 - 5.2.2 ME/DE Hardware Architecture Pipeline Scheduling..... 137
 - 5.3 Dynamic Search Window Formation 140
 - 5.3.1 ME/DE Memory Access Pattern Analysis 140
 - 5.3.2 Search Map Prediction 142
 - 5.3.3 Dynamic Search Window Formation 143
 - 5.4 On-Chip Video Memory..... 145
 - 5.4.1 On-Chip Memory Design..... 145
 - 5.4.2 Application-Aware Power Gating 146
 - 5.5 Hardware Architecture Evaluation 148
 - 5.5.1 Dynamic Window Formation Accuracy..... 148
 - 5.5.2 Hardware Architecture Evaluation..... 148
 - 5.6 Summary of Energy-Efficient Algorithms for Multiview Video Coding 150
- 6 Results and Comparison..... 151**
 - 6.1 Experimental Setup 151
 - 6.1.1 Software Simulation Environment 151
 - 6.1.2 Benchmark Video Sequences 152
 - 6.1.3 Fairness of Comparison..... 155
 - 6.1.4 Hardware Description and ASIC Synthesis 155
 - 6.2 Comparison with the State of the Art..... 156
 - 6.2.1 Energy-Efficient Algorithms 156
 - 6.2.2 Video Quality Control Algorithms..... 161
 - 6.2.3 Energy-Efficient Hardware Architectures..... 163
 - 6.3 Summary of Results and Comparison..... 166

7 Conclusion and Future Works	169
7.1 Future Works	171
7.1.1 Remaining MVC Challenges	172
7.1.2 3D-Video Pre- and Post-processing	172
7.1.3 Next-Generation 3D-Video Coding.....	172
Appendix A: JMVC Simulation Environment	175
A.1 JMVC Encoder Overview	175
A.2 Modifications to the JMVC Encoder	178
A.2.1 JMVC Encoder Tracing	178
A.2.2 Communication Channels in JMVC	178
A.2.3 Mode Decision Modification in JMVC.....	179
A.2.4 ME/DE Modification in JMVC.....	179
A.2.5 Rate Control Modification in JMVC.....	179
Appendix B: Memory Access Analyzer Tool	181
B.1 Current Macroblock-Based Analysis	182
B.2 Search Window-Based Analysis	182
Appendix C: CES Video Analyzer Tool	185
References.....	189
Index.....	199

List of Figures

Fig. 1.1	Video scaling trend	4
Fig. 1.2	(a) Mobile systems performance trend and (b) Li-ion battery capacity growth.....	6
Fig. 2.1	Macroblocks and slices organization.....	12
Fig. 2.2	Multiview video sequence	13
Fig. 2.3	Multiview video capture, (de)coding, transmission, and display system.....	13
Fig. 2.4	Neighborhood correlation example	15
Fig. 2.5	Prediction comparison between simulcast and MVC.....	17
Fig. 2.6	MVC encoder block diagram	18
Fig. 2.7	MVC prediction structure example	19
Fig. 2.8	Nine prediction directions for intra-prediction 4×4	20
Fig. 2.9	Four prediction directions for intra-prediction 16×16	20
Fig. 2.10	Block processing order in the transform module	21
Fig. 2.11	Zigzag scan order for a 4×4 block.....	21
Fig. 2.12	Order of macroblock borders filtering	22
Fig. 2.13	Temporal and disparity similarities	23
Fig. 2.14	Motion and disparity estimation.....	24
Fig. 2.15	MVC rate control actuation levels.....	28
Fig. 2.16	SRAM voltage-scaling infrastructure.....	35
Fig. 2.17	Energy/power dissipation sources	39
Fig. 2.18	Fast mode decision example.....	41
Fig. 2.19	ME/DE search conceptual example.....	43
Fig. 2.20	Model predictive control (MPC) conceptual behavior	47
Fig. 2.21	Markov decision process (MDP).....	48
Fig. 2.22	Variance-based region of interest map (Flamenco2)	50
Fig. 3.1	MVC energy consumption and battery life	54
Fig. 3.2	MVC component blocks energy breakdown	55
Fig. 3.3	MVC energy breakdown for multiple search window sizes.....	55

Fig. 3.4	MVC energy for distinct mode decision schemes	56
Fig. 3.5	ME/DE energy breakdown	56
Fig. 3.6	MVC vs. Simulcast complexity	57
Fig. 3.7	MVC computation breakdown	58
Fig. 3.8	Memory bandwidth for 4-views MVC encoding	59
Fig. 3.9	Frame-level energy consumption for MVC	61
Fig. 3.10	Memory requirements for motion estimation at MB level	61
Fig. 3.11	Objective video quality in relation to coding modes	64
Fig. 3.12	Energy-efficient Multiview Video Coding overview	66
Fig. 4.1	Coding mode distribution	74
Fig. 4.2	Visual analysis of the coding mode correlation	75
Fig. 4.3	Coding mode hits in the 3D-neighborhood	77
Fig. 4.4	Variance PDF for different coding modes	78
Fig. 4.5	(a) PDF for RDCost difference (between the current and the neighboring MBs) for SKIP <i>hit</i> and <i>miss</i> ; (b, c) Surface plots of RDCost difference for the SKIP coding mode hit and miss; (d) RDCost prediction error for spatial neighbors	79
Fig. 4.6	PDF of RDCost for different prediction modes in Ballroom sequence	80
Fig. 4.7	Average RDCost prediction error for spatial neighbors in Vassar Sequence	81
Fig. 4.8	MVC prediction structure and 3D-neighborhood details	82
Fig. 4.9	MV/DV error distribution between predictors and optimal vector (Ballroom, Vassar)	83
Fig. 4.10	View-level bitrate distribution (Flamenco2, QP=32)	85
Fig. 4.11	Frame-level bitrate distribution for two GGOPs (Flamenco2, QP=32)	86
Fig. 4.12	Basic unit-level bitrate distribution (Flamenco2, QP=32)	86
Fig. 4.13	PDF showing the area of high probability as the shaded region	87
Fig. 4.14	PDF of RDCost for SKIP MBs	88
Fig. 4.15	Threshold curves for RDCost	89
Fig. 4.16	PDF of variance for different prediction modes	89
Fig. 4.17	Overview of the multilevel fast mode decision	91
Fig. 4.18	Early SKIP threshold curves for (a) RDCost and (b) Variance	92
Fig. 4.19	Evaluation of thresholds for early termination (Ballroom, QP=32)	93
Fig. 4.20	Early termination threshold plots for Relax (<i>blue</i>) and Aggressive (<i>red</i>) complexity reduction	94
Fig. 4.21	MVC coding structure for asymmetric coding	96
Fig. 4.22	Energy-aware MVC complexity adaptation scheme	97
Fig. 4.23	Pseudo-code of mode decision for different QCCs	98

Fig. 4.24 Probability density function for RDCost..... 98

Fig. 4.25 Run-time complexity adaptation state machine 99

Fig. 4.26 Average tested modes (QP= {22,27,32,37,42}, GOP=8, Views=8) 102

Fig. 4.27 View-level time savings and Δ PSNR comparison of Relax and Aggressive levels (Exit sequence, QP=32) 102

Fig. 4.28 Average tested modes for all sequences 102

Fig. 4.29 Detailed number of evaluated modes for (a) *Relax* and (b) *Aggressive* (Exit Sequence) 103

Fig. 4.30 Frame-wise PSNR loss comparison of Relax and Aggressive levels (Exit, QP=32)..... 104

Fig. 4.31 Frame-wise time saving comparison of Relax and Aggressive levels (Exit, QP=32)..... 104

Fig. 4.32 Overhead of our scheme 105

Fig. 4.33 Flow diagram of the adaptive fast ME/DE 108

Fig. 4.34 Rate-distortion comparison with full search 110

Fig. 4.35 View-level execution time savings compared to TZ Search..... 110

Fig. 4.36 Comparison of the number of SAD operation..... 111

Fig. 4.37 Hierarchical rate control system diagram..... 112

Fig. 4.38 MPC-based RC horizons 115

Fig. 4.39 Frame-level rate control diagram 115

Fig. 4.40 Basic unit-level rate control diagram..... 119

Fig. 4.41 View-level bitrate distribution (Flamenco2)..... 124

Fig. 4.42 Bitrate and PSNR distribution at frame level (GOP #8)..... 125

Fig. 4.43 Bitrate distribution at BU level (GOP #8) 126

Fig. 5.1 ME/DE hardware architecture template 129

Fig. 5.2 SAD Calculator architecture..... 131

Fig. 5.3 Programmable Search Control Unit (a) FSM and (b) program memory 132

Fig. 5.4 On-chip video memory organization..... 134

Fig. 5.5 On-chip video memory cache tag..... 134

Fig. 5.6 Address Generation Unit (AGU) for dedicated video memory 135

Fig. 5.7 MVC prediction structure in our fast ME/DE algorithm..... 136

Fig. 5.8 Pipeline processing schedule of our ME/DE architecture..... 137

Fig. 5.9 GOP-level pipeline schedule 138

Fig. 5.10 MB-level pipeline schedule for (a) TZ module and fast ME/DE module in (b) Ultra fast and (c) Fast operation modes..... 139

Fig. 5.11 ME/DE search pattern for TZ search and Log search 141

Fig. 5.12 Number of pixels accessed in external memory..... 141

Fig. 5.13 ME/DE search window wastage..... 142

Fig. 5.14 Memory usage variation within one video frame..... 142

Fig. 5.15 Search Map prediction for the Log search 143

Fig. 5.16	Algorithm for Search Map prediction and the dynamic formation of the Search window	144
Fig. 5.17	Analyzing the memory requirements for ME/DE of different MBs in Ballroom sequence	146
Fig. 5.18	Search window memory organization with power gating	147
Fig. 5.19	Search Map prediction accuracy and on-chip memory misses	148
Fig. 5.20	ME/DE hardware architecture block diagram.....	149
Fig. 6.1	Spatial–temporal–disparity indexes for the benchmark multiview video sequences.....	154
Fig. 6.2	Time savings comparison with the state of the art	156
Fig. 6.3	Time savings considering the multiple QPs	157
Fig. 6.4	Time savings distribution summary.....	157
Fig. 6.5	Rate-distortion results for fast mode decision algorithms.....	158
Fig. 6.6	Complexity adaptation for MVC for changing battery levels.....	159
Fig. 6.7	Complexity reduction for the fast ME/DE	160
Fig. 6.8	Average number of SAD operations.....	161
Fig. 6.9	Fast ME/DE RD curves	161
Fig. 6.10	Bitrate prediction accuracy.....	162
Fig. 6.11	Accumulated bitrate along the time.....	163
Fig. 6.12	Rate-distortion results for the HRC.....	164
Fig. 6.13	ME/DE architecture with application-aware power gating physical layout.....	165
Fig. 6.14	Memory-related energy savings employing dynamic search window technique.....	166
Fig. A.1	JMVC encoder high-level diagram.....	176
Fig. A.2	Mode decision hierarchy in JMVC.....	177
Fig. A.3	Inter-frame search in JMVC.....	177
Fig. A.4	Communication in JMVC	178
Fig. B.1	MVC viewer main screen	182
Fig. B.2	Current macroblock-based analysis screenshot.....	183
Fig. B.3	Output example: four prediction directions and their respective accessed areas.....	183
Fig. B.4	Current macroblock-based analysis screenshot.....	184
Fig. B.5	Output example: reference frame access index considering two block matching algorithms: full search and TZ search	184
Fig. C.1	CES video analyzer user interface.....	186
Fig. C.2	CES video analyzer features.....	186
Fig. C.3	Coding mode analysis using CES video analyzer	187
Fig. C.4	ME/DE analysis using CES video analyzer	187

List of Tables

Table 4.1	Predictors hit rate and availability	84
Table 4.2	Quality states	99
Table 4.3	Detailed results for Δ PSNR, Δ Bitrate, and time savings compared to the exhaustive RDO-MD	101
Table 4.4	Comparison between the Quality States (QS)	106
Table 4.5	Comparison of our fast ME/DE algorithm to TZ Search	109
Table 4.6	Variables definitions	114
Table 4.7	Comparison of frame-level HRC Bitrate accuracy	118
Table 4.8	Comparison of BD-PSNR	118
Table 4.9	Comparison of proposed HRC Bitrate accuracy	122
Table 4.10	BD-PSNR and BD-BR comparison	123
Table 5.1	Search Pattern Memory example	133
Table 5.2	Comparison of our fast ME/DE algorithm	149
Table 6.1	Video encoder settings	152
Table 6.2	Simulation infrastructure	153
Table 6.3	Benchmark video sequences	153
Table 6.4	Bjøntegaard PSNR and BR for fast mode decision algorithm	159
Table 6.5	Bjøntegaard PSNR and BR for the HRC	163
Table 6.6	Motion and disparity estimation hardware architectures comparison	164

Abbreviations

3D	Three-Dimensional
3DTV	Three-Dimensional Television
3DV	Three-Dimensional Video (future video standard)
ASIP	Application-Specific Instruction-Set Processor
AVC	Advanced Video Coding
BR	Bitrate
BU	Basic Unit
CABAC	Context-Based Adaptive Binary Arithmetic Coding
CAVLC	Context-Based Adaptive Variable Length Coding
CIF	Common Intermediate Format
CODEC	Coder/Decoder
DC	Direct Current
DCT	Discrete Cosine Transform
DDR	Double Data Rate
DE	Disparity Estimation
DF	Deblocking Filter
DMV	Differential Motion Vector
DPB	Decoded Picture Buffer
DPM	Dynamic Power Management
DSP	Digital Signal Processing
DV	Disparity Vector
DVS	Dynamic Voltage Scaling
EPTZ	Early Prediction Terminator Zone
FIR	Finite Impulse Response
FPGA	Field Programmable Gate Array
FPS	Frames Per Second

FRExt	Fidelity Range Extensions
FSM	Finite State Machine
FTV	Free-Viewpoint Television
GB	Giga Bytes
GDV	Global Disparity Vector
GGOP	Group of <i>Group of Pictures</i>
GIPS	Giga Instructions per Second
GOP	Group of Pictures
HBP	Hierarchical Bi-Prediction
HD1080p	High Definition 1920 × 1080 Progressive
HDTV	High-Definition Digital Television
HEVC	High Efficiency Video Coding
HRC	Hierarchical Rate Control
HVS	Human Visual System
IC	Integrated Circuit
IEC	International Electrotechnical Commission
IEEE	Institute of Electric and Electronics Engineers
IQ	Inverse Quantization
ISO	International Organization for Standardization
IT	Inverse Transform
ITU-T	International Telecommunication Union—Telecommunication
JM	Joint Model for H.264
JMVC	Joint Model for MVC
JVT	Joint Video Team
KIT	Karlsruhe Institute of Technology
MB	Macroblock
MC	Motion Compensation
MD	Mode Decision
MDP	Markov Decision Process
ME	Motion Estimation
MPC	Model Predictive Controller
MPEG	Moving Picture Experts Group
MSE	Mean of Square Errors
MV	Motion Vector
MVC	Multiview Video Coding
MVP	Motion Vector Predictor
PC	Personal Computer
PDF	Probability Density Function
PID	Proportional-Integral-Differential Controller
PMV	Predictive Motion Vector
POC	Picture Order Counter

PSM	Power-State Machine
PSNR	Perceptible Signal-to-Noise Ratio
Q	Quantization
QCC	Quality-Complexity Class
QCIF	Quarter Common Intermediate Format
QHD	Quad HDTV
QP	Quantization Parameter
QS	Quality State
RC	Rate Control
RD	Rate-Distortion
RDO	Rate-Distortion Optimization
RDO-MD	Rate-Distortion Optimized Mode Decision
RGB	Red, Green, Blue
RL	Reinforcement Learning
RoI	Region of Interest
RTL	Register-Transfer Level
SAD	Sum of Absolute Distances
SATD	Sum of Absolute Transformed Distances
SI	Switching I
SIMD	Single Instruction Multiple Data
SoC	System on Chip
SP	Switching P
SRAM	Static Random Access Memory
SSE	Sum of Square Errors
T	Transform
UFRGS	Universidade Federal do Rio Grande do Sul
UVLC	Universal Variable Length Code
VCEG	Video Coding Experts Group
VGA	Video Graphics Array
VHDL	VHSIC Hardware Description Language
VHSIC	Very High Speed Integrated Circuit
VLIW	Very Large Instruction Word
VP	Viewpoint
YUV	Luminance, Chrominance Component 1, Chrominance Component 2

Chapter 1

Introduction

The consumers' thirst for new and more immersive multimedia technologies allied to the industry interest to boost the entertainment market has driven the fast popularization of 3D-video content generation, 3D-capable devices, and 3D applications. Although the first 3D-video device was developed in 1833 and the first 3D film demonstration dates from 1915, this format only became worldwide known in the 1980s through IMAX technology. The real 3D-video hype, however, was noticed in the late 2000s through the massive popularization and availability of 3D movies followed by the 3D-capable televisions dedicated to home cinema. For a better perspective of this popularization, more than 10 % of the televisions sold in USA in 2011 were 3D capable. The latest field to be affected by the 3D-video popularization is exactly the field responsible for the biggest IC (integrated circuits) industry growth after the popularization of personal computers: the mobile embedded systems. Smartphones, tablets, personal camcorders, and other mobile devices shipments already surpassed PC shipments. For instance, more than 650 million smartphones are expected to be shipped in 2013 compared to 430 million PCs in the same year. Jointly, the popularization of 3D videos and mobile devices is leading to a scenario where a large amount of such 3D-capable smart devices is reaching the users every day, resulting in a large amount of 3D-video content being generated, encoded, stored, transmitted, and displayed. According to CISCO, video content already represents 51 % of the current Internet traffic and is envisaged to touch the 90 % mark due 2014. It is also important to consider that the 0.6 Exabytes per month mobile traffic in 2011 is expected to reach 10.8 Exabytes per month in 2016.

To cover the gap between 3D-video content generation and network and storage capabilities there is a need to efficiently encode 3D videos and reduce the amount of data required for their representation. The multiview video coding (MVC), an extension to the H.264/AVC, is the state of the art on 3D-video coding. Based on the multiple views paradigm, as the majority of current 3D-video technology, the MVC reduces the 3D videos representation in 20–50 % compared to H.264/AVC simulcast. The cost of this efficiency improvement comes from an increased coding complexity and increased energy consumption, mainly at the encoder side.

The energy consumption incurs from multiple processing units working in parallel to attend throughput constraints (processors, DSPs, GPUs, ASICs) and intense memory access. In a scenario dominated by mobile devices, the increase in energy consumption goes against the battery restrictions posed by these mobile embedded systems. This conflict of interests between coding efficiency and energy constraints brings the main challenge related to 3D-video realization on embedded systems: jointly *design algorithmic and architectural energy-efficient solutions to enable real-time high-definition 3D-video coding, while maintaining high video quality under severe energy constraints*. The main goal of this monograph is to address this challenge by presenting novel algorithms and hardware architectures designed to show the feasibility of 3D-video encoding on embedded battery-powered devices.

In the next sections, after this introduction, an overview of 3D-video applications that make the 3D-video field so promising is presented. After that, a brief introduction on the trends for 3D-video coding and multimedia embedded systems is presented, followed by the related issues and research challenges. This chapter is finalized by a summary containing the contributions of this work.

1.1 3D-Video Applications

The adoption of 3D videos is directly associated with the existence of new applications requiring the deepness sensation in order to improve the users' immersion experience. From here onwards an overview of the main 3D-video applications is presented. These applications share the same concept of capturing multiple views in the same 3D scene. To give the depth illusion, distinct views are displayed to each eye with displays that employ technologies based on parallax barriers, lenticular sheets, color polarization, directional polarization, or time interleaving; more details on this phenomenon are provided in Chap. 2.

- *Three-dimensional video personal recording*: Popularized by the 3D-capable mobile devices and the 3D-video sharing services the 3D-video personal recording is the most massive 3D-video service in terms of video content availability. With a 3D-video recorder device the users are free to create and publish their own video content.
- *Three-dimensional television (3DTV)*: 3DTV is an extension of the traditional 2D with the depth perception. In this kind of application two or more views are decoded and displayed simultaneously where each viewer sees two views, one for the right eye, and other for the left eye. The simplest 3D displays, which are the stereoscopic displays, show two simultaneous views requiring the use of special glasses (polarized or active shutter glasses) to provide 3D sensation. The evolution of stereoscopic displays is the auto-stereoscopic display, which eliminates the need for glasses. In this case, parallax barriers and lenticular sheets are the most common solutions. Multiview displays are able to display higher number of views at the same time increasing the observer freedom by supporting head parallax, i.e., the viewpoint changes when the observer changes its position.

- *Free-viewpoint television (FTV)*: In this application, the user is able to select the desired viewpoint in a 3D scene. It provides realism and interactivity to the user, i.e., the focus of attention can be controlled. The display technology used may vary from 2D televisions to multiview displays.
- *Three-dimensional telepresence*: Allows the user to communicate and interact with interlocutors as if they were in the same location. Telepresence has been widely used for video teleconferencing, mainly in corporate environments, and for the implementation of the so-called virtual offices. The evolution towards 3D represents a meaningful step in order to improve the perception and interaction level between the conference attendees.
- *Three-dimensional telemedicine*: Telemedicine was defined to surpass physical limitations and make it possible for a doctor to attend patients or perform surgeries while in a distinct location by using telecommunications methods. The 3D-video capability brings the telemedicine to a whole new level where the specialist can precisely perceive the 3D space and proceed accurately through robotic actuators. This technology enables a better health care quality in remote places that do not count on qualified specialists.
- *Three-dimensional surveillance*: Traditional video surveillance systems rely on 2D videos and pose difficulties to authorities if precise depth information is required. Employing 3D videos for surveillance provides a much richer information once it is possible to accurately extract depth and angulations data for all objects in the 3D scene. Therefore, a better description on the interaction between objects, such as possible criminals and victims, is obtained.

Among these applications, some are not designed for mobile use (e.g., 3D Surveillance and 3D Telemedicine) or require only decoding at the mobile device (e.g., 3DTV, FTV). For other applications, however, the capability to encode 3D videos is mandatory. For instance, 3D-video personal recording requires real-time and energy-efficient 3D-video encoding. 3D Telepresence, when running on embedded devices, demands real-time, energy-efficient, and low-delay 3D-video encoding. Aware of the challenges posed by the presented set of applications, this work focuses on the MVC video encoder.

1.2 Requirements and Trends of 3D Multimedia

Although the processing power of computational systems, mainly for embedded systems, has increased meaningfully (as detailed in Sect. 1.3), the multimedia applications' performance and energy requirements are increasing in a significantly higher pace due to increased video resolutions, frame rates, sampling accuracy, and number of views in case of 3D videos. In other words, the amount of data to be processed in a video sequence has been increasing in multiple axes simultaneously.

Figure 1.1 relates the number of macroblocks (MB— 16×16 image block used as basic coding unit in MVC—for details refer to Chap. 2) to be processed per second considering the different video resolutions, frame rates, and number of views.

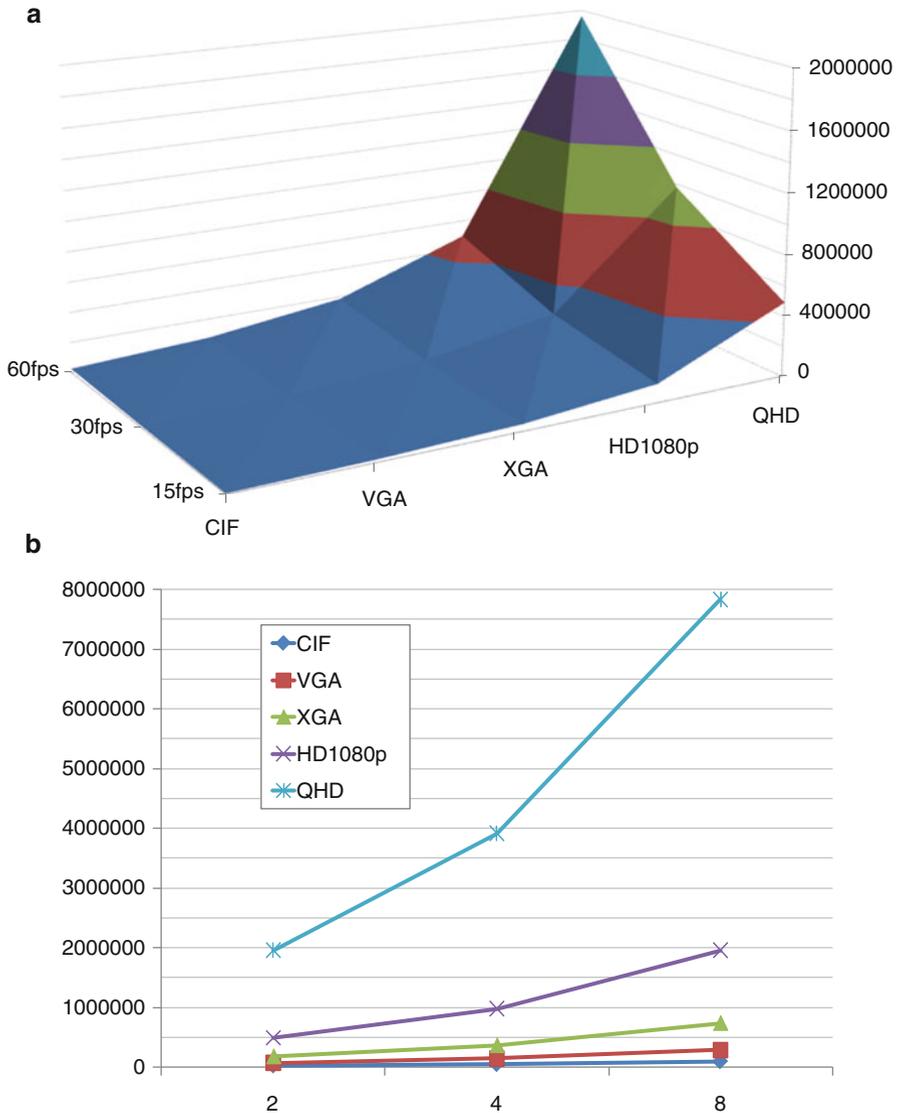


Fig. 1.1 Video scaling trend

Previous coding standards, for instance MPEG-2, were designed and typically used in videos with low-medium resolutions and low-medium frame rates such as CIF (352×288), VGA (640×480), and SDTV (768×576) at 15–30 fps (frames per second) (note that these numbers refer to the typical use and main target operation profiles; the standards define a very high operation range). The H.264 additionally targets high resolutions and high frame rates such as 720p (1240×720) and

HD1080p (1920×1080) at 30–60 fps. The next generation of coding standards, represented by H.265/HEVC (High Efficiency Video Coding), will also target on high and ultrahigh resolutions and frame rates including QHD (3840×2160) and UHD TV (7680×4320) videos at 60–120 fps. To quantify this growth, the relation between the corner cases shown in Fig. 1.1a, CIF@15 fps and QHD@60 fps, is equivalent to a 327× factor. Also, targeting improved quality, the samples' bit depth is increasing from 8 bits up to 14-bit samples, requiring wider data operators. At the complexity and energy consumption perspective, the scenario is even worse once there is a nonlinear relation with the data amount. The increase in resolution, for instance, leads to higher processing effort per MB, higher memory traffic, and larger on-chip memory related to the Motion Estimation (ME; see Chap. 2), resulting in energy consumption increase. Moreover, the video coding standards evolution severely contributes to the increase of complexity and energy requirements. For example, the H.264 encoder is approximately 10× more complex than the MPEG-4 encoder, while the HEVC is expected to bring additional 2–10× complexity increase factor in relation to H.264.

Considering 3D videos, the scaling scenario becomes more dramatic, as shown in Fig. 1.1b. Besides the resolution and frame-rate increase, it is necessary to deal with the linear data growth in relation to the number of views. As MVC includes new coding tools the complexity and energy consumption increase in a nonlinear (above-linear) fashion, as quantified in Sect. 3.1. The impacts of the fast 3D multimedia requirements scaling on embedded systems are discussed in the next section.

1.3 Overview on Multimedia Embedded Systems

The fast evolution of multimedia embedded systems has been driven by the so-called smart devices (smartphones, tablets, and other mobile devices capable of data, audio, and video communication) popularization. Meaningful progress has been done by the major players in the field, in terms of performance boost and energy efficiency. The progress, however, is not enough to fill the gap between multimedia application requirements and technology evolution. The ARM specialists, whose processors equip about 90 % of the current embedded devices, predict a performance increase in the order of 10× when comparing the state of the art in 2009 to the predicted one for 2016, as shown in Fig. 1.2a. Energy restrictions related to slow battery evolution is the major factor limiting the performance of embedded systems. According to Panasonic, the capacity of Li-ion batteries has been increasing, on average, 11 % annually since 1994, as shown in Fig. 1.2b.

The high performance and energy efficiency required by the current 3D-video applications are not met by generic embedded solutions such as embedded processors, GPUs, and DSPs. There is a need to implement application-specific hardware accelerators to deliver the required throughput while minimizing energy consumption at the cost of a flexibility drawback. The latest high-end embedded SoCs (System on Chip) already implement this approach for multimedia processing,

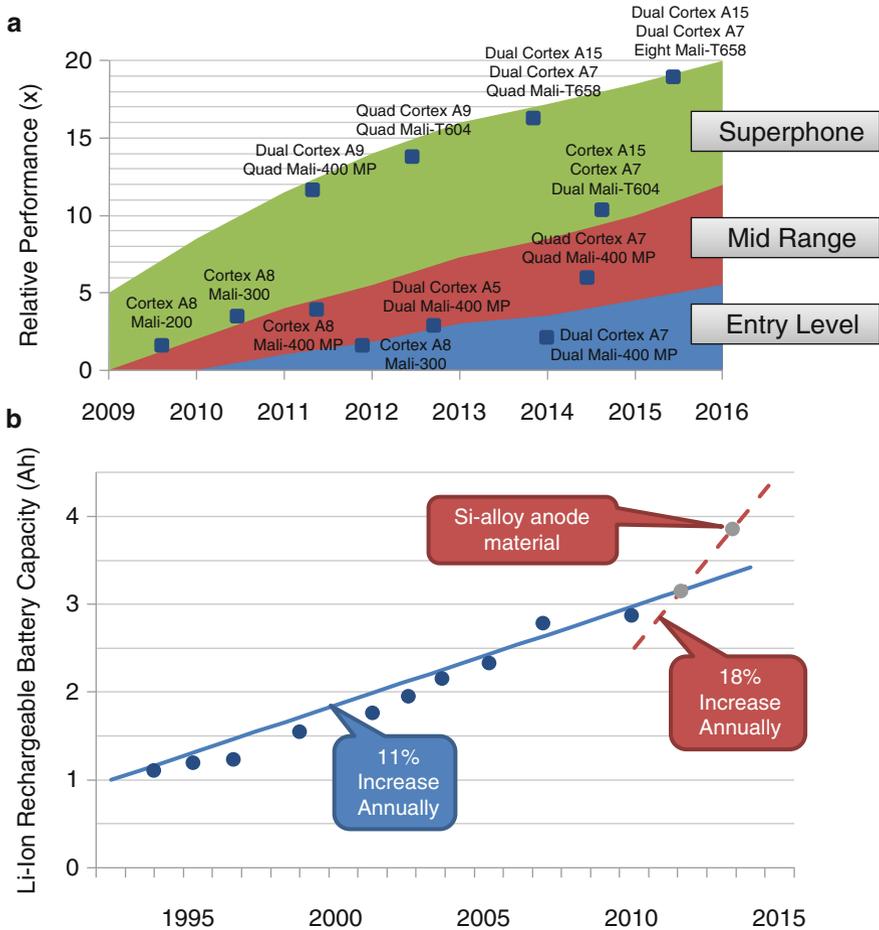


Fig. 1.2 (a) Mobile systems performance trend and (b) Li-ion battery capacity growth

e.g., H.264 video encoding and decoding, as detailed in Sect. 2.5. Some examples are Qualcomm Snapdragon, Nvidia Tegra, Samsung Exynos, and Texas Instruments OMAP. The hardware support, however, needs to be extended in order to efficiently handle 3D videos.

1.4 Issues and Challenges

The demand for mobile 3D multimedia processing allied to high-performance demands and severe embedded devices energy constraints poses serious challenges to the researchers and developers actuating in the embedded multimedia systems

field. In this scenario, employing hardware accelerators optimized for specific MVC application is mandatory. Given the gap between 3D multimedia processing and the embedded processing reality, there is a need to further reduce the complexity and energy consumption at algorithmic and architectural levels. Such optimizations are only possible by employing deep application knowledge to perform a coupled and integrated optimization of the algorithms employed and the underlying hardware architecture.

In addition to the varying coding settings and battery state, multimedia applications are susceptible to input content variations that significantly change the system behavior and requirements. For instance, videos with higher motion intensity require more processing and memory accesses resulting in more processing units and larger on-chip memory finally leading to increased energy consumption. Such variations are only detected at run time. Therefore, energy-efficient MVC encoding systems require algorithmic and hardware run-time adaptivity that employ application and video content characteristics knowledge. The adaptation schemes must be able to handle the energy efficiency vs. video quality trade-off in order to find the optimal operation point for each given system state and video input.

Energy reduction algorithms and energy-oriented optimizations might lead to rate-distortion (RD) performance losses, i.e., video quality reduction for the same bitrate. To avoid or minimize this drawback, there are mechanisms able to control the losses through the optimization of the bit distribution among different views, frames, and image regions.

The in-depth study of the issues and challenges related to MVC encoding is presented in Chap. 3. In the following section, the contribution is summarized.

1.5 Monograph Contribution

The goal of this monograph is to understand the run-time behavior of the MVC encoder at the energy consumption perspective and propose algorithms and hardware architectures able to jointly attend the performance constraints and respect the energy envelope restrictions for state-of-the-art embedded devices. In this section, a summary of the contributions of this monograph is presented, highlighting the main innovations proposed. A deeper description of these contributions is found in Chap. 3, while the technical details are presented in Chaps. 4 and 5, and results in Chap. 6.

1.5.1 3D-Neighborhood Correlation Analysis

The novel energy-efficient algorithms and hardware architectures proposed in this work are designed upon a strong MVC application knowledge including all MVC encoder algorithms and their run-time response to distinct input data. Along this work, the application knowledge, for many cases, is studied in terms of the

correlation within the 3D-neighborhood. The 3D-neighborhood concept is a space domain defined in this monograph that contains the MBs belonging to the neighboring regions in the spatial, temporal, and disparity axes. Due to the redundancies existing within this neighborhood (see discussion in Sect. 2.2), the 3D-neighborhood provides valuable information to predict video encoding side information, algorithms behavior, memory access pattern, etc. Therefore, the offline and online 3D-neighborhood data are used to define and control the energy-efficient algorithms, hardware design, memory architecture and sizing, and adaptation schemes.

1.5.2 *Energy-Efficient MVC Algorithms*

The energy-efficient algorithms for MVC are concentrated in three MVC encoding blocks: mode decision, motion and disparity estimation (ME/DE), and rate control (RC). Mode decision (MD) and ME/DE units are responsible for the dominant energy consumption in the MVC encoder, as discussed along Chap. 3. The proposed fast MD and ME/DE target energy reduction through complexity reduction. These algorithms interact with the novel energy-aware complexity adaptation algorithm that controls the energy consumption by changing the coding efficiency considering battery state. The drawback posed by the energy-efficient algorithms comes in terms of quality drop under certain coding conditions. To minimize this negative impact a hierarchical rate control (HRC) solution to optimize the bit utilization while maximizing and smoothen video quality in spatial, temporal, and disparity domains is proposed.

- *Multilevel mode decision-based complexity adaptation*: Incorporates an Early SKIP prediction technique to a sophisticated mode decision scheme composed of six decision steps and bad-prediction protection. This fast MD employs multiple MD aggressiveness strengths (to control energy vs. quality losses), 3D-neighborhood knowledge, coding modes ranking, video properties-based prediction, and Rate-Distortion cost (RDCost) prediction. Quantization parameter (QP)-based thresholding is employed to react to QP changing scenarios. The complexity adaptation algorithm employs asymmetric view coding to maximize the perceived video quality in face of battery discharging and provides graceful quality degradation along the time.
- *Fast motion and disparity estimation*: The proposed Fast ME/DE widely exploits the motion and disparity vectors correlation within the 3D-neighborhood in order to avoid the search for non/key frames in the MVC prediction structure. According to the confidence in the neighboring MBs, the algorithm selects the Fast or Ultra-Fast prediction mode.
- *Hierarchical rate control*: This innovative solution for the MVC rate controller employs two actuation levels, frame-level and basic unit-level rate control, with coupled feedback loop. The frame-level RC uses the Model Predictive Controller (MPC) to estimate the bitrate for future frames and decide the best QP. Markov decision process (MDP) with reinforcement learning (RL) and regions of interest (RoI) weighting is employed at BU level to further optimize the QP selection within the frames.

1.5.3 Energy-Efficient Hardware Architectures

The energy-efficient hardware architectures target the motion and disparity estimation processing, which represents the most complex and energy-intense coding block of the MVC encoder. A ME/DE architecture is proposed aiming to reduce the energy consumption for 4-views real-time HD1080p encoding through on-chip memory and external memory accesses reduction and efficient dynamic power management for the processing path and memory architecture. The architectural innovations are introduced in the following and detailed in Chap. 5.

- *Motion and disparity estimation hardware architecture*: Along this monograph is proposed an architectural solution for the ME/DE block in the MVC encoder. This architecture features techniques to improve the performance and reduce the overall energy consumption. Our description defines each main building block composing the proposed architecture and the interaction between them. The hardware blocks are designed to provide support to multiple search algorithms, throughputs, and memory hierarchy.
- *Multibank on-chip video memory*: This proposal enables a reduced on-chip video memory and sector-level power gating in order to reduce the energy consumption through leakage current lowering. The on-chip memory works in a cache fashion and employs multiple banks for high throughput. Distinct dynamic power management (DPM) techniques are proposed based on the memory prediction using the 3D-neighborhood information.
- *Memory design methodology*: A study of the memory requirements under different coding scenarios and video contents is presented to provide the basis for defining the memory size and organization. Based on this study, an offline statistical analysis is used to define the memory hierarchy considering on-chip memory size and number of external memory access.
- *Dynamic search window formation-based data reuse*: Macroblocks previously encoded in the 3D-neighborhood are used to create a search map that tracks the search pattern behavior. From the search map, a prefetch scheme named Dynamic Search Window formation is employed. This technique focuses on the reduction of external memory accesses and the reduction of active memory sectors in the on-chip video memory.
- *Application-aware power gating*: This proposal implements a memory requirements prediction scheme to accurately control power states of the on-chip video memory sectors. Once again, the MBs within the 3D-neighborhood are used as source of information for decision making.

1.6 Monograph Outline

This monograph is organized as follows:

Chapter 2 presents an overview of the background knowledge required to understand this work along with the related works published in academic channels

and state-of-the-art industrial solutions. The basics of 2D and 3D digital video concepts, 3D-video systems, multimedia architectural options, and MVC are provided. Afterwards, a state-of-the-art revision is presented including the latest reduced-complexity and energy-efficient solutions for the MVC encoding.

Chapter 3 brings a deep study and discussions on the requirements and challenges related to the realization of MVC real-time encoding on embedded devices. The discussions are centered on the energy consumption and encoded video quality. Chapter 3 also presents the overview of contributions presented along this monograph. For simplicity, the monograph contribution is also summarized using a high-level diagram.

In Chap. 4 all the novel energy-efficient algorithms proposed in this work are thoroughly explained. They are classified and described in three sections: coding mode decision with complexity adaptation, motion and disparity estimation, and video quality management. Technical details ranging from case studies down to implementation level are followed by algorithm-specific results.

The architectural contribution for motion and disparity estimation is presented in Chap. 5. Firstly, an architectural overview is provided followed by a detailed description of each hardware component block. At this point onwards, Chap. 5 brings the discussion on the processing scheduling, the novel dynamic search window formation algorithm, and the on-chip memory design and application-aware power gating techniques. Additionally, the architectural specific results are presented in this chapter.

Chapter 6 brings the overall results for the proposed novel algorithms and architectures compared to state-of-the-art related works. Chapter 7 depicts the conclusions of this work and points to future research opportunities and challenges related to the next generations of 3D multimedia processing and 3D-video coding.

Additional tools and simulation environments are presented in the appendices. Appendix A presents the MVC reference software, the JMVC, and details the modifications applied to the JMVC in order to enable software experimentation. The in-house developed Memory Access Analyzer tool and its graphic interface is presented in Appendix B. Appendix C presents the CES Video Analyzer tool highlighting the extensions implemented to support multiview videos.

Chapter 2

Background and Related Works

In this chapter the basic notions on digital videos, multiview video systems, and the multiview video coding (MVC) standard are presented. The mode decision, motion and disparity estimation, and rate control modules are detailed since they are the main foci of this monograph. Detailed state-of-the-art review is presented considering 3D-video systems, multimedia architectures, energy-efficient algorithms, and architectures for video coding.

2.1 2D/3D Digital Videos

A video is formed by a sequence of frames (or pictures) of a scene captured in a given frame rate providing to the spectator the sense of motion. Usually, the frame rate goes from 15 to 60 frames per second (fps) depending on the application requirements. Each frame is formed by a number of points named picture elements, i.e., pixels. The number of pixels in each frame is called resolution, i.e., the number of horizontal and vertical pixel lines. The typical resolutions also depend on the target application. For instance, mobile devices use to handle relatively lower resolution and lower frame rate sequences if compared to home cinema that targets high resolution and high frame rates.

Different color spaces are used to represent raw and decoded videos; the most usual ones are RGB (red, green, blue) and YUV. Most monitors operate at the RGB space, while most of video coding standards work over the YUV space. The YUV space is composed of three color channels: one luminance (Y) and two chrominance channels (U and V). The main reason for using YUV space for video coding is related to its smaller correlation between color channels, making easier to independently encode each channel. Since the human visual system (HVS) is less sensible to chrominance when compared to luminance, it is possible to reduce the amount of chroma information without affecting the overall perception. The reduction of chroma information is made using color subsampling (also known as pixel

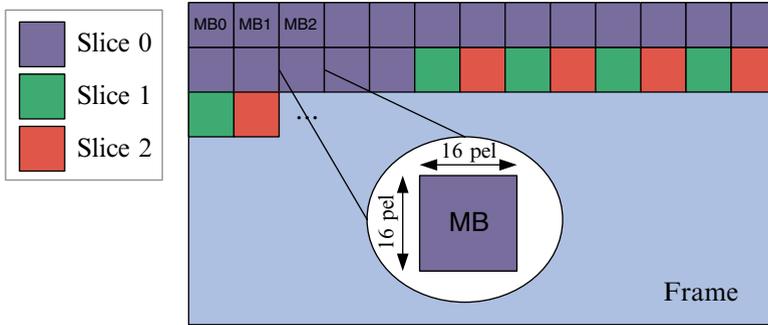


Fig. 2.1 Macroblocks and slices organization

decimation). The most used color subsampling pattern is the YUV 4:2:0 that stores one U and one V sample for each four luminance samples reducing in 50 % the total amount of raw video data.

All current widely used video coding standards are based on block coding. In other words, they divide each frame in pixel blocks to encode the video. These blocks are named macroblocks (MB). In the H.264, the latest video coding standard, the MBs are blocks of 16×16 luma pixels and its associated chroma samples (see Fig. 2.1). A group of MBs is called slice. The slice can be formed by one or more MBs that may be contiguous or not. One frame is formed by one or more slices. In turn, each slice is classified in one of three different types (here the SI and SP slices are not considered): Intra (I), predictive (P), and bi-predictive (B) slices. The example in Fig. 2.1 is composed of three slices: one contiguous (Slice 0) and two noncontiguous slices (Slices 1 and 2). Note, the terminology used here is based on the H.264 standard and is directly applicable to the MVC standard.

For a better comprehension on the different slice types it is necessary to understand the two basic prediction modes used by the state-of-the-art video encoders: intra-frame and inter-frame prediction. The intra-frame prediction only exploits the spatial redundancy by using surrounding pixels to predict the current MB. The inter-frame prediction exploits the temporal redundancy (similarity between different frames) by using areas from other frames, called reference frames, in order to better predict the current MB. Intra (I) macroblocks use the intra-frame prediction, while predictive (P) and bi-predictive (B) macroblocks use the inter-frame prediction. While P macroblocks only use past frames as reference (in coding order), the B macroblocks can use reference frames from past, future, or a combination of both. Intra slices are formed only by I MBs. Predictive (P) slices support I and P macroblocks and bi-predictive (B) slices support I and B macroblocks.

Multiview video sequences are composed of a finite number of single-view video sequences captured from independent cameras in the same 3D scene. Usually these cameras are carefully calibrated, synchronized, and positioned. They are typically aligned in a parallel 1D-array or 2D-array; however, there are systems where the cameras are positioned in arch or cross shapes. The typical spacing between

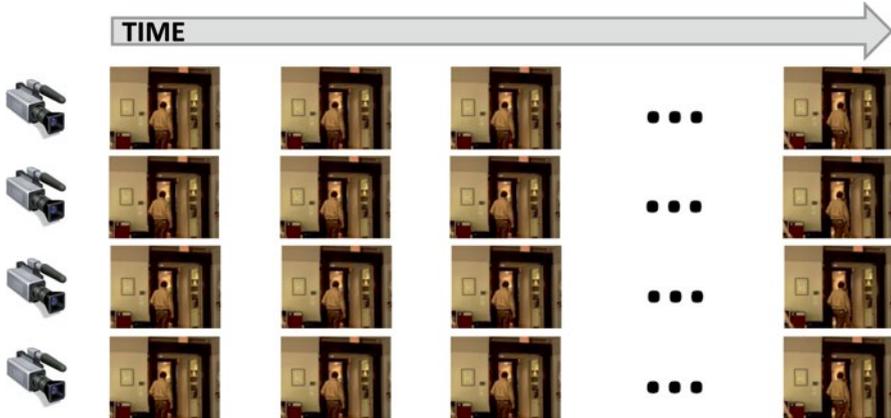


Fig. 2.2 Multiview video sequence

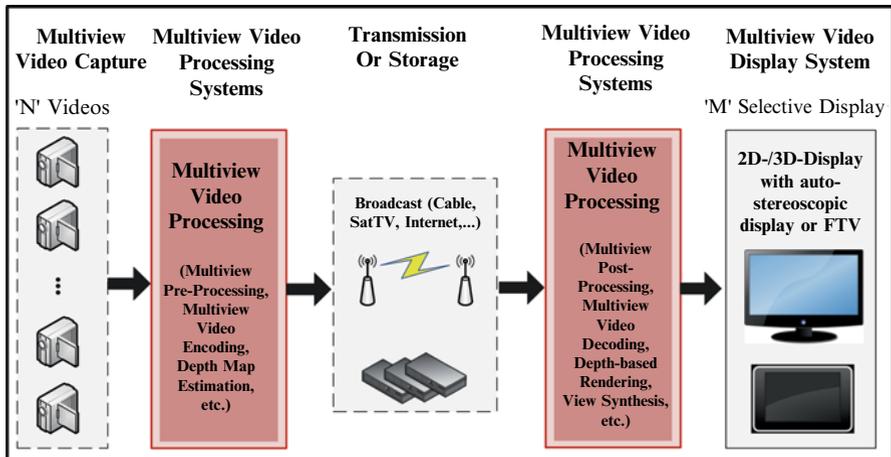


Fig. 2.3 Multiview video capture, (de)coding, transmission, and display system

cameras is 5 cm, 10 cm, or 20 cm for most of the available test sequences. In Fig. 2.2 a multiview video with four views and the captured frames along the time axis is presented. At the video encoding perspective, the MVC, as detailed in Sect. 2.3, extends the concept of inter-frame prediction to inter-view prediction where the correlation between different views is exploited. A deeper discussion regarding the spatial, temporal, and view/disparity correlations is provided in Sect. 2.2.

Figure 2.3 depicts the complete system required to capture, encode, transmit, decode, and display multiview videos. The captured sequence is encoded by an MVC encoder in order to reduce the amount of data to be transmitted. The generated bitstream may be transmitted using broadcast or Internet or stored in media servers

or local storage. At the decoder side the bitstream, or part of it, is decoded and displayed according to the displaying technology available at the receiver end. In a simple single-view display the decoder considers only the base view that is decodable with a regular (H.264/AVC) video decoder. In the case of stereoscopic displays (two views) only two views are decoded and displayed. In free viewpoint television (FTV) systems the user selects the desired viewpoint within the 3D scene and the video decoder selects which views to decode. For multiview displays all views displayed must be decoded along with the reference views used to reconstruct them.

2.2 Multiview Correlation Domains

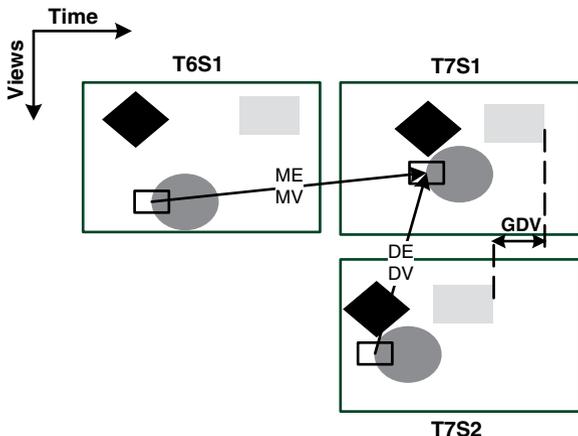
This section defines the three types of redundancies or correlations present in multiview video sequences in order to provide the background required for a better understanding of the MVC coding tools, detailed in Sect. 2.3, and for the 3D-neighborhood concept presented in Sect. 3.5.1. Here we discuss the correlation at pixel level, i.e., the similarities used to predict the image pixels, and at coding information level, i.e., how neighboring blocks share coding properties such as coding modes, vectors, etc. To have a more general description we present independently the three correlation dimensions (1) spatial correlation, (2) temporal correlation, and (3) view/disparity correlation. Single-view video coding standards are able to exploit (1) and (2), while MVC incorporates (3) to provide improved prediction for multiview videos.

2.2.1 Spatial Domain Correlation

The spatial correlation is the similarity within regions in the same frame. Previous image and video coding standards, such as JPEG2000 and H.263, were already able to exploit this similarity through MB prediction based on neighboring pixels (see Sect. 2.3). Neighboring MBs tend to belong to the same image region and share similar image properties. For this reason, the surrounding pixels typically are good block predictors for the intra-frame prediction process. Exception cases happen in object borders where the image properties may change abruptly. Consider the example in Fig. 2.4, all the MBs in the white background share similar image properties. The same happens for the MBs within one of the objects. The discontinuity occurs when an object border is found leading to increased prediction error. Note that, for simplicity, the spatial correlation is referred as one dimension, but it is actually composed of two dimensions, the width and height of a frame.

On average, the current coding standards are able to efficiently employ the intra-frame prediction for pixel data. However, the correlation of coding side information (coding mode, motion vectors, disparity vectors, etc.) is just superficially exploited. In H.264, a few simple techniques exploit this kind of correlation. The differential

Fig. 2.4 Neighborhood correlation example



coding of intra prediction modes inside a macroblock exploits spatial correlation of coding information. In this technique, the intra coding mode is coded considering the coding mode of the previous block. Another example is the motion vector prediction process that uses the neighboring vectors to predict the current vector. By employing the motion vector prediction, only the differential motion vector needs to be coded and transmitted. These examples show that there is also significant correlation at coding side information level.

2.2.2 Temporal Domain Correlation

The temporal correlation represents the similarities between different frames in the same view of a video sequence. That is, the objects of a given frame are usually present in neighboring temporal frames with a displacement that depends on its motion. Consider the frames T6S1 (view 1, time 6) and T7S1 (view 1, time 7) in Fig. 2.4, the same objects are seen in both frames with a small displacement. Thus, frame T7S1 may be accurately predicted from the reference frame T6S1. The displacement between the two frames is found using the motion estimation (see Sect. 2.3.2). Besides the pixel-level prediction, the coding data are also similar for the same object along the time. In other words, for the same object in distinct time instants the same set of coding modes and motion behavior tend to be employed. The correlation is lost when there is an occlusion or the object moves out of the captured scene.

Analogous to the spatial correlation, there are tools able to exploit the temporal correlation at pixel level, i.e., the motion estimation (Sect. 2.3.2). At coding side information level, an attempt to exploit this correlation was proposed in the H.264 standard by using the temporal direct prediction for motion vectors. This prediction uses the collocated MB (MB sharing the same relative position in the frame) motion vector in order to predict the current one.

2.2.3 Disparity Domain Correlation

The disparity is a complete new domain introduced by multiview videos. It refers to the similarities between frames in different views. The similarities or redundancies at pixel level are exploited by the disparity estimation tool (Sect. 2.3.2). However, no tool is able to exploit this correlation at the coding information level.

As depicted in Fig. 2.4, frames T7S1 (view 1, time 7) and T7S2 (view 2, time 7), the same objects are present in the neighboring views displaced by the so-called disparity vector. Since they are the same objects, the same image properties are shared and similar coding information tends to be used in different views. The disparity neighborhood correlation is lost when a given object is out of the area captured by a given camera or there is an object occlusion for a given camera point of view.

In order to obtain an accurate evaluation of the available correlation, we have carried out an extensive analysis of multiview videos. For this analysis we have used different multiview video sequences following the MVC test recommendation by Joint Video Team (JVT) (Su et al. 2006). These sequences have coding structures similar to the one presented in Fig. 2.7. Our analysis, discussed in Sect. 3.5.1, constitutes an in-depth exploration of coding mode distribution, video statistics, motion and disparity vectors, coding mode, and RDCost correlation in the so-called 3D-neighborhood (spatial, temporal, and view neighborhood).

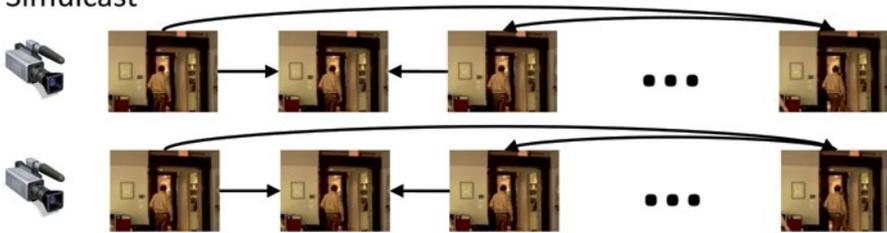
2.3 Multiview Video Coding

Encoding multiview video sequences can be performed using different techniques. The most primitive one is the simulcast approach, where a single-view video coding standard (usually H.264/AVC) is used to encode independently each view. As presented in Fig. 2.5, the simulcast approach considers the intra-frame prediction and inter-frame prediction (a.k.a. motion estimation) exploiting the spatial and temporal redundancies. However, the disparity or inter-view redundancy (i.e., the redundancy between frames of different views) is not considered. The MVC standard uses the inter-view prediction (a.k.a. disparity estimation) to take advantage of the similarities between views from the same scene. The inter-view prediction represented by the red arrows in Fig. 2.5 is responsible for a bitstream reduction of 20–50 % for the same video quality (Merkle et al. 2007). Details on the MVC new tools, coding efficiency, and complexity are discussed along this section.

In a strict definition, the MVC is not a coding standard but an extension of the H.264/AVC or MPEG-4 Part 10 (JVT 2003). The MVC was defined by the JVT in March 2009 (JVT 2008 and JVT 2009b). The JVT is the group of experts formed by the Motion Picture Experts Group (MPEG) from ISO/IEC and the Video Coding Experts Group (VCEG) from ITU-T.

The standard usually works over the YUV (or YCbCr) (Miano 1999) color space that is composed by one luminance channel and two chrominance channels (red and blue chrominance), but other color spaces are supported, such as RGB and YCoCo

Simulcast



MVC

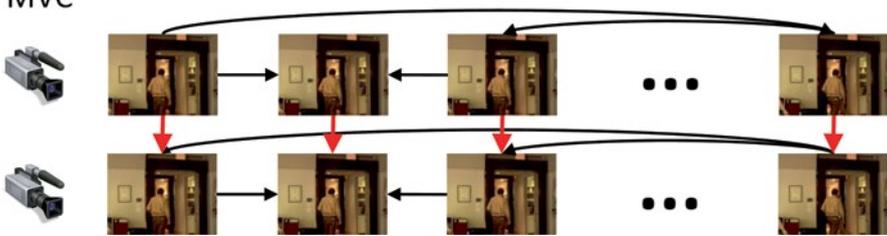


Fig. 2.5 Prediction comparison between simulcast and MVC

(orange and green chrominance). The MVC also supports different subsampling patterns including 4:2:0 (four luminance samples for one sample of each chrominance channel), 4:2:2 (two luminance samples for one sample per chrominance channel), and 4:4:4 (one luminance channel for one sample in each chrominance channel). The supported color space/subsampling and coding tools depend on the profile of video coding operation (JVT 2009a, b).

Originally three profiles were defined in the H.264 standard: Baseline, Main, and Extended. The Baseline profile focuses on video calls and videoconferencing. It supports only I and P slice and the context-adaptive variable length coding (CAVLC) entropy coding method. The Main profile was designed for high-definition displaying and video broadcasting. Besides the tools defined by the Baseline profile, it also includes the support to B slices, interlaced videos, and CABAC entropy coding. The Extended profile targets video streaming on channels with high package loss and defines the SI (Switching I) and SP (Switching P) slices (Richardson 2010). In 2005 the Fidelity Range Extension (FRExt) defined the High profiles: High, High 10, High 4:2:2 and High 4:4:4 targeting high fidelity videos (JVT 2009a, b).

The MVC extension introduced to the standard a new set of CABAC contexts and new supplemental enhancement information (SEI) messages to simplify parallel decoding and the transmission of sequence parameters (JVT 2009a, b). Additionally, the disparity estimation or inter-view prediction was proposed (Merkle et al. 2007). This is the most important innovation in the MVC that allows the exploration of similarities between different views. Its function is to find the best matching for the current macroblock in a reference frame within the reference view. The possible search criteria, search patterns, and objective are similar to the motion

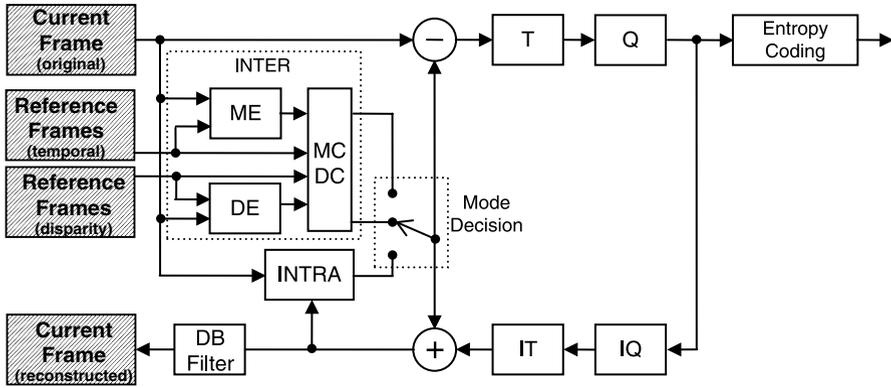


Fig. 2.6 MVC encoder block diagram

estimation. However, the dynamic behavior of the disparity estimation differs significantly with respect to the ME. In the following section, details of the MVC encoding process are presented.

2.3.1 MVC Encoding Process

In Fig. 2.6 the high-level block diagram of the MVC encoding process is presented. As a hybrid coding standard it is composed of three phases: prediction, transforms, and entropy coding. The transform and entropy phases are similar to H.264/AVC, except for the new syntax elements to be encoded by the entropy encoder. The main innovation is in the prediction phase, which incorporates the inter-view prediction tool, the disparity estimation (DE).

The base view, the first one to be encoded, is encoded in compliance to the H.264 standard. Then, the prediction has two options, the intra-frame or the inter-frame prediction. Other views are named dependent views and additionally employ inter-view prediction. The complete encoding process is described in this section, considering the Main profile tools in YUV color space with 4:2:0 subsampling, while further extensions available in the High profiles are omitted for simplicity.

The MVC prediction structure inherits all the possibilities for temporal references and coding orders defined by the H.264. In addition, distinct possibilities of view coding order may be employed. The most used view coding orders are IPP and IBP (Merkle et al. 2007). The prediction structure depicted in Fig. 2.7 employs IBP view coding order using hierarchical bi-prediction (HBP) structure in temporal domain for eight views and group of pictures (GOP) size equals to 8. The set of GOPs for all views are referred in MVC as GGOP (group of groups of pictures). The frames located in the GGOP borders are called anchor frames while all others are the non-anchor frames.

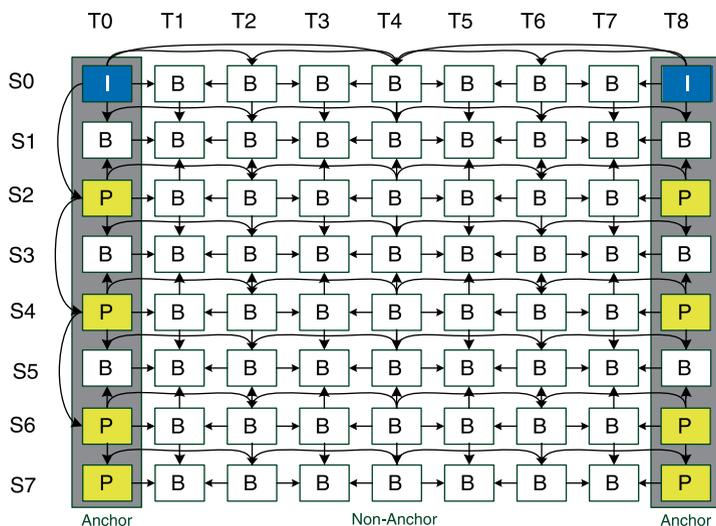


Fig. 2.7 MVC prediction structure example

The intra-frame prediction uses the neighboring pixels within the frame to predict the samples in the current MB. The MVC supports two MBs partitioning sizes for intra-frame prediction. The size 4×4 has nine prediction directions, as presented in Fig. 2.8, where modes 0 and 1 apply a simple copy of the neighboring blocks and modes 3–8 perform a weighted interpolation according to the prediction direction. Mode 2 (DC) replicates the average of the neighboring samples to the entire block. Each one of the 16 blocks inside the MB may use different prediction directions in order to find the best prediction.

The intra-prediction can also be performed using the 16×16 block size. However, in this mode the number of prediction directions is restricted. Figure 2.9 presents the four prediction directions. Modes 0–2 are analogous to the modes 0–2 of the 4×4 block size. The plane mode (3) applies one linear filtering (Richardson 2010) to the neighboring samples resulting in a gradient texture. The 4×4 and 16×16 presented predictions are used for luminance samples. The chrominance prediction uses the same four directions present in 16×16 intra-prediction. The block size depends on the color subsampling; for the 4:2:0 color subsampling, 8×8 chroma blocks are used.

The inter-frame prediction or motion estimation (ME) provides other possibility of prediction. Its function is to perform a search in the past or future previously encoded frames to find the best matching candidate in order to provide a good prediction. The ME features bi-prediction, multiple block sizes, motion vector prediction, $\frac{1}{4}$ sample motion vector accuracy, weighted prediction, and other tools that help to improve the prediction quality (Richardson 2010), as defined in the H.264/MVC video coding standard and detailed in Sect. 2.3.2.

For the dependent views (all views except the base one), the inter-view prediction or disparity estimation (DE) is also available (Merkle et al. 2007). This MVC

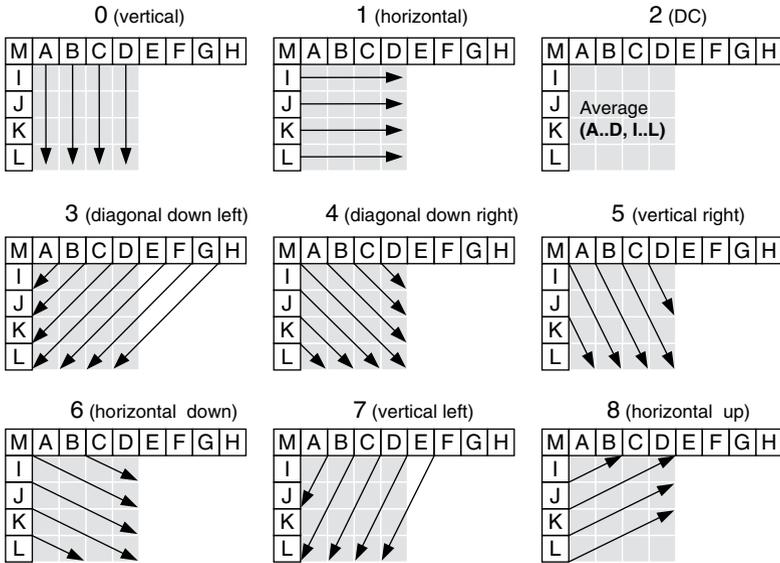


Fig. 2.8 Nine prediction directions for intra-prediction 4 × 4

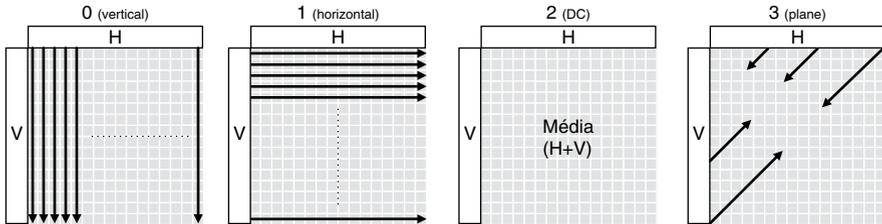


Fig. 2.9 Four prediction directions for intra-prediction 16 × 16

extension searches for the best matching candidate in the frames belonging to previous encoded views (left, right, up, or down, depending on the cameras arrangement and view prediction structure). All features from ME are supported in DE; more details about these features and how they influence the encoder efficiency and complexity will be discussed in Sects. 2.3.2 and 3.1.1.

The output of the prediction phase is a large set of prediction candidates. Among all different block sizes for intra-prediction, inter-frame prediction, and inter-view prediction, the best prediction mode must be selected by the mode decision (MD) in order to provide the optimal rate–distortion (RD) trade-off (Richardson 2010). The rate is the number of bits required to encode the MB and distortion is the objective video quality measured in peak signal-to-noise ratio (PSNR). To have the optimal solution all modes must be completely encoded, reconstructed, and evaluated according to an RD optimization equation. Therefore, the MD (represented by the selection key in Fig. 2.6) is of key importance since it controls the quality vs.

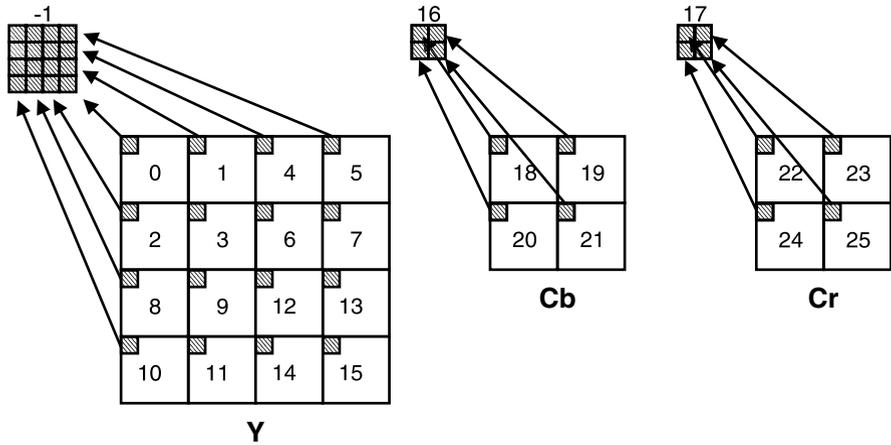
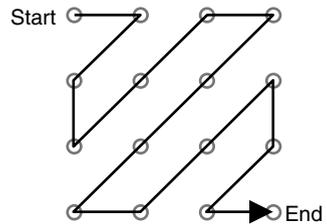


Fig. 2.10 Block processing order in the transform module

Fig. 2.11 Zigzag scan order for a 4x4 block



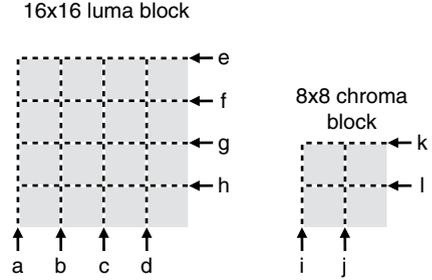
efficiency trade-off and the complexity of the encoder. The MD optimization process is discussed in Sect. 2.3.3.

After the prediction phase is completed, the predicted macroblock and the original macroblock are subtracted to generate the image residues. To reduce the energy in a few coefficients the residues are transformed from the space domain to the frequency domain using an integer approximation of the 4x4 2D-DCT transform. If the intra-prediction 16x16 is selected, an additional Hadamard transform is applied after the DCT. In this case, the DC coefficients of each 4x4 block (left upper corner of each block as depicted in Fig. 2.10) compose another 4x4 coefficient block and are submitted to a 4x4 Hadamard transform. The values inside each block in Fig. 2.10 represent the double-Z processing order of the blocks in the transform.

Once the transforms are concluded, each block is quantized to reduce the dynamic range of the coefficients for the entropy coding. In MVC a linear quantization is used. The quantization step is defined by the H.264/MVC standard (Richardson 2010).

Finally, the quantized coefficients are sent to the entropy encoder. Each block is scanned in zigzag order, according to Fig. 2.11, and encoded by one of the two standard entropy encoders: CABAC or CAVLC. The CAVLC use predefined tables

Fig. 2.12 Order of macroblock borders filtering



depending on the syntax element being encoded. The coding method is an evolution of variable length coding to better adapt to multiple contexts. The context-adaptive binary arithmetic coding (CABAC) is a new tool defined by the H.264/AVC standard and implements a novel coding technique able to reduce the bitstream size by about 5–15 % (Wiegand et al. 2003) in comparison to the CAVLC encoder. The tables of probability used in CABAC are updated at bit level and present strong data dependencies. For further information please refer to (JVT 2009a, b; Richardson 2010).

After the entropy coding, the bitstream is assembled and the encoding is complete. However, every macroblock has to be reconstructed to work as reference for further MBs. For that, the inverse quantization and inverse transforms are applied to the quantized coefficients (the same data previously sent to the entropy).

Once the residues are inversely quantized, they are added to the predicted block in order to reconstruct the decoded MB. The reconstruction loop guarantees the consistency between encoder and decoder sides avoiding drifting between encoder and decoder. To reduce the blocky effect (due to different prediction modes) in the reconstructed frames, the standard defines an in-loop deblocking filter (DF). The filtered MBs are used for displaying and to generate the reference for inter-frame and inter-view predictions. Intra-prediction uses unfiltered macroblocks inside a frame. The DF has five filtering strengths and filters the borders of each 4×4 block of the image following the order presented in Fig. 2.12 (Richardson 2010).

2.3.2 Motion and Disparity Estimation

Multiview video sequences are usually captured using a high sample rate, over 30 fps, to improve the motion flow and give the observer a sense of smoother motion. This high frame rate implies in a high redundancy or similarity between neighboring frames in the time axis. As noticed in Fig. 2.13, frames S0T0 and S0T1 are very similar; hence only the differences between them have to be transmitted. The algorithm that exploits these inter-frame similarities is the motion estimation (ME). It searches in the temporal neighboring frames, known as reference frames (see Fig. 2.14), the region that represents the best match for the current block or macroblock. Once the best matching block is found, a vector pointing to that

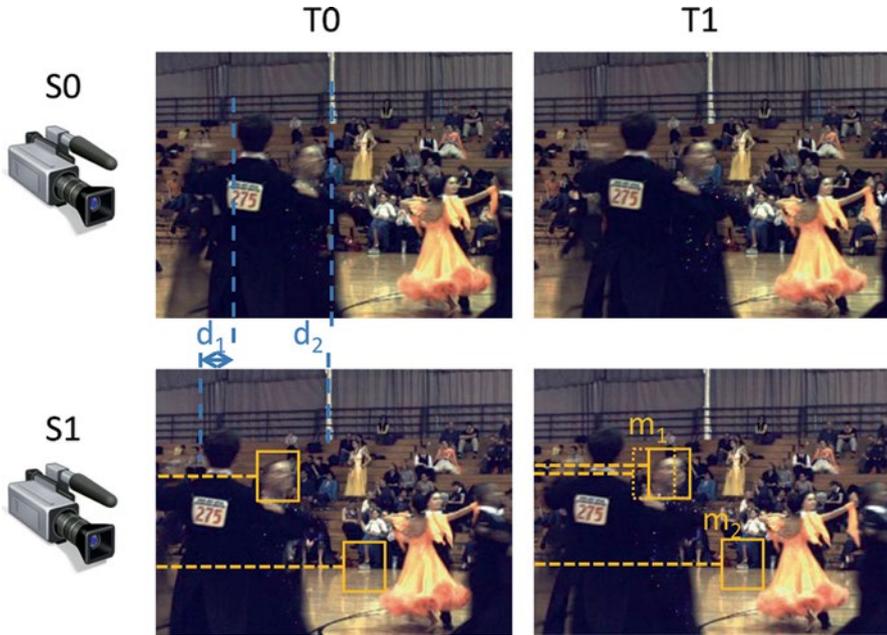


Fig. 2.13 Temporal and disparity similarities

position, the motion vector (MV) in Fig. 2.14, is generated. Consider, for example, a background region (one of the yellow boxes in Fig. 2.13), there is no motion between T0 and T1, so the motion vector m_2 is probably zero. The dancers moving (woman’s face in the yellow box) present a displacement along the time; this displacement is represented by m_1 . The set of motion vectors of a given frame are called motion field and represent valuable information to understand the motion of an object as time progresses.

The cameras that capture the different sequences in a given 3D scene are located near each other (typically, about 5–20 cm apart) (Su et al. 2006); thus there are many regions that are shared between neighboring cameras. A very high similarity is perceived between neighboring cameras as exemplified in frames S0T0 and S1T0 of Fig. 2.13. The MVC defines the disparity estimation (DE) to exploit the redundancy between different views and minimize the transmission of replicated information multiple times. The approach of the DE is similar to the ME. It searches for the best matching candidate block within frames of the neighboring views. The frame used for search is called reference frame while the view is called reference view, as shown in Fig. 2.14. Once the matching block is found the position is pointed by the so-called disparity vector (DV); see Fig. 2.14. The set of DVs in a frame are referred as disparity field and represent the disparity of the objects between views. While the length of motion vectors (MV) represents the speed an object is moving (or the camera is moving) the disparity vectors denote the displacement of a given object between two views. The disparity depends on the distance between cameras,

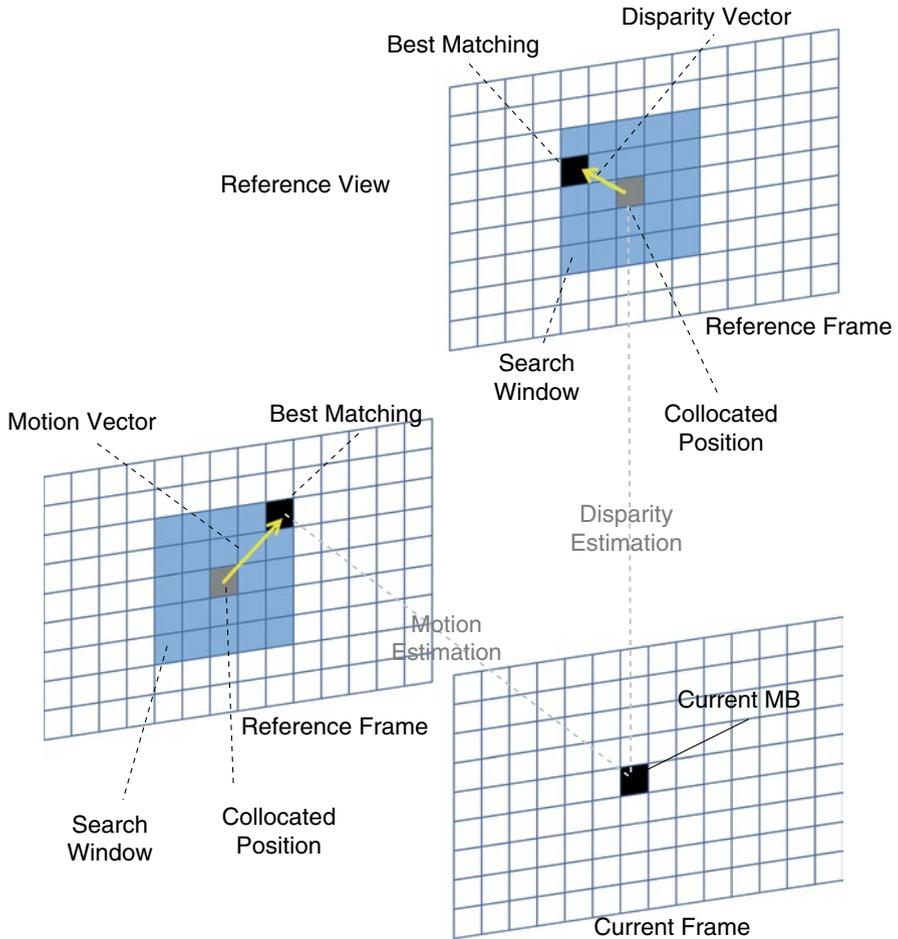


Fig. 2.14 Motion and disparity estimation

and the distance between the camera and the object (Kauff et al. 2007). The closer the object is the larger the displacement or disparity. For instance, in Fig. 2.13, the background presents almost no disparity between S_0 and S_1 (d_2) while the dancers have a much larger disparity vector (d_1). The average disparity vector between two views considering all objects and background is named global disparity vector (GDV) (Kauff et al. 2007; Han and Lee 2008); see Fig. 2.7.

The ME/DE search is not performed over the complete reference frame but in a region called search window (SW) defined by a search range (SR), as shown in Fig. 2.14, for instance an SR $[\pm 16, \pm 16]$ covers an SW of 33×33 samples. Many search schemes for ME were proposed along the last two decades and their characteristics are well known. The exhaustive search algorithm, the *Full Search* (FS) (Yang 2009), provides the optimal results at the cost of a very high computational

effort. Many fast algorithms focusing on complexity reduction with small quality loss are found such as *Log Search* (JVT 2009a, b), *Diamond Search* (DS) (Kuhn 1999), *Three Step Search* (TSS) (Jing and Chau 2004a), *UMHexagon Search* (Chen et al. 2002) and *Enhanced Predictive Zonal Search* (EPZS) (Tourapis 2002), to list a few. These algorithms are based on multiple search steps oriented by geometric shapes. The most recent schemes also consider the neighboring MBs as predictors to define the search starting point. Using predicted starting point is an evolution compared to the previous search schemes that use the collocated MB as starting point. Recalling, the collocated MB is the macroblock in the reference frame that belongs to the same relative position of the current MB.

Despite the similarity between ME and DE there are behavioral differences that make solutions defined for ME inefficient when applied to DE. For instance, most of the traditional ME fast search patterns perform badly for DE. The reason is that motion vectors are usually located in a relative small length range while disparity vectors usually are much longer. The disparity vectors frequently have 50–100 samples length. For this reason, the recommended search range is at least $[\pm 96, \pm 96]$ for SD resolutions (Xu and He 2008). In this scenario most of the fast algorithms tend to fall in local minima and do not find the optimal candidate. For this reason the JMVC, the reference software for MVC (JVT 2009a, b), implements the *TZ Search* that is more complex in comparison to DS and EPZS, for example, but is still 23× times faster than FS (Yang 2009). The TZ employs predictor centered search start and a larger geometric shape search pattern. It performs well for both ME/DE with negligible or no quality loss in comparison to FS.

However, once the conceptual tasks of ME and DE are similar, the available features are the same and together they represent the most computational and memory intensive tasks in the video encoder; see discussion in Sect. 3.1. For this reason, ME/DE have to be jointly considered in order to propose smart fast algorithms and efficient *architectural solutions for real-time MVC encoding*.

In the following, the motion and disparity estimation features are detailed. Note that all these tools are mandatory at the decoder side depending upon the operation profile but are optional for the encoder.

Bi-prediction: In MVC, two types of MBs employ the ME/DE: Predictive (P), which is coded using inter-frame prediction referencing only past frames and backward views, in display order, or bi-predictive (B), which is coded using reference frames both from past/backward and from future/forward (this is possible due to the out-of-order coding and decoding allowed by the standard). In a B macroblock, each partition can be predicted from one or two reference frames (Sullivan and Wiegand 2005). In case of bi-prediction the final prediction is generated by calculating the average of the prediction from past/backward and future/forward.

The reference frames are stored in two lists: List 0 and List 1. List 0 orders the frames from the past and backward views and List 1 orders the frames from the future and forward views (JVT 2003). Both lists can be ordered using temporal references first or disparity references first. For temporal references first, in List 0 the reference index 0 is the closest past encoded frame. For disparity references

first, the index 0 in List 0 is the closest backward view reference frame. Analogous organization is observer in List 1.

Multiple block sizes: MVC allows ME/DE blocks of several sizes. The 16×16 MB can be segmented in two 16×8 , two 8×16 , or four 8×8 partitions (JVT 2009a, b). Each 8×8 partition can be segmented in other two 8×4 , two 4×8 , or four 4×4 sub-partitions. Each partition may point to one reference frame per list (List 0 and List 1) while each sub-partition may use only the frames referenced by the partition that it belongs. Each partition or sub-partition may have a single MV or DV.

Multiple reference frames and reference views: Differently from earlier standards, in MVC the past and future reference frames are not only fixed to the immediate ones. Therefore, to reconstruct one given macroblock, temporally distant frames can be used in the prediction process and this distance is limited only by the size of the decoded picture buffer (DPB) (Sullivan and Wiegand 2005). The reference frames are managed in List 0 and List 1 as previously cited. Analogously, the reference views are not restricted to the closest backward or forward views, any previously encoded views may be used as reference depending on the coding settings.

Quarter-sample motion vector accuracy: In general, the motion of blocks does not match exactly in the integer grid of pixels in a frame, and then fractional-sample motion vector accuracy is used to reach a better match. The MVC (JVT 2003) defines the use of a quarter-sample motion compensation for the reference frame blocks. For luma samples, a six-tap FIR filter is used to interpolate half-samples, and then a simple average of integer and generated half-samples is used to generate the quarter-sample interpolation (JVT 2003). When working with 4:2:0 subsampling, the chroma samples interpolation applies 1/8 sample accuracy.

Weighted prediction: The MVC defines a weighted prediction in the inter-frame coding process to apply a multiplicative weighting factor and an additive offset to each interpolated sample of a given reference frame. For single directional prediction from List 0 or List 1 this tool is defined as presented in Eq. (2.1), where “ x ” is replaced by the list number (0 or 1), “ w ” is the weighting factor, “ $\log WD$ ” is a scaling factor, and “ o ” is the additive offset. P represents the interpolated pixels and P' the weighted sample. For bi-predictive prediction the weighted prediction is defined as presented in Eq. (2.2):

$$P'(i, j) = ((P_x(i, j) \times w_x + 2^{\log WD - 1}) \gg \log WD) + o_x, \quad (2.1)$$

$$P'(i, j) = ((P_0(i, j) \times w_0 + P_1(i, j) \times w_1 + 2^{\log WD - 1}) \gg (\log WD + 1)) + ((o_0 + o_1 + 1) \gg 1). \quad (2.2)$$

Motion/disparity vector prediction: Exploiting the neighboring blocks correlation, the MVC standard defines that motion vectors and reference indexes (pointer to the reference frame in List 0 or List 1) have to be inferred from the reference index and motion/disparity vectors of neighboring blocks. The inferred vectors are called predicted motion vectors (PMV). Differential motion vectors (MVD) are coded in the bitstream and summed up to the PMVs, obtaining the current motion vector (MV) or disparity vector (DV). The PMVs are normally obtained applying the median to

the spatial neighbor blocks vectors. However, SKIP macroblocks and direct predicted macroblocks (macroblocks with no transmitted residue or motion vectors) are differently processed using the direct spatial or direct temporal predictions. The motion/disparity vector prediction is one example of using the video correlation to predict coding side information, as previously mentioned in Sect. 2.2.

2.3.3 MVC Mode Decision

The MVC provides a big number of options for the macroblocks prediction. Intra-prediction defines two prediction sizes (three in case FRExt is considered), 16×16 and 4×4 , with four and nine prediction modes, respectively. ME evaluates multiple candidate blocks for seven different block sizes. Additionally, the new disparity estimation adds a set of coding possibilities as large as the motion estimation possibilities.

The mode decision (MD) module is the responsible to deal with this large optimization space. For that it implements an optimization algorithm and defines a cost function called RDCost, the rate–distortion cost (a.k.a. J cost). The objective is to evaluate the coding modes and to find the one that minimizes the RDCost to obtain the best coding relation between rate and distortion. Equation (2.3) presents the J function where c and r represent the current original MB and the reconstructed one, $MODE$ is the prediction mode used, and QP is the quantization parameter. D represents the distortion measured after the complete MB reconstruction according to a distortion metric and R is the number of bits used to encode the current MB; this number is available once the entropy encoding is completed. λ is the Lagrange Multiplier used to control the rate–distortion trade-off. The Lagrange Multiplier value is not defined by the standard; however, typically it is defined by the Eq. (2.4) and depends upon the QP. To quantify the distortion, different metrics may be used; some examples are sum of absolute differences (SAD), sum of absolute transformed differences (SATD), and sum of square errors (SSE). The SSE is mostly used in the mode decision step since it provides better PSNR results. The reason is that PSNR is calculated using mean square errors (MSE) which is only a division of SSE value, so the SSE is directly related to PSNR. It is important to understand that PSNR is currently the widely most accepted objective video quality metric. However, SAD is widely used in real-time systems due to its light-weight computation:

$$J(c, r, Mode | QP) = D(c, r, Mode | QP) + \lambda_{Mode} \times R(c, r, Mode | QP), \quad (2.3)$$

$$\lambda = 0.85 \times 2^{(QP-12)/3}. \quad (2.4)$$

Although the algorithm to find the mode that minimizes the RDCost is not defined by the standard, the MVC reference software JMVC implements an exhaustive search by completely encoding all possible coding modes and selecting the best mode. It is known as rate–distortion optimized mode decision (RDO-MD) also referred as Full RDO or Exhaustive RDO. The RDO-MD guarantees the optimal MB encoding but drastically increases the encoder computational effort and makes the same approach for real-time MVC encoding unfeasible for the current technology.

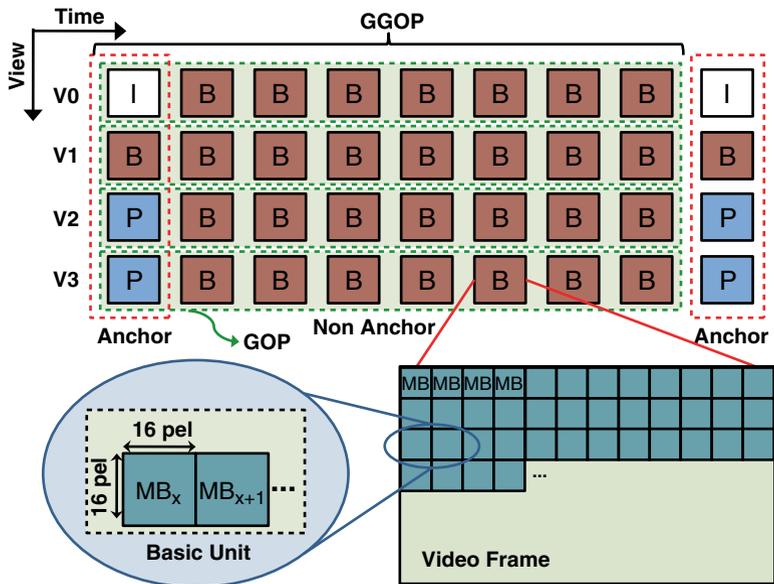


Fig. 2.15 MVC rate control actuation levels

2.3.4 MVC Rate Control

According to, Li et al. (2003) the rate control (RC) is a block of the video encoder that aims to regulate the output-coded bitstream to produce high video quality at a given target bitrate. In the MVC scenario, an efficient RC scheme must be able to provide increased video quality for a given target bitrate with smooth visual quality variation along the time, for different views and within the frames. Most importantly, the RC should keep the bitrate as close as possible to the target bitrate (optimizing the bandwidth usage) while avoiding sudden bitrate variations.

The rate control unit typically controls the quality vs. bitrate through QP adaptation. The bitrate and/or the video distortion metric are predicted using a prediction model. According to the prediction and the target bitrate (amount of bits per second used to encode the video), an adequate QP is selected. As QP grows higher, more residual data are quantized (video details lost) and more quality losses are inserted. The actual generated bitrate and the video quality may be used as feedback for the RC unit in order to update the prediction and QP definition. The QP adaptation may be performed in distinct actuation levels. In general, the RC for MVC can be classified in at least three actuation levels (1) GOP level (group of pictures—set of frames), (2) frame level, and (3) basic unit (BU—set of one or more macroblocks *MB*) level, as shown in Fig. 2.15. It is possible to combine GOP level and frame level, and based on this observation, for simplicity, they are jointly discussed in this monograph.

In the following sections we present the state of the art related to 3D-video systems, MVC encoders, and multimedia processing. Also, a literature overview on the latest low-complexity and energy-efficient solutions focusing on mode decision,

ME/DE and rate control for the MVC standard is presented. An overview on low-power techniques is also provided to give the technical background required for our energy-efficient architectural solutions.

2.4 3D-Video Systems

The advances in video coding techniques targeting multiview videos have been driven by the increasing set of commercial systems employing 3D-video capabilities. These systems range from high-end cinemas and 3DTVs to mobile devices including content suppliers. Wider adoption is expected for the upcoming years with the increase in the available video content through 3D-capable television broadcasters, optical media (Blu-ray Disc Association 2010), popularization of personal 3D camcorders, 3D-video stream services (YouTube 3D 2011; Vimeo 2012), etc. All these commercial systems are, currently, based on stereoscopic videos (two views). An increase in the number of views is expected for the near future (Fujii 2010) to improve the observer freedom and provide a more immersive experience. Some experimental and academic multiview systems are already available or under development to support the next generations of 3D-video systems. In this section we start presenting the most prominent commercial 3D-video systems.

3D-cinema systems are based on three market-leader technologies based on stereoscopic videos (IMAX 2012; RealD 2012; Dolby 2012). The technology employed in IMAX (2012) requires the use of linear polarized glasses to block the light for one eye at a time allowing each eye to see only the frames intended for that eye. Two projectors are used to display 48 fps where each eye is able to effectively see 24 fps in a time-sharing strategy. RealD (2012) is also based on time-sharing between the two eyes; however, the glasses are circular polarized glasses where each glass is polarized in an opposite direction. Also, RealD (2012) requires a single projector able to display 144 fps. Each effective frame for each eye is displayed thrice resulting in effective 24 fps per eye. Finally, Dolby (2012) employs passive glasses with dichroic filters where each view is displayed with a distinct chromatic filter and perceived by a single eye. With this strategy both views are simultaneously displayed allowing the use of regular 24 fps projectors. At the video coding perspective, all these high-end applications support the stereoscopic MVC (The Digital Entertainment Group 2009).

Stimulated by the content available in the 3D Blu-Ray (Blu-ray Disc Association 2010)—optical media that supports the MVC coding standard—the 3D televisions already exceed 10 % of the televisions sold in the United States, in 2011, and this number is expected to reach 37 % of the market share in 2014 (Research and Markets 2010). Other countries are expected to follow this trend. The majority of those 3DTVs are based on stereoscopic displaying and require active shutter or passive polarized glasses to provide the 3D sensation. Many devices employ built-in decoders supporting the MVC standard. Along with the cinema solutions, the 3D televisions are not energy-critical and typically implement only the video decoder side, less complex in relation to the encoder.

Currently, portable devices capable of handling digital video are available everywhere for a reasonable cost. The omnipresence of these gadgets implies a very large amount of data being produced. In this scenario the coding efficiency is a key issue in order to reduce the storage and transmission costs for digital video. Various devices are also capable of real-time 3D-video recording, such as Panasonic (2011), Fujifilm (2011), Sharp (2011), and Sony (2011). Most of them feature two cameras and encode the video sequences independently (simulcast). However, the increase in number of views from 2 up to 4–8 views (Fujii 2010) in order to provide enhanced 3D experience freedom is envisaged for the next 3–5 years. In this scenario it is simple to conclude that the large amount of data generated requires the use of the state-of-the-art MVC standard. The first personal camcorder to fully support stereo MVC was released by Sony in 2011.

Although 3D-capable mobile devices are already available, attending MVC performance and energy constraints remains a big challenge for industry and academia, as discussed in Sect. 3.2. The current multimedia processing systems based on processors, DSPs, and non-MVC-optimized application-specific integrated circuits (ASIC) implementations are not efficient to provide the required throughput with the required energy efficiency while sustaining video quality and coding efficiency. In the following section we present an overview of the main multimedia architectural approaches and solutions in the current state of the art.

2.5 Multimedia Architectures Overview

In this section we present the multimedia processing architectures classified in four main classes: Multimedia Processors/DSPs, Reconfigurable Multimedia Processors, ASIC Multimedia cores, and Heterogeneous Multicore SoCs. On the one hand, ASIC solutions provide the highest performance and energy efficiency at the cost of reduced flexibility limiting the applicability to upcoming video standards. Still, the current lack of MVC-oriented ASIC optimizations prohibits further increase in both performance and energy efficiency. On the other hand, multimedia processors/DSPs allow high flexibility to multiple standards while providing reduced performance and poor energy efficiency if compared to ASICs. Additionally, reconfigurable processors may allow significant increase in performance and flexibility through instruction set architecture (ISA) extensions. The reconfigurable processors, however, present reconfiguration energy issues and are unable to reach the ASIC-like performance and energy efficiency required by the 3D multimedia applications.

2.5.1 *Multimedia Processors/DSPs*

Aware of multimedia processing characteristics, the Multimedia Processors/DSPs are designed to exploit the parallelism inherent to these applications. Massive multicore architectures are proposed to target task parallelism by supporting multiple

parallel threads. Data-level and instruction-level parallelisms are exploited by employing single instruction multiple data (SIMD) and very large instruction word (VLIW) architectures, respectively. Some proposals are able to implement hybrid parallelism by handling multiple cores with SIMD and/or VLIW instruction sets.

In Abbo et al. (2008), the Xetal-II employs 320 SIMD processing elements with a dedicated 10 Mb on-chip frame memory. It is able to provide 107 GOPS with a 60 W power consumption with instructions designed targeting video analysis applications. A multicore system for video decoding is proposed in Finchelstein et al. (2009) employing a caching mechanism to reduce the memory reads. The work presented in Khailany et al. (2008) describes a processor with 16 parallel lanes where each lane is a 5-ALU VLIW core. At 800 MHz, this solution delivers 512 GOPS (82 pJ/MAC) and guarantees baseline HD1080p H.264 encoding at 30 fps. The multi-streaming SIMD multimedia engine proposed in Chiu and Chou (2010) claims a 3.3–5.5× performance increase compared to MMX architecture (Intel Multimedia Extension) by employing 12 multimedia kernels. These parallel architectures provide a relative high performance but are still far below MVC requirements and the power envelope is out of embedded devices boundaries.

A 2-issue VLIW stream processor is presented in Chien et al. (2008) with throughput for CIF encoding at 30 fps. Stereo processing-oriented optimizations for VLIW processors are presented in Payá-Vayá et al. (2010). The authors claim performance improvements by implementing a new register file access mechanism and disparity functional unit to calculate disparity map. Also, an application-specific instruction processor (ASIP) based on a VLIW DSP architecture is described in Zhang et al. (2009) and delivers increased performance if compared to traditional DSP and SIMD.

2.5.2 Reconfigurable Processors for Video Processing

In Otero et al. (2010) an architectural template for run-time scalable systolic coprocessors is presented. It focuses on run-time adaptation to time-variable tasks or changing system conditions. It exploits replacing and relocation of basic processing elements of the array using FPGAs dynamic reconfiguration. In Beck et al. (2008) is employed a coarse-grained reconfigurable array with a run-time mechanism designed to translate MIPS instruction to be executed in the reconfigurable array. Berekovic et al. (2008) present the mapping of MPEG-2 and H.264/AVC to the ADRES (coarse-grain reconfigurable processor) delivering throughput for real-time CIF decoding at 50 MHz with a 4×4-core array. CRISP, a coarse grain reconfigurable stream processor (Chen and Chien 2008), implements an image processing pipeline reaching 55 fps for HD1080p resolution. Aggressive performance losses are expected for video coding due to increased complexity compared to the implemented image processing algorithms.

In Bauer et al. (2007), the rotating instruction set processing platform (RISPP) is presented bringing more flexibility to extensible processors. It features a special instruction forecasting algorithm able to predict the hotspots and allows to adapt at

run time the different *Molecules* (implementation of the special instructions). This architecture was evaluated using some H.264 processing hotspots (SATD, DCT, etc.) and demonstrated high flexibility to deal with the performance vs. hardware trade-off. This concept was extended in Bauer et al. (2008a) by integrating a special instruction run-time scheduler able to outperform the state of the art in 2.38× for the H.264 application. When integrated to a transmutable embedded processor (Bauer et al. 2008b), the RISPP concept was able to present up to 7.19× speedup in relation to related works for H.264.

Compared to regular processors, reconfigurable processors target to increase the overall performance by adapting, at run time, to distinct applications properties. Also, the adaptivity can be efficiently exploited within the same application. Considering multimedia applications, the performance/energy requirements may vary with the video content, user settings, battery level, etc. It brings a big optimization potential at the system perspective. However, when considering a single application for a given description, in this case real-time encoding for MVC HD1080p, the profit of this adaptive behavior is not perceived. Moreover, in this scenario, the energy and time costs for reconfiguration pose additional difficulties in terms of throughput and energy efficiency if compared to processors, DSPs, and ASIPs.

2.5.3 *Application-Specific Integrated Circuits*

Multiple ASIC hardware architectures were proposed targeting real-time high-definition (de)coding in accordance to the latest video coding standards trying to reduce the total energy consumption for embedded devices. The architecture presented in Chen et al. (2009c) delivers H.264 encoding for D1 (720×480) resolution with 43.5–67.3 mW consumption. In Lin et al. (2008) an H.264 video encoder able to process HD1080p sequences at 242 mW is presented. Chang et al. (2009) propose a real-time 720p H.264 encoder at 183 mW consumption. It implements a 3-stage pipeline, 8-pixel intra-prediction parallelism, and a parallelized subsampling algorithm. A 59.5 mW AVC/SVC/MVC decoder for up to three-view HD1080p videos is presented in Chuang et al. (2010). The first complete video encoding solution for MVC encoding was presented in Ding et al. (2010). The AVC/MVC 3D/Quad Full HDTV supports three-views HD1080p and consumes 522 mW.

Compared to other approaches, ASIC provides high throughput and energy efficiency. Considering the state-of-the-art IC manufacturing technology, ASIC implementation is the only solution able to encode high-definition MVC for an increased number of views at real time, as shown in the related work overview presented. Still, further optimizations are possible in relation to the presented ASIC-related works. Having Ding et al. (2010a, b) as comparison basis, the future MVC encoding systems will require increased number of views. Moreover, in Ding et al. (2010a, b) single-view-based optimization techniques (such as search window reduction that seriously affects the disparity estimation) are employed leaving a high potential for multiview-aware optimizations.

2.5.4 *Heterogeneous Multicore SoCs*

Heterogeneous multicore architectures are also proposed targeting multimedia applications. In Kollig et al. (2009), the proposed systems handle an ASIC HW video codec, audio codecs, VLIW processors, MIPS host CPU, DSP, and other HW accelerators. The system proposed by Kondo et al. (2009) is composed of two specific accelerators (video decoder and descriptor), one general accelerator (MX), three RISC CPUs, and caches. Woo et al. (2008) describe a 195 mW mobile multimedia SoC with ARM 9, AVC/MPEG decoder, JPEG codec, fully programmable 3D engine, and multiple peripheral interfaces. The heterogeneous multicore approach is the most used in the current set-top boxes, digital TV decoders, and smart mobile devices. It takes advantage of some degree of flexibility along with the performance of specific accelerators.

The SoCs in current commercial mobile devices such as smartphones and tablets implement heterogeneous multicore SoCs employing processors with SIMD extensions, DSPs, ASIC codecs or hardware accelerators, and programmable embedded GPUs. Qualcomm Snapdragon S4 (Qualcomm Inc. 2011) is composed of up to 4 ARM cores, Hexagon DSPs, video coding hardware accelerators, and the Adreno embedded GPU. Nvidia Tegra 3 (Nvidia Corp. 2012) is based on up to 4 ARM cores and employs dedicated video encoder/decoder and ULP GeForce GPU. Samsung Exynos 4 (Samsung Electronics Co. Ltd. 2012) is composed of quad-core ARM processor, video/image/audio ASIC codecs, and the ARM Mali GPU. Texas Instruments OMAP 5 (Texas Instruments Inc. 2012) employs 2 ARM Cortex-A15, 2 ARM Cortex-M4, DSPs, video/audio accelerators, and the PowerVR GPU. Note, even with efficient ARM processors, SIMD extensions, DSPs, and programmable massively parallel GPUs, the embedded SoCs require ASIC codecs/acceleration units to deliver the throughput and energy efficiency for real-time high-definition video encoding. To deal with multiview videos and attend performance/energy requirements on embedded battery-powered devices, these SoCs will require MVC-oriented optimizations at algorithmic and architectural (including datapath and application-aware units/memory management optimizations) levels.

2.6 Energy-Efficient Architectures for Multimedia Processing

In this section we introduce the state of the art on energy management along with an overview on video memories, energy-efficient techniques, and architectures for multimedia processing. Additionally, the infrastructure to support dynamic voltage scaling (DVS) on SRAM memories and the dynamic power management (DPM) schemes are presented. This technique is extensively used in the literature and in the solutions proposed along this monograph.

2.6.1 *Video Memories*

On-chip memories are becoming dominant part of current systems, mainly for signal processing systems. In the scope of video coding, the video memory represents the main on-chip memory component responsible for storing frames used as reference to encode other frames. In the current literature are found solutions specific for video/frame memories or generic solutions for any video/image processing tasks. Some of these solutions are described in the following.

The work in Grun et al. (1998) proposes a memory size estimation method for applications containing multidimensional arrays such as video processing. The memory estimation is generated from the application algorithm specification. The paper also addresses the discussion relating parallelism to the memory size. Zhu et al. (2006) present a memory size computation method for multimedia algorithms. The solution uses algebraic techniques and the theory of integral polyhedral to compute exactly the memory size for multimedia algorithms. The authors in Yamaoka et al. (2005) use a triple-mode SRAM to implement an on-chip memory for mobile phone application. The on-chip memory is composed of four SRAM banks that can be managed by a leakage state controller.

In terms of specific video memories, the authors of Shim and Kyung (2009) propose a video memory composed of multiple on-chip memories employing a data reuse to reduce the external memory access. A memory switching method is defined to increase the utilization of on-chip memory. Tsai et al. (2007) present a low-power cache for the reference frame memory of H.264/AVC. This work uses the block translation cache architecture and a search trajectory prediction prefetching scheme. The authors claim a 35 % memory writing power reduction with 67 % memory static power reduction.

2.6.2 *SRAM Dynamic Voltage-Scaling Infrastructure*

The static energy due to leakage current represents a significant source of the total energy consumption in deep submicron technologies. Also, current integrated circuit footprints are dominated by embedded memories which are typically implemented as SRAM (static random access memory). Therefore, reducing SRAM static consumption is a key challenge to reach overall energy reduction.

The fabrication technology evolution has provided meaningful contribution to leakage reduction by employing high-K oxides (Huff and Gilmer 2004), FinFET transistors (Pei et al. 2002), etc. Along this monograph we assume the use of an on-chip SRAM memory featuring multiple power states with data retention capabilities. The high-level memory organization is presented in Fig. 2.16. An implementation for this memory organization including a picture of the silicon die is demonstrated in Zhang et al. (2005). Still, there is a need to further reduce the leakage at architecture and system levels through techniques such as power gating, DVS, and DPM. In Sects. 2.6.3 and 2.6.4, we present an overview of power/energy management techniques for memories and multimedia systems, respectively.

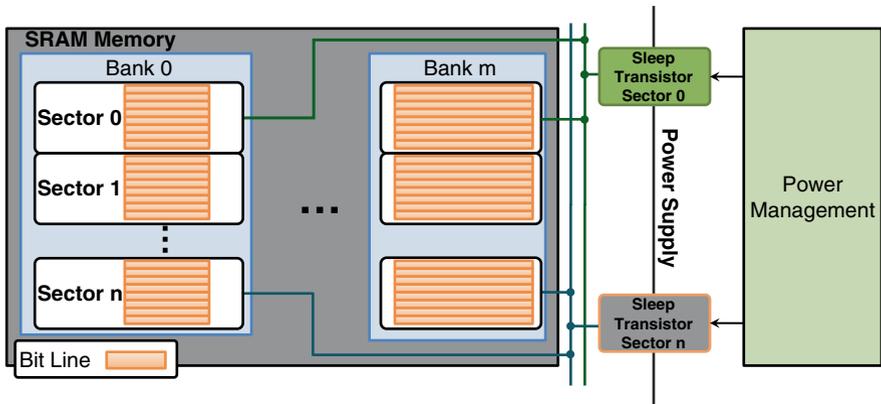


Fig. 2.16 SRAM voltage-scaling infrastructure

2.6.3 Dynamic Power Management for Memories

In order to take advantage of the on-chip memory organizations (such as those presented in Sect. 2.6.1) and multistate SRAM infrastructures (see Sect. 2.6.2) in order to convert these features to actual energy savings, efficient DPM scheme is required. Some DPM solutions in the context of video processing and embedded systems are discussed below.

Panda et al. (1997) present a local memory optimization technique for embedded systems based on memory performance analytical estimation for a given application. The base architecture employs cache and a scratch-pad memories with parameters defined by the proposed technique. In Cong et al. (2009) the memory performance and energy optimization are performed through automatic memory partitioning and scheduling. Firstly, the solution considers the cycle accurate application schedule to meet the memory performance requirements. Secondly, the memory banks are partitioned to reduce dynamic power consumption. The work in Wang and Mishra (2010) implements a dynamic cache reconfiguration technique along with DVS for memory system energy minimization.

Generic techniques for reducing the on-chip SRAM leakage [like (Singh et al. 2007; Agarwal et al. 2006)] propose memories with multiple sleep modes in order to better exploit the leakage vs. wake-up penalty trade-off. State-retentive power gating of register files featuring multiple sleep modes is presented in Roy et al. (2011). In Mondal and Ogrenci Memik (2005) the hardware power gating is controlled by monitoring the underlying hardware. These observation-based techniques may lead to mispredictions, especially in case of sudden variations. The techniques in Liu et al. (2008) and Rajamani et al. (2006) consider application knowledge for a video decoder case study, but they only exploit the knowledge at frame level. These techniques consider longer periods and may not cope with severe variations at the MB level.

In Fukano et al. (2008) a DVS using dual power supply is used to implement a 65 nm SRAM memory employing three operation modes (1) high speed, (2) low power, and (3) sleep mode. In this work the low power and sleep modes are data retentive avoiding data refetching, but it does not support partial DVS for specific sectors of the SRAM. Yamaoka et al. (2004) present a similar solution employing three operation modes while supporting bank-level DVS. It achieves leakage reduction through adapting the virtual supply voltage using PMOS transistors. Finally, the 65-nm SRAM design presented in Zhang et al. (2005) provides more flexibility through adopting multiple power states and fine grain power control. The DVS is controlled at sector level using a custom NMOS sleep transistor to control the virtual ground voltage.

2.6.4 Energy Management for Multimedia Systems

Energy and power management for multimedia systems has been studied in many research works mostly targeting embedded applications. The authors in Cao et al. (2010) employ DVS with five distinct voltage levels. It is controlled using the application-specific knowledge through workload modeling for a wavelet video encoder. In Kamat (2009) a battery level-aware MPEG-4 video encoder with a notification manager and an application-specific controller is presented. Some solutions exploit the energy vs. video quality trade-off at run time to adapt to the system scenario. Ji et al. (2010) partition the input data in distinct profiles used for energy budgeting generating scalable video quality according to the energy scenario. Similar work is presented in Ji et al. (2009) applying game-theory algorithms to control the video encoder. Liang and Ahmad (2009) propose a rate–complexity–distortion model to progressively adjust the H.263+ encoder behavior considering the video content. It employs DVS providing and reaches up to 75 % energy reduction. A power–rate–distortion model (He et al. 2008) is used for energy minimization in video communication devices by exploring energy trade-off between video encoding and wireless communication providing up to 50 % energy reduction. A dynamic quality adjustable H.264 encoder is proposed in Chang et al. (2009). It defines quality states to change the number of coding modes considering the power vs. quality trade-off. The implemented ASIC provides real-time 720p encoding at 183 mW consumption. The proposals summarized in this section are useful at the MVC scenario but lack the MVC-specific knowledge such as workload model, quality states, rate–distortion behavior, etc. Thus, the simple application of these approaches leads to inefficient energy management performance.

Authors in Javed et al. (2011) presented an adaptive pipelined MPSoC for H.264/AVC with a run-time system that exploits the knowledge of macroblock characterization based on their spatial and temporal properties (Shafique et al. 2010; Shafique et al. 2010a) to predict the workload. Based on this knowledge, unused processors are clock-gated or power-gated. These techniques provide limited energy-efficiency in MVC as they cannot exploit the MVC-specific knowledge such as (a) distribution of memory usage at frame and MB levels and (b) memory usage correlation in the 3D-neighborhood.

2.6.5 *Energy-Efficient Video Architectures*

The work of Shafique et al. (2010) presents an energy budgeting scheme for the H.264 ME. This solution considers the total energy available along with the video properties to dynamically define a search pattern able to deal with the energy vs. quality trade-off. Each frame is classified into one of six energy classes and further classification refinement is performed at MB level. The highest complexity class performs a search composed of three search patterns (Octagon Star, Polygon and Diamond) without samples subsampling. The lowest complexity class employs a Diamond-shaped search using 4:1 subsampling. The highest complexity class requires 17× more energy when compared to the lowest complexity class.

The authors in Chen et al. (2006) evaluated different state-of-the-art data reuse schemes (Level-A, Level-B, Level-C, and Level-D) and proposed a new search window-level data reuse for H.264 ME (Level-C+) in order to reduce the energy consumption related to external memory access and on-chip memory storage. Level-A and Level-B solutions are based on candidate blocks. While Level-A fetches and stores on-chip a single candidate block, Level-B fetches a whole candidate stripe (inside the search window). They require frequent external memory access and only fit with regular search patterns which is not the case for state-of-the-art ME/DE algorithms. Level-C and Level-D follow the same logic but at search window level. Level-C stores one search window (avoiding the retransmission of overlapping search window regions accessed by neighboring MBs in the same line) and Level-D a search window stripe for the whole frame. Observe that Level-D requires a extremely large on-chip memory for large search window or frame size. As Level-C presents a reasonable trade-off between external memory access and on-chip memory size it was extended in Level-C+. Level-C only exploits the data reuse between horizontal neighboring MBs. Level-C+ proposes to increase the vertical on-chip storage to include the search window of the MB line below. This allows exploiting the vertical data reuse at the cost of increased on-chip memory and out-of-order processing (two MB lines are processed using double-Z order).

In Wang et al. (2009) a bandwidth-efficient H.264 ME architecture using binary search is proposed. This solution employs a frame-level preprocessing that downsamples the image twice in a factor 2. It results in three images (or three layers), the original image, the downsampled image, and the twice downsampled image. After that, a search is performed in the three layers. This technique is also modified to allow parallel processing and easy hardware implementation. A hardware architecture is presented targeting low power through low memory access, efficient hardware utilization, and low operation frequency.

A complete MVC encoder targeting low-power operation is presented in Ding et al. (2010) employing eight pipeline stages, dual CABAC, and parallel MB interleaving. A cache-based solution is used for the search window reading along with a specific prefetching technique. The cache tags are formed by the frame index and x and y block position. Also, each cache entry stores an image block (instead of words like in generic caches) following the same concept proposed in Zatt et al. (2007). The search is constrained to a $[\pm 16, \pm 16]$ search window with a predicted center point. The ME/DE architecture is described in more details in the previous work

from the same group (Tsung et al. 2009). This approach might lead to quality loss when the center point prediction is not accurate. Also, the authors ignored the fact that fast ME/DE schemes already consider this information to start the search. The MVC encoder is able to real-time encode four views HD720p at 317 mW.

Generally, the search window-based data reuse approaches suffer from excessive leakage resulting from big on-chip SRAM memories required to store the complete rectangular search windows. This point becomes crucial for MVC as the DE requires relatively large search windows (mainly for high resolutions) such as $[\pm 96, \pm 96]$ to accurately predict high disparity regions (Xu and He 2008). In this case, even considering asymmetric search windows incurs in large on-chip storage overhead, thus suffering from significant leakage.

The authors in Shim and Kyung (2009) use multiple on-chip memories to realize a big logical memory or multiple memories (one for each reference frame) according to the frame motion. A search window centered prediction is employed for data prefetching while the size of search window is dynamically adjusted at frame level using the size of motion vectors found in previous frames. The data reuse scheme Level-C is employed.

A data-adaptive structured search window scheme is presented in Saponara and Fanucci (2004). An adaptive window size approach is proposed considering the spatial/temporal correlation of the motion field. If the vectors of the current and past frames do not exceed a given value, there is no need to search in a region larger than this vector size and the fetching of a reduced window is necessary. In case the window is too small and the error starts to increase, a test detects it and the search window is increased regardless of the neighborhood. This solution leads to reduced external memory access, but its potential for on-chip memory reduction is not discussed.

The work in Chen et al. (2007) proposed a candidate-level data reuse scheme and a Four-Stage Search algorithm for ME. Firstly, multiple search start points are predicted from the neighboring MBs motion activity. The predicted points are evaluated and the best one is selected for a Full Search around its position. A ladder-shaped data arrangement is also proposed in order to support random access for the proposed algorithm. The candidates parallel processing is performed using a systolic array.

In Tsai et al. (2007) a caching algorithm is proposed for fast ME. Additionally, a prefetching algorithm based on search path prediction is proposed in order to reduce the number of cache misses. The work (Tsai et al. 2007), however, is limited to a fixed Four Step Search pattern and it does not consider disparity estimation and power gating.

2.7 Energy/Power Consumption Background

Before moving to the discussion related to energy-efficient algorithms and architectures it is necessary to understand the sources of energy consumption and how they might be reduced. Moreover, the energy consumption is directly related to the hardware implementation and only indirectly related to the algorithms. However, it is possible to design algorithms able to result in energy reduction at the hardware level by reducing computational effort, processing time, memory access, etc.

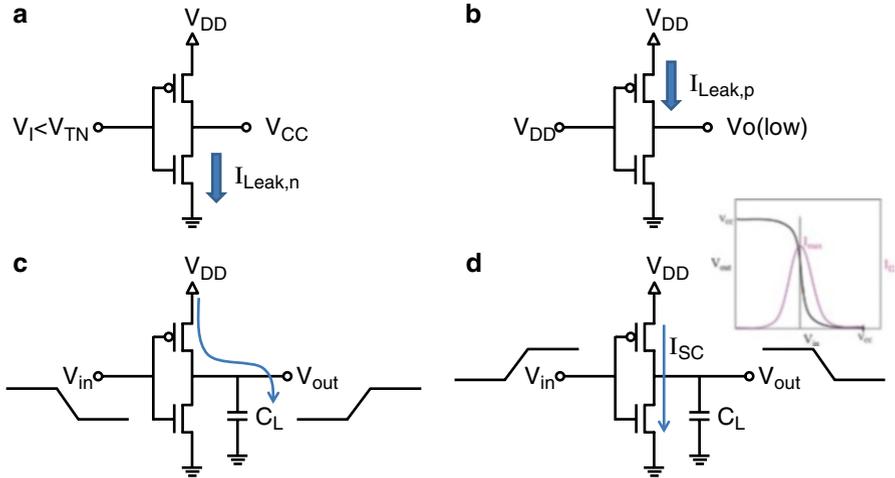


Fig. 2.17 Energy/power dissipation sources

In Fig. 2.17 are represented the three main power dissipation sources for CMOS circuits using an inverter as example: leakage current (static), switching power (dynamic), and short circuit current (dynamic). Eq. (2.5) shows the total power in terms of these three components. The static power dissipation is a result of the leakage currents. Consider Fig. 2.17a where the input voltage (V_I) is lower than the NMOS transistor threshold voltage (V_{TN}). In this case, an ideal inverter NMOS transistor does not conduce any current. However, real MOSFET transistors cannot completely block this current, the so-called leakage current. The closer V_I is to V_{TN} , the stronger the leakage. The same happens to PMOS transistors when a $V_I > V_{TP}$ is applied to the gate (Fig. 2.17b). The leakage power for the case represented in Fig. 2.17b is calculated by Eq. (2.6).

The dynamic power is composed of two components: the switching power (Fig. 2.17c) and the short circuit power (Fig. 2.17d). Equation (2.7) defines the switching power that linearly depends on the load capacitance (C_L , that depends on the fanout of the device), the source voltage V_{DD} , the frequency of operation (f), and the frequency of switching of that device (α). It represents the energy that is charged in the load capacitance and later drained to the ground. Note that only after two switches the energy is actually drained; in the first time instant (shown in Fig. 2.17c) the capacitance C_L is charged and in the second time instant (after another switch) the energy is drained from C_L to ground. It justifies the $\frac{1}{2}$ factor in Eq. (2.7). The short circuit current happens while the input signal changes V_{DD} -GND or GND- V_{DD} . There is a given input voltage where both PMOS and NMOS transistors are conducting and a current is drained directly from V_{DD} to the ground. It is depicted by the current in Fig. 2.17d and the short circuit power is defined by Eq. (2.8). The total energy drained is the total power along the time (t) as represented in Eq. (2.9). Other power dissipation sources (such as gate leakage) exist in the CMOS devices, but they are omitted in this short overview for simplicity reasons.

As can be seen from this overview it is possible to reduce both static and dynamic power. For instance, reducing the computation reduces the dynamic power once α is reduced. If frequency scaling is used, f is also reduced. Moreover, if voltage scaling is used the dynamic energy is reduced in a quadratic order because V_{DD} is reduced. For leakage reduction, circuits featuring multiple thresholds are used. Hardware support is required; however, application knowledge and energy-aware control algorithms are required to accurately control thresholds, frequency, and voltage:

$$P_{Total} = P_{Leak} + P_{Switch} + P_{Short}, \quad (2.5)$$

$$P_{Leak} = I_{Leak} \times V_{DD}, \quad (2.6)$$

$$P_{Switch} = \frac{1}{2} \lambda \times f \times C_L \times V_{DD}^2, \quad (2.7)$$

$$P_{Short} = I_{Short} \times V_{DD}, \quad (2.8)$$

$$E_{Total} = P_{Total} \times t. \quad (2.9)$$

2.8 Energy-Efficient Algorithms for Multiview Video Coding

2.8.1 Energy-Efficient Mode Decision

The mode decision is one of the main contributors for the MVC high complexity and consequent energy consumption. The optimal solution using the exhaustive RDO-MD requires the evaluation of all possible inter-prediction and intra-prediction modes defined by the standard. Such solution is not feasible for real-world implementations. Thus, there is a need to reduce the number of evaluated modes during the coding process. Statically defining modes to be tested does not perform well due to changing coding parameters and video input characteristics. For this reason, it is necessary to dynamically define the most probable coding modes using the run-time available data. Figure 2.18 shows a hypothetical fast MD scheme which selects a few candidate modes out of all possible modes. Current solutions, as detailed along this section, use information extracted from the video content (texture, edges, brightness), coding mode history, video geometry, etc.

Several fast MD schemes have been proposed to reduce single-view H.264 complexity, such as fast I-MB MD, fast SKIP MD, fast P-MB MD, and the combination of the above. These fast mode decisions of H.264 may be deployed for MVC. However, they will perform inefficiently for the non-anchor frames as they do not exploit the inter-view correlation and the knowledge of GDV.

Recently, multiple fast MD schemes have been proposed for MVC (Peng et al. 2008a; Lee et al. 2008; Han and Lee 2008; Shen et al. 2009a, b, 2010a; Ding et al. 2008a; Zeng et al. 2011; Chan and Tang 2012) considering the GDV, camera geometrical properties, inter-view correlation, and early SKIP prediction.

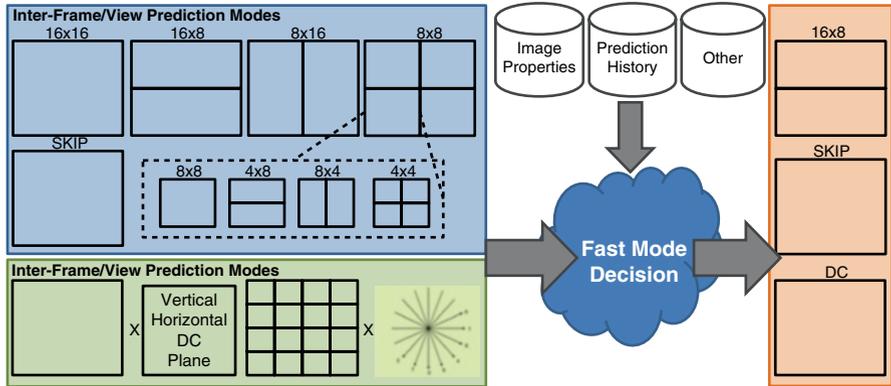


Fig. 2.18 Fast mode decision example

The authors in Lee et al. (2008) proposed an object-based mode decision that uses image segmentation to evaluate different prediction modes for foreground and background regions. The image is segmented using a motion-based approach, considering the vectors size and the SAD in relation to the collocated block (in the same relative position). In case the motion vector significantly defers from the vector average (respecting a threshold) and the SAD exceeds a given value the region is considered as a foreground object; otherwise it is a background. A region growing is used to merge the foreground objects. The foreground regions are coded using DE, while the background are coded using ME. The boundary MBs are coded using the exhaustive RDO-MD.

A fast mode decision based on GDV is presented in Han and Lee (2008). In this scheme the base view—encoded using exhaustive RDO-MD—is used to segment other views in foreground and background regions. The coding modes of the base view are used to classify the image regions. SKIP and Inter 16×16 MBs are defined as background while the remaining modes are considered foreground objects. As the objects present displacement between views, the GDV is used to displace the classified regions as well. Finally, the foreground regions are encoded using exhaustive RDO-MD and the background regions are encoded using big block sizes.

The fast mode decision scheme of Shen et al. (2009a, b) considers the information of reference view to classify the current MB in three complexity classes. For that, the authors propose a mode complexity metric (MDC) defined as the sum of each mode complexity in a 3×3 MBs window. SKIP and Inter 16×16 have “0” complexity, Inter 16×8 and 8×16 have “1” complexity, and Inter 8×8 (or smaller) and Intra have “2” and “3” complexity values, respectively. If the MDC is smaller than a given threshold (T_0), that regions is classified as *simple*. In the opposite, if MDC exceeds another threshold ($T_1 > T_0$) it is classified as *complex*. Regions presenting MDC between these thresholds are defined as *medium* complexity. The *simple* regions test only Inter 16×16 mode. *Medium* regions evaluate Inter 16×16 , 16×8 , and 8×16 modes. *Complex* MBs are encoded using the exhaustive RDO-MD.

In Zeng et al. (2011) a fast mode decision approach is proposed based on the classification of the current MB according to its motion activity based on the coding

modes of the base view. Firstly, the five motion-activity classes are defined in relation to the coding modes. SKIP belong to the motionless class (1). Slow motion class (2) is defined for SKIP and ME 16×16 . ME 16×8 and 8×16 are considered Moderate Motion (3). Fast motion regions (4) are defined by ME 8×8 or smaller. Finally, DE and Intra define High texture with fast motion or scene cuts (5). The mode correlation-based mode decision (MCMD) metric is defined and calculated using the 3×3 collocated MB window. Within this 3×3 neighboring MBs, each neighbor MB has an offline defined weight. This MCMD metric is used to classify the current MB motion activity in one of the classes described above. Independent of the motion activity, the SKIP mode is firstly evaluated and an early termination test is employed. If the SKIP prediction was not effective other modes are evaluated according to the motion class. The same classification described above is used here. For instance, A slow motion MB evaluates only the ME 16×16 .

The work proposed in Chan and Tang (2012) exploits the statistical behavior of the RDCost for the different coding modes along with the motion vectors difference in order to speed up the MVC encoding. In this solution, an interactive mode decision is employed. Based on statistical knowledge showing the ME is used more frequently than DE, the first interaction evaluates only the ME modes (all sizes). If the ME-based prediction is not satisfactory, a second interaction is used to evaluate the ME modes. However, only the block sizes that presented better coding performance for ME are evaluated for DE in the second interaction.

State-of-the-art schemes mainly achieve the complexity reduction via fast MD. However, they do not exploit the full space of neighborhood correlation in all spatial, temporal, and view domains. These schemes deploy fixed thresholding (Han and Lee 2008; Shen et al. 2009a, b) and, consequently, are unable to react to the changing QPs (i.e., changing bitrates). Moreover, in their worst case, state-of-the-art schemes—like (Han and Lee 2008; Shen et al. 2009a, b)—check all prediction modes, thus falling back to the exhaustive RDO-MD. As a result, these schemes provide limited complexity reduction.

In general, state-of-the-art schemes consider reference view encoded using the exhaustive RDO-MD and employ their fast MD scheme on the other views. These schemes prioritize the frames from the base view and the encoded quality of other views relies on the one encoded using exhaustive RDO-MD. This might lead to meaningful prediction error increase for the last views.

2.8.2 *Energy-Efficient Motion and Disparity Estimation*

To find a single optimal/good matching block the ME/DE performs several block-matching operations in multiple reference frames. Additionally, this search is replicated for multiple block sizes defined by the MVC standard. However, there are search directions (ME or DE), reference frames, and reference regions that are highly unlikely to provide a good matching. Also, there are suboptimal points that provide similar results at the cost of much reduced searching effort. See the example in Fig. 2.19a. A good matching for the diamond object is available in just one of the

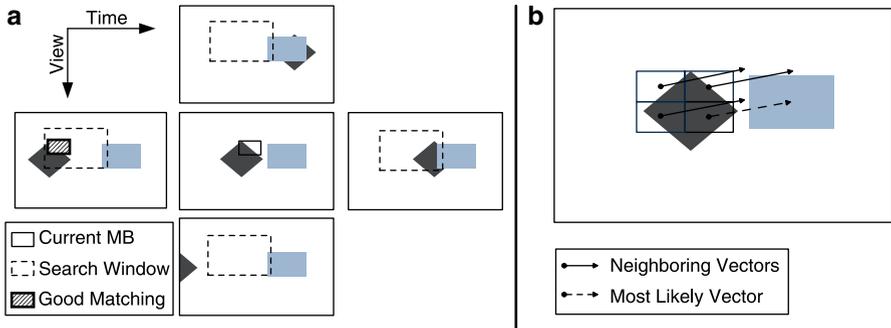


Fig. 2.19 ME/DE search conceptual example

four reference frames, the past temporal reference. In the future temporal reference the diamond is partially occluded by a second object (rectangle). In the disparity references the diamond is either occluded or out of the captured scene. In this scenario there is no need to perform searches in all reference frames, resulting in complexity/energy reduction. Other example is depicted in Fig. 2.19b. Note that the previously encoded neighboring MBs share a similar motion/disparity vector since they belong to the same object. Therefore, the current MB, which also belongs to the same object, is very likely to share a similar vector. This knowledge may be used to reduce the number of search operations by reducing the number of candidate blocks. A wide range of techniques to reduce the ME/DE complexity are available as presented in the following.

State-of-the-art fast ME/DE algorithms employ variable search range based on disparity maps (Xu and He 2008) taking into account the distinct behavior between ME and DE. The work presents an study on how the search window size impacts in the coding efficiency showing the importance of big search windows. However, disparity maps show that it is possible to reduce the effective search window by monitoring the disparity maps. From the disparity maps two parameters named vertical and horizontal scales (VS, HS) are defined. From the parameters the search window is reduced or increased in an asymmetric way, i.e., the search window may assume rectangular shapes. The increase and reduction are done in a factor 2.

In Kim et al. (2007) two strategies are used to predict motion and disparity vectors. One vector predictor used the traditional spatial median predictor from upper, left, and upper right neighboring MBs. The other predictor used the camera geometry and vectors from previously encoded frames to estimate the current vectors. The difference between the two predicted vectors is used to calculate the search window size. A small difference means accurate predictors and a small search window is required. Otherwise, for big differences a larger window is needed.

A fast direction prediction (ME or DE) based on the blocks motion intensity is proposed in Lin and Tang Angela (2009). It exploits the inter-view correlation to predict a search direction for reducing the ME/DE complexity. The base view is encoded using the ME and the frame regions are classified as slow motion if the SAD is smaller than a threshold. Similarly, the anchor frames of all views are

classified according to this strategy. For each MB to be encoded, the collocated MBs from base view and anchor frame are tested. In case both are slow motion, the current MB probably is also a slow motion MB and will be encoded using ME only. If only the base view collocated MB is not slow motion, DE is employed. Other cases require ME and DE processing.

The schemes in Han and Lee (2008), Ding et al. (2008a, b), and Deng et al. (2009) exploit the information from the base view and classify MBs into foreground and background regions. In Ding et al. (2008a, b) a fast ME based on complete DE is proposed. The DE is used to locate the correlated MB in the base view. After that, the coding information extracted from the base view is used to predict the motion vectors and partition sizes for the current MB.

The view-temporal correlation is exploited in Deng et al. (2009) by using the motion information of the first encoded view in order to reduce the complexity of further views. Additionally, disparity vectors from anchor frames are also taken into consideration. Using the geometric relation between the vector from base view and anchor frames, the authors predict the motion and disparity vectors that are used as search start point. A 2×2 refinement is applied around the predicted point. This process is repeated for each search direction.

The inter-view correlation is also evaluated in Shen et al. (2009a, b, 2010a, b) to reduce ME/DE search window. The so-called motion homogeneity is calculated using the collocated motion field from previous frames. If the MB presents a complex motion (homogeneity higher than a threshold) the complete search window is used for searching. Homogeneous motion MBs use a search window reduced in a factor of 4, $1/4$ of vertical size and $1/4$ of horizontal size. For the intermediate case, the search window is reduced in a factor of 2. Simultaneously, this solution employs a search direction selection. Homogeneous regions employ only ME search while complex motion regions employ both ME and DE. Moderate motion regions use the RDCost information to enable DE search.

Algorithm and architecture for disparity estimation with minicensus-adaptive support is proposed in Chang et al. (2010). A minicensus transform is applied over a pair of frames in two neighboring views to define a matching cost at pixel level. Weights are additionally generated using color distance. The cost and weights are aggregated to find the best disparity between the pair of frames. According to the authors, the two-pass strategy reduces the complexity if compared to a direct approach. The architecture proposed is discussed in Sect. 2.6.

The main drawback of these fast ME/DE algorithms resides in the fact that they do not exploit the full potential of the 3D-neighborhood correlation available in spatial, temporal, and disparity domains. Moreover, even the more sophisticated techniques are dependent on the complete first view encoding. However, it does not scale well for a large number of views as the prediction quality degrades in a hierarchical prediction structure. By encoding one view, the motion field information can be extracted but not the disparity field information (as no inter-view prediction is performed in this case). Therefore, it potentially limits the speedup of disparity estimation. Additionally, most of the techniques use fixed thresholding and thus perform inefficient under varying quantization parameters (QPs).

2.9 Video Quality on Energy-Efficient Multiview Video Coding

Techniques to reduce the complexity and energy consumption of the video encoder (such as fast mode decision and motion/disparity estimation) typically lead to video quality losses. To control the quality losses rate control methods may be employed through QP adaptation. Several rate control schemes are found in the current literature. Mostly they are developed targeting single-view encoders such as H.264. Recently, a few works specific to the MVC standard have been proposed focusing on frame- and BU-level RC. In this section we present an overview of the state of the art on rate control.

In the single-view domain the majority of recent proposals are extensions of the RC implemented in the H.264 reference software that employs a quadratic model for mean absolute differences (MAD) distortion prediction (Li et al. 2003). However, the quadratic model leads to limited control performance, as discussed in Tian et al. (2010). Aware of this limitation, the authors in Jiang et al. (2004) and Merritt and Vanam (2007) propose improved MAD prediction techniques. The scheme presented in Kwon et al. (2007) implements both distortion and rate prediction models while in Ma et al. (2005) the RC exploits rate–distortion optimization models. An RC based on a PID (proportional–integral–derivative control) feedback controller is presented in Zhou et al. (2011). In Wu and Su (2009), an RC scheme for encoding H.264 traffic surveillance videos using regions of interest (RoI) to highlight regions that contain significant information is proposed. In Agrafiotis et al. (2006), RoI is used to highlight preset regions of interests using priority levels. However, single-view approaches do not fully consider the correlation available in the spatial, temporal, and view domains and, consequently, cannot efficiently predict the bit allocation or distortion resulting in inefficient RC performance.

The early RC proposals targeting the MVC encoder are based on simple extension of single-view approaches (Li et al. 2003) and are still unable to fully exploit multiview properties. Novel solutions, however, have been proposed and most of them are limited to frame-level actuation. The solution in Yan et al. (2009b) uses an improved MAD prediction that differentiates the frame types. Intra frames, P and B frames with only temporal prediction, P and B frames with only disparity prediction, and B frames with both temporal and disparity prediction feature distinct MAD prediction equations. Once the MAD is predicted, the target bitrate is predicted for the GOP, refined to the GOP, and finally defined for each frame. An appropriate QP for each frame is defined based on the target bitrate. This work is extended in Yan et al. (2009a) by employing a technique to define the first QP in the GOP; it is used to encode the I frame. But these solutions are unable to properly handle the complex HBP structure of MVC limiting the number of input samples and the rate control learning.

The authors of Xu et al. (2011) define an pyramid-based priority structure extracted from the MVC HBP. The higher pyramid levels are used as reference to encode lower pyramid level, e.g., I and P frames belong to the highest level, B frames that refer to I and P frames belong to the second highest level, and so on. The higher levels are prioritized and are encoded using lower QPs (high quality) in

order to reduce error propagation. This solution, however, considers a fixed HBP structure and does not exploit the inter-GOP correlation.

To deal with distinct image regions within a frame there is a need for a BU-level RC. Moreover, in order to find a global optimal solution, a joint frame- and BU-level rate control scheme must be designed. Recent works have proposed solutions for the BU-level RC in MVC. In Park and Sim (2009) is presented a solution that deals with the frame-level and Macroblock (or BU)-level rate control. Firstly, the rate for each view is calculated based on weight parameters defined by the user. After that, the QP for each GOP is defined using the traditional H.264-based approach (Li et al. 2003) followed by a QP refinement for each frame. The frame-level QP definition considers the HBP coding structure to prioritize frames in higher hierarchical levels. A MAD-based strategy is used to calculate the target bitrate at MB level and a rate-distortion model (not described in the paper) is employed to define the QP for each MB.

The authors in Lee and Lai (2011) consider the HVS properties to propose a BU-level rate control solution that prioritizes the regions that are visually more important to the observer. For that, they define regions of interest using the Just-noticeable difference (Liu et al. 2010) metric along with luminance difference and edge information. Depending on the relation between these metrics, the QP is increased or decreased in relation to the initial QP (maximum QP in the neighborhood). However, this solution does not employ feedback-based control and just considers the coding information from one reference frame.

Generally, the available rate control techniques cannot fully exploit the correlation potential available in the spatial, temporal, and view domains of MVC. In addition, they are unable to adapt to multiple HBP structure and cannot employ the inter-GOP periodic behavior for RC optimization. Moreover, at the best of our knowledge, no work has proposed a Rate Control scheme for MVC able to jointly consider frame and BU level in a hierarchical and integrated fashion.

2.9.1 Control Techniques Background

In this section are presented the background concepts required to understand the rate control solution proposed in this monograph. Firstly, are presented the control theory basics behind the model predictive control (MPC) used for the frame-level RC. On the following, we present the statistical foundation supporting the Markov decision process (MDP) that is implemented in our BU-level RC. Finally, the concepts related to reinforcement learning (RL) are introduced.

2.9.1.1 Model Predictive Control

The control theory is a subfield of mathematics originated in engineering to deal with influences in the behavior of dynamic systems (Tatjewski 2010). Several control methods have been proposed ranging from very general to application-specific solutions to cope with a wide range of applications. Control problems specifications

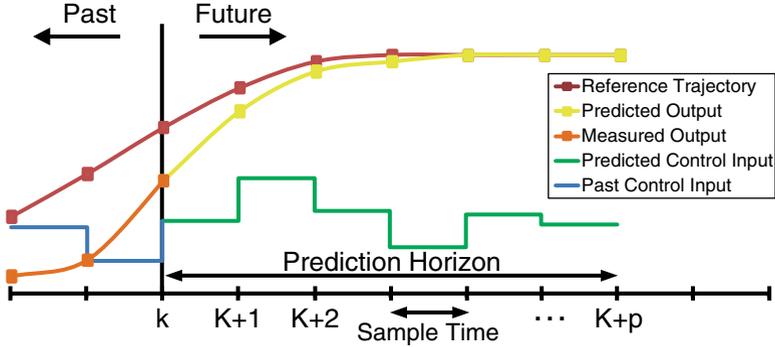


Fig. 2.20 Model predictive control (MPC) conceptual behavior

may significantly vary and the selected control method must ensure the stability of the given system. Thus, the selection of a control method for a given dynamic system may be very challenging. In case the controller does not fit the system it may compromise the stability of the entire system.

Among state-of-the-art control methods, the MPC has gained prominence by being able to accurately predict and actuate on a dynamic multivariable systems. It represents not a single control algorithm but a controller design scheme applicable to distinct systems including continuous or discrete in time, linear or nonlinear, integrated or distributed systems. MPC outperforms conventional feedback controllers (like PID) by keeping explicit integration of input and state constraints while considering state space constrains. Also, MPC can dynamically adapt to new contexts by employing rolling input and output horizons (see more details below).

The main goal of the MPC is to define the optimal sequence of actions to lead the system to a desired and safe state by considering the system feedback to previous states and previously taken actions (see conceptual MPC behavior in Fig. 2.20). To define this sequence of actions the MPC minimizes the performance function presented in Eq. (2.10). It minimizes the cost by defining a set of outputs y based upon a set of inputs u . Where $u[k+i-1|k]$, $i = \{1, \dots, m\}$ denotes the set of process inputs with respect to which the optimization is performed; u is known as the control horizon or input horizon in the MPC theory. Analogously, $y[k+1|k]$, $i = \{1, \dots, p\}$ is the set of outputs, named prediction horizon or output horizon (see Fig. 2.20). The control horizon determines the number of actions to find. The prediction horizon determines how far the behavior of the system is predicted. m and p are the size of control/input and prediction/output horizons, respectively. m is the number of measured outputs (history size) used for the optimization process, while p defines how many outputs are predicted; that is, how many future actions are considered in the optimization processes. k is the horizons index and represents the k th input/output horizon. y^{sp} defines the output set point that limits the prediction horizon:

$$\min_{u[k|k]..u[k+p-1|k]} \sum_{i=1}^p w_i (y[k+i|k] - y^{sp})^2 + \sum_{i=1}^m r_i \Delta u[k+i-1|k]^2. \quad (2.10)$$

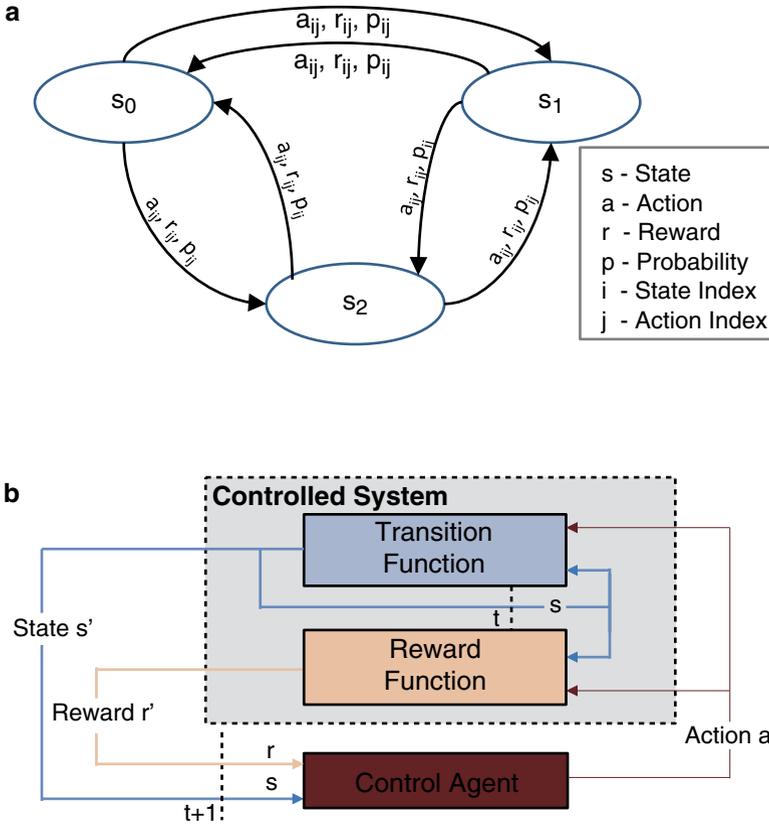


Fig. 2.21 Markov decision process (MDP)

2.9.1.2 Markov Decision Process

The MDP is a mathematical decision-maker framework for systems that outcome partly random and partly controlled by a decision maker (Arapostathis et al. 2003). MDP is a time discrete stochastic control process based on the extension of Markov Chains that adds the concepts of actions and rewards.

The symbolic representation of the MDP is a state machine or an automaton, as depicted in Fig. 2.21a, which evolves in response to the occurrence of events. It is formally defined by 4-tuples $(S, A, P(.,.), R(.,.))$ composed of a finite set of states $S = \{s_0, s_1, \dots\}$, actions $A = \{a_0, a_1, \dots\}$, rewards $R = \{r_0, r_1, \dots\}$, and transition probabilities $P = \{p_0, p_1, \dots\}$. The S includes all possible states assumed by the controlled system; actions A are the possible acts to be taken by the decision maker in face of a given system state. $P(S)$ is the probability distribution of transitions between system states and, finally, $R(S)$ is the reward related to a given action for a given state. At each discrete time step t the process lays in a state $s \in S$ and the decision maker may choose any action $a \in A$ that will lead the process to a new state $s' \in S$ providing a

shared reward $R_{at}(s, s')$, as shown in Fig. 2.21b. The rewards are used by the decision maker in order to find an action that maximizes, for a given policy, the total accumulated reward, as shown in Eq. (2.11) (where $0 \leq \gamma \leq 1$ denotes the discount factor):

$$\sum_{t=0}^{\infty} \gamma^t R_{at}(s_t, s_{t+1}). \quad (2.11)$$

By definition, the Markov process is considered a controlled Markov process if the transition probabilities $P(S)$ can be affected by an action. Equation (2.12) defines the probability P_a that an action a in the state s at time t will lead to state s' at time $t+1$:

$$P_a = R_{at}(s_{t+1} = s' \mid s = s_t, a_t = a). \quad (2.12)$$

Multiple extensions have been proposed to the MDP in order to best fit to distinct problem classes. For systems where the transitions probabilities or the rewards are unknown a priori, the Reinforcement Learning method may be applied to solve the MDP, as detailed in the following section.

2.9.1.3 Reinforcement Learning

Reinforcement learning model is an agent to improve autonomous systems performance through trial and error by learning from previous experiences instead from specialists, that is, the agent learns from the consequences of actions. In reinforcement learning model the agent is linked to the system to observe its behavior and take actions. RL theory is based on the Law of Effect, that is, if an action leads to a satisfactory state the tendency to produce this action increases. For each discrete time step t the RL agent receives the system state $s \in S$ and rewards $R(S)$ to take an action $a \in A$ that maximizes the reward $R_{at}(s, s')$. This action may lead the system to a new state $s' \in S$ and produce a system output, in terms of a scalar reinforcement value, used to define the new reward $R_{a(t+1)}(s, s')$ according to Eq. (2.13). The general representation of reinforcement learning value given by RL in Eq. (2.14), where U denotes the function that changes the system state from s to s' and h_R denotes the history of reinforcement learning:

$$RL_{a(t+1)}(s, s') = RL_{at}(s, s') + RL, \quad (2.13)$$

$$RL = U(s, s') + h_R. \quad (2.14)$$

2.9.1.4 Region of Interest

Within a video frame there may exist multiple regions or objects with distinct image properties and distinct importance for the observer. The image regions that are considered, for some reason, more important are called Regions of Interest. In this monograph, we consider all regions of semantically equal importance



Fig. 2.22 Variance-based region of interest map (*Flamenco2*)

leaving space for application-specific optimizations such as for 3D-surveillance, 3D-telemedicine, etc. However, at the encoding perspective, textured regions tend to have different coding properties at the mode decision and bit allocation perspectives if compared to homogeneous regions. To classify the image regions we use the variance map (Fig. 2.22) to characterize the texture complexity. Variance depicts the degree of dissipation of a given population [see definition in Eq. (4.1)]. In this case, how the pixel values of an image region are distributed. High variance define textured regions (represented by brighter points in Fig. 2.22) while low variance define homogeneous regions (dark regions in Fig. 2.22).

2.10 Summary of Background and Related Works

The MVC is the most efficient video coding standard focusing on 3D-video coding. It is able to provide 20–50 % of coding efficiency increase, if compared to H.264 simulcast, by employing inter-view prediction, the disparity estimation. Mode decision and motion and disparity estimation represent the most complex modules in the MVC encoder and bring big challenges for their real-world implementation.

The implementation of MVC encoders may exploit different multimedia processing architectural solutions. Currently, the most preeminent alternatives are multimedia processors/DSPs, reconfigurable processors, ASICs, and heterogeneous multicore SoCs. Each solution presents positive and negative points. On the one hand, ASICs provide the highest performance and energy efficiency at the cost of no flexibility. On the other hand, multimedia processors/DSPs are totally flexible but deliver low performance and reduced energy efficiency. Heterogeneous multicore and reconfigurable processors provide trade-off points between ASICs and processors. By employing units specialized in each kind of task, the heterogeneous multicore SoCs improve the performance in relation to multimedia processors but typically present issues related to programming and portability. Reconfigurable processors can cover this gap by employing extensible instruction set and defining, at

run time, if regular or custom instructions should be used in that specific time instant. Still, these solutions are unable to meet the performance and energy efficiency required for MVC encoding without application-specific ASIC acceleration. Therefore, considering the current technology, a complete ASIC encoder or heterogeneous SoCs with hardware specific accelerators are seen to be the most feasible solutions for embedded mobile devices.

Multiple proposals targeting on complexity and energy reduction for the MVC are available in the current literature. These contributions are centered in two abstraction levels: the algorithmic and architectural levels. At the coding algorithms perspective, complexity reduction is most frequently addressed at the mode decision and motion and disparity estimation because they represent the most complex MVC blocks. The mode decision solutions used distinct side information in order to reduce the number of coding modes tested during the coding process. Video properties such as texture, edges, luminance, and motion/disparity activity are used to predict the most probable coding modes in each image regions. Additionally, extensive analysis has been done to learn how neighboring views and frames are correlated. This correlation is also useful to predict the coding modes. As the ME/DE spends about 90 % of the total encoding time, the same kind of information is used to predict the most probable motion and disparity vectors and reduce the ME/DE complexity. However, the related works do not fully exploit the correlation available within the 3D-neighborhood and perform badly under content changing scenarios. Moreover, these solutions are not developed considering the energy perspective and cannot react to battery-level changing situations by dynamically adapting the complexity to the available energy.

Generally, the complexity reduction techniques lead to uncontrolled quality degradation and coding efficiency losses. The rate control becomes a key task in order to minimize this complexity reduction drawback. The majority of rate control solutions currently available target the H.264 or are simple extensions from H.264 solutions. The few rate control algorithms designed for MVC focus only on frame-level or basic unit-level actuation levels. Additionally, these algorithms do not use the intra- and inter-GOP bitrate correlation in the 3D-neighborhood.

At the hardware architectural perspective ME/DE is the most studied MVC coding block. The ME/DE is a processing and memory-intensive task requiring massively parallel processing and efficient memory access and management. The resulting high-energy consumption is mainly related to external memory access and on-chip video memory size. Diverse related works propose ME/DE processing hardware architectures, memory hierarchies, and data reuse techniques. However, they share limitations related to the complexity reduction algorithms implemented (leading to quality losses), excessive external memory accesses, or large on-chip memory resulting in high energy. Moreover, most of the available architectures lack the ability to dynamically adapt its operation according to changing coding parameters or video content characteristics.

Therefore, there is a demand for novel and energy-efficient MVC encoding solutions able to significantly reduce energy consumption under changing video and system scenarios. For this reason, this monograph targets on jointly addressing the energy issues at algorithmic and architecture levels while sustaining the video quality.

Chapter 3

Multiview Video Coding Analysis for Energy and Quality

The Multiview Video Coding (MVC) standard brings high coding efficiency gains reducing the bitrate in 20–50 % for similar video quality if compared to the H.264/AVC simulcast. The coding efficiency gains are driven by novel high-complexity coding tools that drastically increase the overall encoding processing effort and, consequently, the energy consumption. In this section an extensive analysis of the energy requirements for real-time MVC encoding and the energy consumption breakdown are presented. The goal is to provide a better comprehension on the MVC performance and energy requirements. Additionally, the requirements in terms of objective video quality are discussed in the following.

3.1 Energy Requirements for Multiview Video Coding

Encoding MVC at high definitions has shown to be an unfeasible task for mobile devices when all coding tools are implemented without energy-oriented optimizations. State-of-the-art embedded devices are unable to provide the processing performance or to supply the energy required by the MVC encoder. To demonstrate the energy-related challenges to MVC encoding a case study is presented in the following.

Figure 3.1 presents the energy consumption to encode a 4-view HD1080p video sequence using the MVC encoder while considering four fabrication technologies. Also, the battery draining time for the following state-of-the-art smartphone batteries are presented: Apple iPhone4 (5.25 Wh), Nokia N97 (5.6 Wh), and Samsung Galaxy S3 (7.8 Wh). Note, these smartphones are unable to attend the MVC constraint. Despite the technological scaling that provides meaningful energy reduction for deep-submicron technologies, the energy consumption remains high considering embedded devices constraints. For instance, let us analyze the best-case scenario where a device is fabricated with a state-of-the-art 22 nm fabrication node and features a 7.8 Wh battery as available in the latest Samsung Galaxy S3 (Samsung 2012)

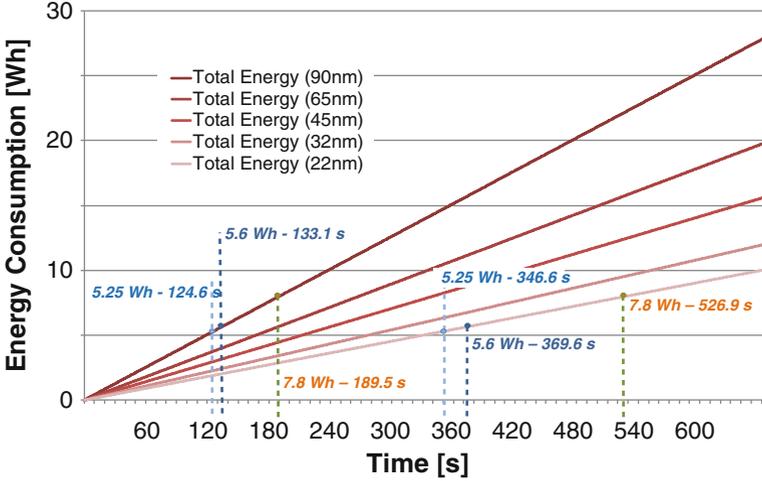


Fig. 3.1 MVC energy consumption and battery life

released in Q3 2012. For a scenario where MVC encoder is the only task draining the battery, only 526.9 s (8 min, 46.9 s) of recording would be possible before the battery was completely drained. The presented battery life is not acceptable and does not attend market and user requirements. Meaningful energy reduction techniques are required to bring the MVC consumption to a feasible energy envelope. For this, a better understanding on the energy consumption sources is required.

Figure 3.2 demonstrates the motion and disparity estimation task (ME/DE) is responsible for about 90 % of the total energy. These numbers were measured using the Orinoco (Krolikoski 2004) simulation environment and might present discrepancies in the actual numbers. This simulation, however, is worth for a relative comparison of the energy consumed by MVC encoding tools. This numbers consider the fast search algorithm TZ Search (Tang et al. 2010) as search pattern. Motion compensation (2.5 %), deblocking filter (2.5 %), and intra-frame prediction encoder (2 %) are the following in terms of energy consumption while representing less than 2.5 % each. Thus, reducing ME/DE consumption is of key importance to reach energy efficiency. ME/DE consumption is directly related to the size of search window (SW), that is, the size of the region to perform the search. The increase in ME/DE search window leads to energy increase due to increased number of matching candidates and larger amount of data required (memory accesses) to perform the task. Figure 3.3 quantifies the energy consumption for five distinct sized SWs. Comparing the corner cases, a small search window $[\pm 16, \pm 16]$ to a big $[\pm 128, \pm 128]$ SW, the energy increases in a factor of 6.5 \times . From single-view knowledge it is possible to affirm that there is no need for using SWs larger than $[\pm 64, \pm 64]$. However, disparity vectors tend to have larger magnitude and the ME/DE task requires increased SW to find these matching candidates. According to Xu and He (2008), for a good disparity estimation performance in HD1080p video sequences, the search window should be at least $[\pm 96, \pm 96]$.

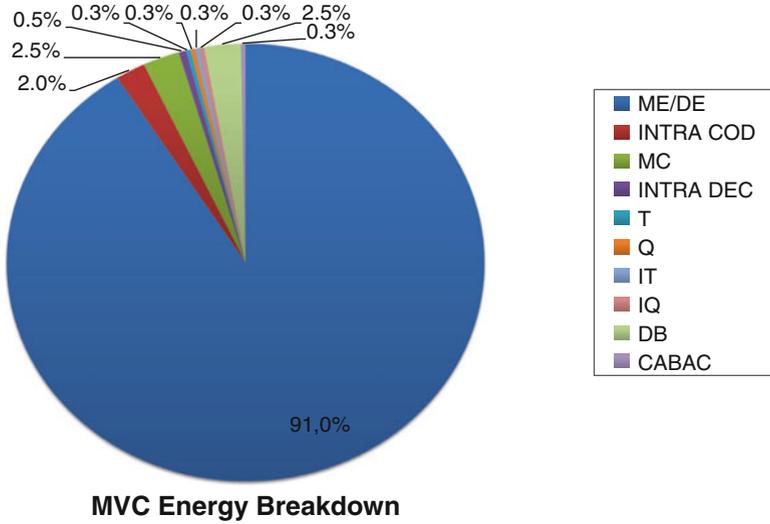


Fig. 3.2 MVC component blocks energy breakdown

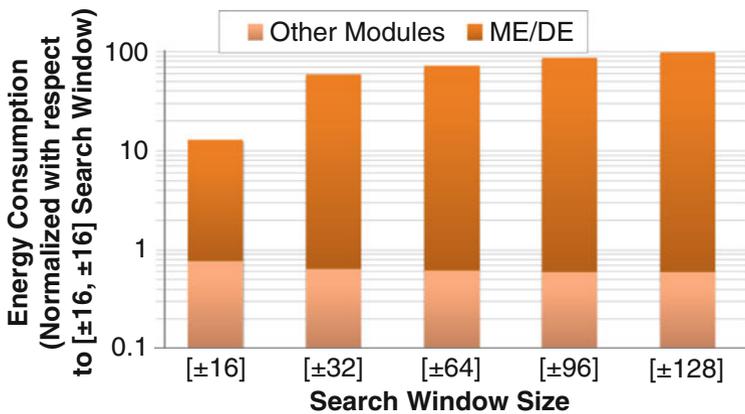


Fig. 3.3 MVC energy breakdown for multiple search window sizes

Although the analysis correctly depicts some sources of energy consumption, a deeper knowledge of the application behavior is mandatory. The MVC encoder hides in the encoding process a control function that controls the complexity of each and single module discriminated in Fig. 3.2. The mode decision (MD) defines how many modes are tested and how many times the ME/DE search is performed, how frequently the intra-frame encoder is used, etc. The relation between the exhaustive mode decision - the Rate-Distortion Optimized MD (RDO-MD) - to the simplest possible MD that tests a single coding mode is in the order of 100x energy consumption, as shown in Fig. 3.4 that tests a single coding mode is in the order of

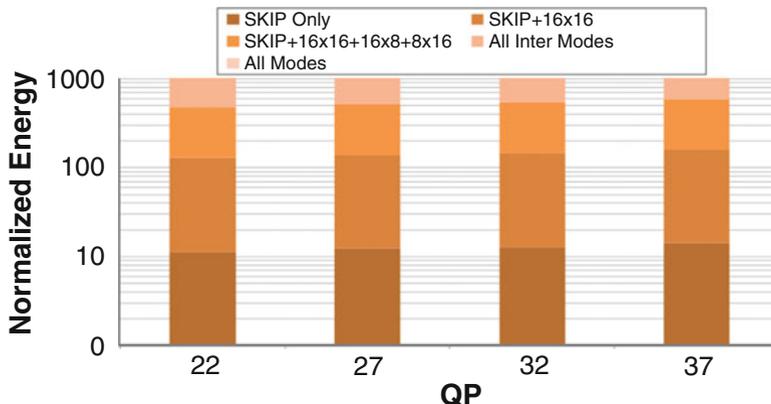


Fig. 3.4 MVC energy for distinct mode decision schemes

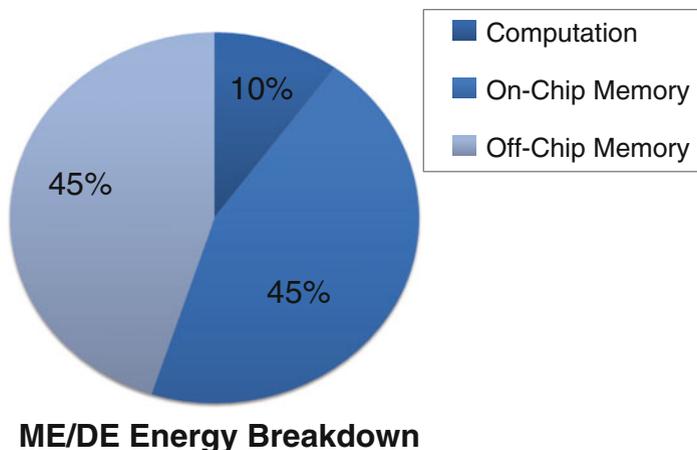


Fig. 3.5 ME/DE energy breakdown

100x energy consumption, as shown in Fig. 3.4. Obviously, the single mode MD is not used in practice under penalty of poor quality and coding efficiency results. Nevertheless, this example highlights the optimization space for energy-efficient solutions in the MD control.

At the architectural perspective, computation and memory (external memory access and on-chip memory storage) are the two energy consumption sources. Here, dynamic and static energies are jointly considered. As shown in Fig. 3.5, the energy breakdown is composed of 90 % memory-related energy consumption while 10 % are represented by the computation itself. Typically, for a rectangular search window on-chip memory using Level-C data reuse (Chen et al. 2006), the on-chip memory energy and external memory access are evenly distributed but may vary according to design options (on-chip memory size, data-reuse scheme, etc.).

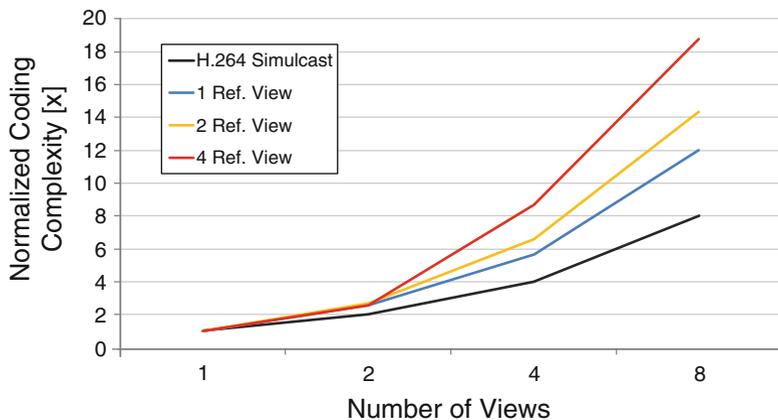


Fig. 3.6 MVC vs. Simulcast complexity

The presented energy breakdown highlights the importance of reducing memory-related energy. Even so, the reduction of complexity stands as key challenge for low-energy MVC. Observe that multimedia processing applications are typically data-oriented applications and require intense memory communication. However, the complexity reduction leads to a win-win situation where both less data is processed and less memory accesses are required. Thus, complexity reduction positively impacts computation and memory energy consumptions.

The following sections discuss on how computational effort and memory access influence the overall energy consumption in MVC and how these components are distributed among the MVC modules.

3.1.1 MVC Computational Effort

The MVC high-energy consumption is driven by the computational effort associated with the MVC video encoder. In this section we compare the complexity in relation to previous standards and quantify the main sources of complexity within the video encoder. The experiments here presented consider the fast search algorithm TZ search for motion and disparity search.

Figure 3.6 compares the MVC encoder in three distinct scenarios compared to the H.264-based simulcast encoding. The 8-view video sequences were encoded independently (simulcast) and using inter-view prediction with one, two, or four reference views (when available). Custom extensions to the JMVC (JVT 2009a) reference software were done to support more than two reference views. Although using more than two reference views is not a common practice in current encoding systems, the increase in reference views is expected for many-view systems, especially for 2D-array camera arrangements where the four surrounding neighbor views are closely correlated to the current view. The measured complexity to encode

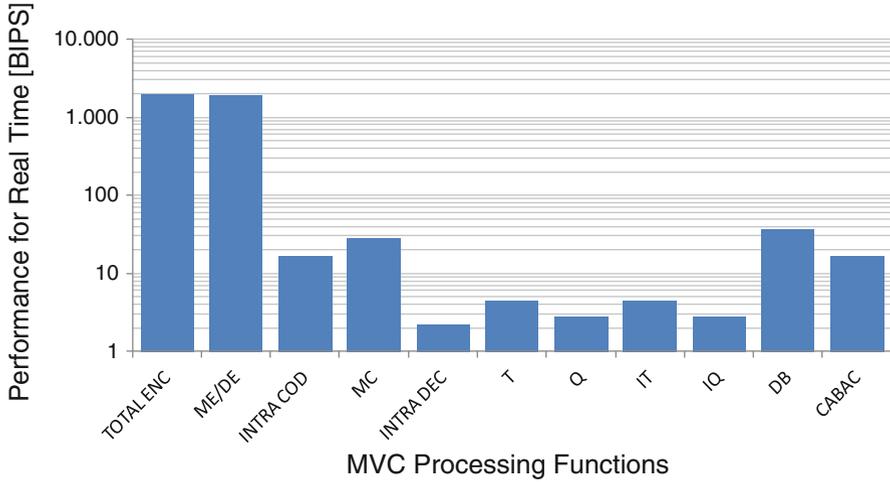


Fig. 3.7 MVC computation breakdown

eight views using four reference views is $19\times$ more complex than encoding a single H.264 view. Even using two reference views, as current multiview systems, the complexity exceeds in $14\times$ the H.264 single-view complexity. To understand what this complexity represents, it is important to consider that real-time H.264 encoding for HD1080p still poses interesting challenges in the embedded devices development and require application specific hardware acceleration (see discussion in Sect. 2.5.4). Moreover, according to Ostermann et al. (2004) the H.264 encoder is about $10\times$ more complex than the MPEG-4 Part 2 encoder. If compared to the simulcast encoding of eight views ($8\times$ compared to H.264 single view), the MVC is $1.75\times$ and $2.37\times$ more complex for two and four reference views, respectively.

The total encoder complexity is mainly concentrated in the motion and disparity estimation (ME/DE) unit responsible for about 90 % of the total processing, as depicted in Fig. 3.7. The deblocking filter (DF) and motion compensation (MC) blocks are the more complex blocks after ME/DE. The MVC encoder complexity measured from the JMVC (JVT 2009a) reference software without optimizations leads to 2 GIPS (Giga Instructions per Second) for only 4-view real-time MVC encoding at HD1080p resolution. This throughput is unfeasible even for high-end desktop computers. For instance, the latest Intel Core i7 3960X (Bennett 2011) processor with six physical cores running at 3.3 GHz is able to provide about 180 MIPS. Thus, the state-of-the-art high-end processors are orders of magnitude below the performance requirements for real-time MVC encoding if no application/architectural optimizations are performed. The task is even more challenging for embedded processors.

For energy-efficient MVC there is a need to drastically reduce the complexity. Based on the presented observations, ME/DE and MD modules have the highest potential for complexity reduction and, for this reason, are explored in this work. Therefore, deep application knowledge is required to design efficient complexity reduction algorithms able to avoid objective and subjective video quality losses.

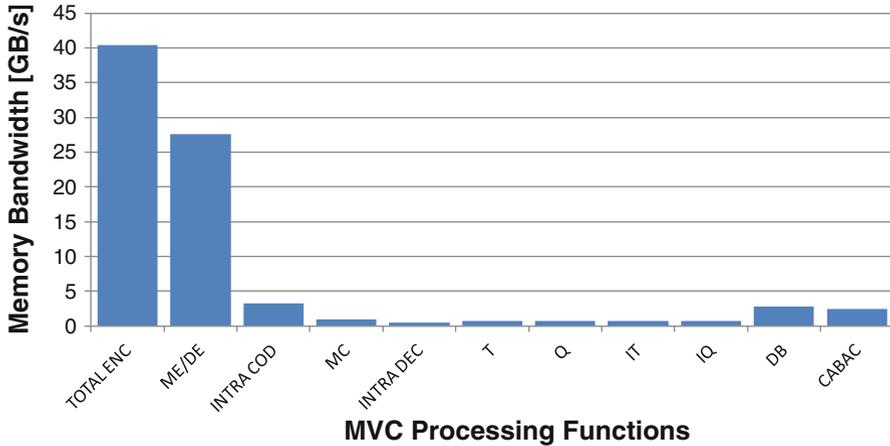


Fig. 3.8 Memory bandwidth for 4-views MVC encoding

3.1.2 MVC Memory Access

The other major component of energy consumption is related to data access. Multimedia applications are known for their data-oriented nature and consequent intense memory communication. The MVC encoder for 4-views HD1080p requires a memory bandwidth of 41 GB/s, as pointed in Fig. 3.8. It can be met using GDDR5 memory interfaces available in high-end GPUs such as Nvidia GeForce GTX 690 (384 GB/s @ 300 W) (Nvidia 2012a) at the cost of high energy consumption. For embedded systems, however, the memories interfaces are limited by power constraints and deliver a reduced bandwidth. Thus, this bandwidth is not feasible for embedded devices. For instance, the Nvidia ULP GeForce embedded in Tegra 3 SoC (Nvidia 2012b) provides a theoretical limit of 4.26 GB/s employing a LPDDR2-1066 memory interface. In this scenario, MVC encoding in embedded devices pose the need for drastically reducing the memory bandwidth through algorithmic and architectural optimizations.

In video encoding systems, mainly the MVC video encoder, the access to the decoded picture buffer (DPB) is the memory bottleneck. The DPB stores all reference frames used for inter-frame and inter-view (ME/DE) prediction. The frames are written in the DPB after the DF processing and the ME/DE block reads the stored data to perform motion/disparity search. ME/DE unit is responsible for about 68 % of the encoder total memory access requiring a 28 GB/s memory bandwidth for 4-view encoding, as shown in Fig. 3.8. The measured memory bandwidth is far higher compared to the raw video data input (355 MB/s) because the reference frame data may be requested multiple times in order to perform the motion/disparity search for distinct MBs. Aware of this behavior, multiple techniques try to reduce external memory accesses through employing on-chip video memories and data-reuse techniques. Even though effective at the external memory perspective, these solutions

significantly increase the on-chip energy consumption. Therefore, energy-efficient external and on-chip memory reduction must be jointly considered at design time and at run time. Moreover, the complexity reduction design, discussed in Sect. 3.1.1, must consider the memory access behavior to optimize the overall energy consumption.

3.1.3 *Adaptivity in MVC Video Encoder*

The high complexity and memory requirements posed by the MVC encoder are not the only challenges related to its realization. MVC energy consumption is unevenly distributed along the time. Processing and memory energy components vary depending upon coding parameters, user's definitions, system state, and video content. These run-time variations make the MVC encoder design even more challenging. If on the one hand, an under-dimensioned encoder leads to performance issues and does not guarantee reduced energy consumption due the need of additional buffering. On the other hand, over-dimensioned encoders face underutilization and unnecessary energy consumption.

The MVC prediction structure is a dominant factor in terms of energy variation once distinct frame types (I, P or B) present distinct processing and memory access behaviors. I frames are the lightest frames once the ME/DE (that represents 90 % of the encoder complexity) is completely skipped. P frames employ ME/DE to a single direction. In this scenario the P frames search in a single reference frame. The B frames require heavy processing, in comparison to I and P frames, and intense memory access while executing ME/DE search in multiple reference frames/views. In Fig. 3.9 the frame-level energy consumption for seven GOPs is presented. Each bar represents the sum of energies spent to encode the frames from all four views that belong to the same time instant (i.e., $S0Tx + S1Tx + S2Tx + S3Tx$). GOP borders (anchors), for the experimented prediction structure, have one I frame, two P frames, and one B frame (that performs only DE once there are no temporal references available), for more details see prediction structure in Fig. 2.7. Consequently, GOP borders drain reduced energy amount (1.5 Ws/frame), as shown in Fig. 3.9. All other relative positions within the GOP are composed only by B frames and the energy consumption drastically increases in comparison to GOP borders. Typically, the center of GOP is the energy hungriest time instant once temporal references are far and more extensive motion search is required to find a good matching. According to the experiment presented in Fig. 3.9, the energy consumption may exceed 7 Ws/time instant in this case. It represents a $4.7\times$ instantaneous energy variation within the same GOP.

Although prediction structure-related energy variations may be easily inferred from the coding parameters, there is another important variation source that may not be easily obtained, the video content-related variations. The video content variations occur at multiple levels (a) view level: distinct views may present distinct video content such as textures, motion and disparity behavior; (b) frame level: video

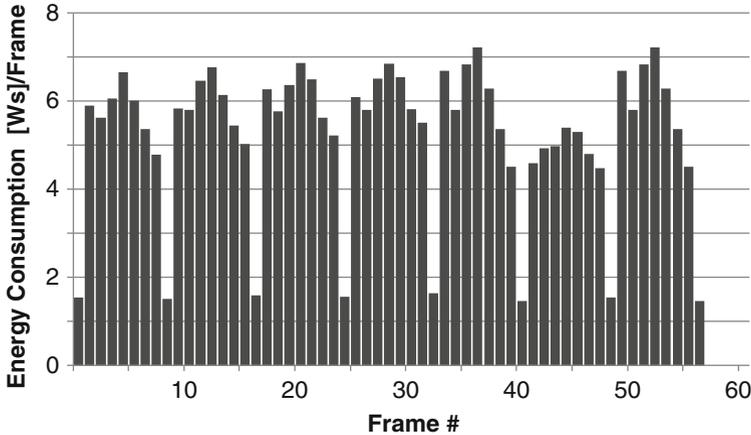


Fig. 3.9 Frame-level energy consumption for MVC

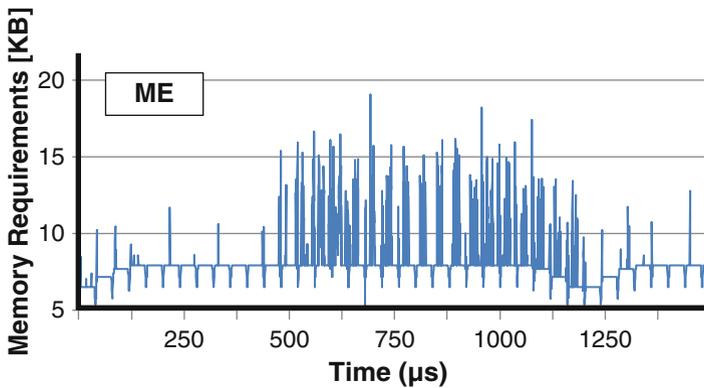


Fig. 3.10 Memory requirements for motion estimation at MB level

properties vary along the time; (c) MB level: within the frame distinct regions or object may present distinct image properties. Figure 3.9 depicts the energy variations along the time. For instance, GOP #6 (frames 41–49) drains reduced energy amount if compared to previous GOPs due to reduced coding effort resulting from easier-to-encode video content. The MB-level variations are shown in Fig. 3.10 in terms of ME memory requirements. The memory usage changes along the time depending on the video content motion intensity. High motion regions present increased memory usage in relation to low motion regions within the same frame.

Therefore, energy efficiency in MVC encoding requires the understanding of the energy sources variations and the design of adaptive architectures able to manage, at run time, the energy consumption while considering dynamically varying parameters (such as video content) and system state.

3.2 Energy-Related Challenges in Multiview Video Coding

The large complexity and intense memory communication related to MVC pose a series of challenges related to real-time encoding for high definitions mainly at the embedded systems domain. Energy consumption represents the most challenging issue related to embedded MVC encoding. Thus, there is a dire need for energy reduction of the MVC video encoder through complexity and memory access reduction. Energy-efficient solutions must jointly consider optimizations at algorithmic and architectural levels. Coupling deep application knowledge to intelligent employment of low-power design techniques is a key enabler for energy-efficient embedded MVC encoder realization.

Based on the discussion presented along Sect. 3.1, the energy-efficient MVC requires the following optimizations at algorithmic level:

- *Energy-efficient mode decision scheme*: The MVC defines an increased optimization space for the optimal prediction mode selection leading to high complexity and energy requirements, as demonstrated in Sect. 3.1. An efficient fast mode decision scheme is needed to reduce the optimization space through heuristics able to accurately anticipate the coding mode selection. The neighborhood information and image/video properties may provide hints to completely avoid the evaluation of unlikely prediction modes
- *Energy-efficient motion and disparity estimation*: ME/DE is the most complex and energy hungry module in the entire MVC encoder. Intelligent optimizations in ME/DE lead to meaningful overall energy reduction. Energy-efficient ME/DE may be reached by applying ME or DE elimination, search direction elimination, motion/disparity vector anticipation, object motion/disparity field analysis, etc.
- *Dynamic complexity adaptation*: The energy-efficient MD and ME/DE can be designed considering distinct strengths in order to handle the energy versus quality trade-off. Additionally, the MVC presents a dynamically varying behavior along the time depending on coding parameters, user's constraints, and video content. An energy-aware complexity adaptation scheme must be able to predict these variations and to react at run-time through reduction/increase of complexity budget by setting MD and ME/DE parameters. The dynamic complexity adaptation scheme may also exploit asymmetric coding properties such as the binocular suppression theory (Stelmach and Tam 1999).

The energy-efficient algorithms described above must be designed considering their impact in the architectural implementation. At architectural level, the energy-efficient solution must employ:

- *Low-energy motion and disparity estimation architecture*: The ME/DE task requires high throughput but typically allows a high level of parallelism. To attend the throughput requirements at a reasonable frequency of operation while reducing energy multiple levels of parallelism must be exploited including (a) pixel-level, (b) MB-level, (c) reference frame-level, (d) frame-level, and (e) view-level parallelisms. It allows operating in a reasonable range of operation

frequency and voltage. The processing units should be designed to enable power gating and/or DVS to adapt to the performance variations.

- *Energy-efficient on-chip video memory hierarchy:* Simply feeding the highly parallel ME/DE processing units while avoiding performance losses is typically a very challenging task. The on-chip video memory, however, has to deal with the high memory-related energy consumption and memory requirements variations. For that, an accurate memory sizing strategy is required. Also, the on-chip video memory must support partial power gating and/or DVS to adapt to memory requirement variations while minimizing static energy consumption.
- *Data-reuse and prefetching technique:* Neighboring MBs tend to access repeated times the same data from reference frames during the ME/DE process. To avoid additional external memory access the reference data must be stored in the local memory. However, increased local memory leads to increased static energy. Hence, only the actually required data must be read from external memory and stored locally. The energy-efficient MVC solution requires a memory-friendly data-reuse technique able to reduce external memory access without employing increased local memory. To avoid performance losses due to local memory misses the required data must be prefetched accordingly. Thus, it demands an accurate memory behavior predictor that understands the ME/DE search pattern. Accurate prefetching becomes even more challenging for state-of-the-art adaptive and customizable search algorithms.
- *Dynamic Power Management:* Supporting power gating and/or DVS in memory and processing units does not directly lead to energy savings. To reach energy efficiency an intelligent dynamic power management scheme is required to define the proper power states at each given time instant. The DPM must apply deep application knowledge including offline statistical analysis, neighborhood history, and image/video characteristics in order to accurately predict performance and memory requirements and take proper action.

Addressing each challenge related to energy-efficient MVC brings a contribution to the overall energy reduction. A balanced combination of energy-efficient techniques may lead to drastic MVC energy reduction. The energy reduction, however, shall not be built upon meaningful coding efficiency/video quality losses. Otherwise, the use of JMVC over simulcast is no more justified. Video quality issues are discussed in details in the following section.

3.3 Objective Quality Analysis for Multiview Video Coding

In the previous subsections the need for energy-efficient MVC encoding was motivated and justified. To reach such efficiency, complexity reduction, efficient architecture, and efficient memory management techniques including run-time adaptations are required. These techniques, however, may lead to undesirable rate-distortion (RD) performance losses. In other words, the optimizations techniques

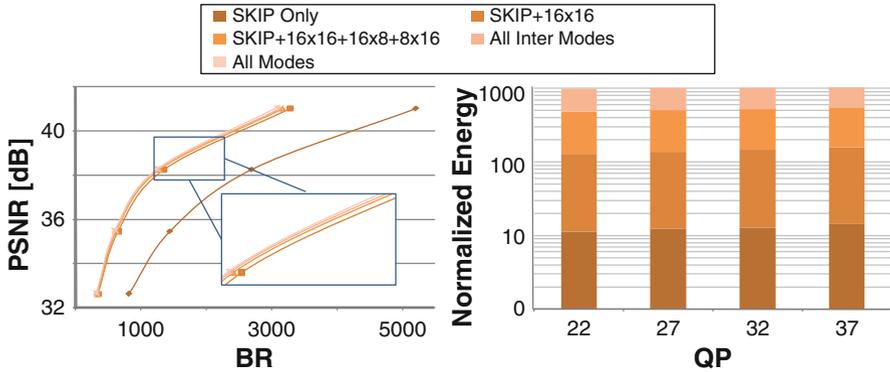


Fig. 3.11 Objective video quality in relation to coding modes

may lead to reduced video quality for the same output bitrate. For simplicity, in this section we discuss the impact of optimizations algorithms in terms of video quality variation. However, it is necessary to keep in mind that, for a more general analysis, the rate-distortion performance must be evaluated by jointly considering the objective video quality and the generated bitrate. The RD trade-off can be managed through Quantization Parameter (QP) adaptation by employing an efficient rate control (RC) scheme (see Sect. 2.3.4).

To enable the use of MVC in real-world solutions, its implementation must be energetically feasible and the resulting video quality (for similar bitrate) must be significantly improved in relation to previous coding standards applying simulcast-based coding. According to this assumption, the energy reduction techniques must aggressively reduce the total energy consumption at the cost of none or reduced quality loss. Figure 3.11 (extended from Fig. 3.4) depicts the impact of some simplified mode decisions in terms of video quality versus bitrate. Take the example of the “SKIP only” MD, which represents 1 % of the total coding energy compared to the exhaustive RDO-MD at the cost of nearly 3 dB quality loss. Remarkably, this is not a reasonable solution due to high-quality loss. According to the experiments presented in Merkle et al. (2009) and Oh et al. (2011), the MVC provides about 1 dB quality increase in relation to H.264 simulcast. In case the energy-efficient optimizations lead to a quality drop at the order of 1 dB there is no reason for using the MVC. In this scenario multiple state-of-the-art H.264 encoders should be employed avoiding the $1.75\times$ – $2.37\times$ complexity increase (Sect. 3.1.1) driven by MVC in relation to simulcast. Intermediary solutions are also presented in Fig. 3.11 dealing with the relation between energy and video quality. The same kind of energy vs. quality observations is noticed in the ME/DE optimizations.

Additionally, the 3D video quality includes additional properties in relation to the regular 2D videos. Blocking artifacts are severely undesirable in 3D videos and must be avoided during the encoding process. Such artifacts may lead to problems for intermediate viewpoints generation and/or to the stereo pair mismatch problem, as described in (Stelmach and Tam 1998). Quality drop due to blurring effect

in certain views, however, is tolerable and is attenuated according to the binocular suppression theory which is based on the psycho-visual studies of stereoscopic vision (Stelmach and Tam 1998). According to it, if the video qualities of left and right eye views differ, the overall perceived quality is close to the high quality of the sharper view. In other words, there is space for controlled quality losses in odd or even views while sustaining the perceived quality and reducing overall energy consumption.

3.4 Quality-Related Challenges in Multiview Video Coding

To reduce the possible quality losses inserted by the energy-efficient optimizations related to the challenges pointed in Sect. 3.2, there is a need to define quality protection mechanisms able to manage the energy versus quality trade-off. Such mechanisms must consider the application dynamic behavior in order to optimize the video quality for a given energy constraint. To sustain the overall video quality the energy-efficient MVC must employ:

- *QP-based thresholding*: Most of the energy reduction schemes proposed in the current literature are unable to react to changing QP scenarios due to fixed thresholding. This limitation leads, for corner case scenarios (low or high QPs), to very high-quality losses or to limited energy reduction. To deliver high video quality while providing meaningful energy reduction an energy-efficient MVC must control the energy reduction schemes through QP-based threshold equations. Moreover, the thresholds must be defined based on extensive statistical analysis to avoid biasing.
- *Frame-level rate control*: Some energy optimizations may prioritize key frames by providing higher energy/processing budgets for such frames. The drawback of these approaches is the uneven quality distribution, at frame level, inside the same view or between neighboring views. If these quality variations are not properly controlled they may lead the observer to experience some discomfort (Stelmach and Tam 1998). In order to avoid such quality variations, a frame-level rate control unit must be implemented. The RC task is to predict and control the bitrate versus quality trade-off and to distribute the amount of bits available (according to a given bandwidth limitation) in such a way to reduce the video quality oscillation and maximize the overall perceived quality.
- *Basic unit-level rate control*: A rate control is also required at basic unit level once energy-efficient optimizations are also defined at MB level. In this scenario, the basic unit RC must be designed to optimize the overall video quality within each frame while considering image/video properties of the image regions.

The challenges described above are critical to deliver high video quality even under a series of energy-restrictive constraints and simplifications along the video coding process. In the following section is presented an overview on this monograph contribution. It describes, at high level, the main energy-efficient algorithms and architectures proposed in this volume along with the video quality control strategies.

3.5 Overview of Proposed Energy-Efficient Algorithms and Architectures for Multiview Video Coding

Figure 3.12 presents the overview of this monograph contribution related to the energy-efficient realization of MVC. The high-level diagram presents the algorithmic and the architectural contributions along with the conceptual contribution related to the 3D-Neighborhood correlation. Each contribution is detailed in the Chaps. 4 and 5, as pointed in Fig. 3.12.

The energy reduction and management algorithms, hardware architecture design, memory designs, and data-reuse schemes are based on the application knowledge to deliver more efficient results. In this monograph we define the 3D-Neighborhood concept that is widely used to guide the algorithmic and architectural contributions of this monograph. The 3D-Neighborhood is defined as the MBs belonging to neighboring regions at spatial, temporal, and view/disparity domains. The analysis of the 3D-Neighborhood space is powerful information to better understand the MBs correlation and to accurately predict the future MBs behavior, as detailed in Chap. 4.

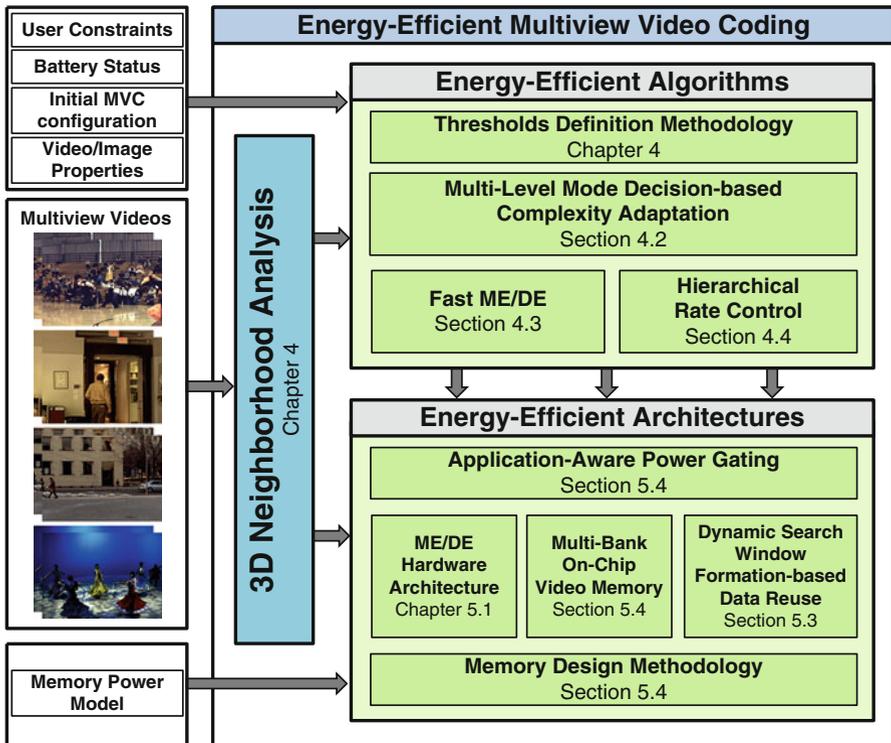


Fig. 3.12 Energy-efficient Multiview Video Coding overview

The algorithm-level contribution is centered in energy reduction through complexity reduction and management. The complexity reduction is reached through a multilevel fast mode decision (see Sect. 4.3) and fast motion and disparity algorithms (for details refer to Sect. 4.4) both based on the 3D-Neighborhood exploitation and image/video properties analysis. The fast MD and ME/DE algorithms are controlled by an energy-aware complexity adaptation scheme (detailed in Sect. 4.3.2) able to handle the energy versus video quality trade-off while considering battery level and encoder state along with external constraints and user's preferences. To avoid the possible quality losses inserted by complexity reduction techniques a hierarchical rate control (HRC) (detailed in Sect. 4.5) featuring both frame- and basic unit-level RC is proposed in order to guarantee a smooth video quality and output bitrate through QP adaptation. Additionally, to provide efficient energy reduction under varying QP scenarios our proposals employ QP-based thresholding according to the methodology presented in Sect. 4.2.

The architectural contribution is focused on the motion and disparity estimation unit and is composed of the ME/DE hardware architecture itself, the application-aware power gating and data-reuse management techniques, and the memory design methodology. A multilevel pipelined ME/DE architectural template is proposed (see details in Sect. 5.1) featuring parallel processing elements, search control, and parallel memory interface initially designed to fit to the fast ME/DE algorithm. The on-chip video memory (see Sect. 5.1.3) sizing and organization were designed considering extensive offline analysis with real video content following our memory design methodology. The on-chip video memory allows sector-level power gating to optimize the energy consumption through implementing a application-aware power gating scheme based on the 3D-Neighborhood knowledge (presented in Sect. 4.1). The external memory communication is optimized while employing reduced on-chip memory through a data-reuse technique based on dynamic search window formation that also exploits the 3D-Neighborhood concept.

3.5.1 3D-Neighborhood

The 3D-Neighborhood is defined as the set of MBs belonging to the neighborhood of the current MB in relation to spatial, temporal, and view/disparity domains. The high coding properties correlation available in the 3D space is discussed and quantified through the statistical analysis presented along Chap. 4. For this reason the 3D-Neighborhood is used to design and control multiple algorithms and architectural decisions proposed in this monograph.

At design time, the offline statistical analysis is used to understand which coding modes are more frequent, the range of motion and disparity vectors, performance and memory requirement variations, bitrate distribution, which neighboring regions are more correlated for distinct cases and video inputs. Such information is used to guide the complete design of the algorithms. Additionally, the offline analysis is used to define the threshold equations according to our thresholding methodology.

The 3D-neighborhood data are also analyzed at run time to perform the actual predictions related to the fast mode decision, fast ME/DE algorithms, and rate control. Also, the data reuse and the memory requirements prediction used to control the power states of the on-chip video memory employ the neighborhood knowledge.

3.5.2 Energy-Efficient Algorithms

In this section are presented the main energy-efficient algorithms proposed in this monograph and detailed along the Chap. 4.

Multilevel Mode Decision-based Complexity Adaptation: We propose a novel dynamic complexity reduction scheme for non-anchor frames in MVC. Our scheme exploits different video statistics and the coding mode correlation in the 3D-Neighborhood to anticipate the more probable prediction modes. It employs a candidate mode-ranking mechanism reinforced with an RDCost-based neighbor confidence level to determine the more probable and less probable prediction modes. The proposed scheme also incorporates an Early SKIP technique that exploits the high occurrence of SKIP MBs in order to reduce the MVC encoder complexity by considering the 3D-Neighborhood correlation and image properties. Two complexity reduction levels named *Relax* and *Aggressive* with different threshold equations are defined. These levels provide a trade-off between energy/complexity reduction and video quality. To limit the propagation of prediction error, the anchor frames are encoded using exhaustive RDO-MD. In this case the prediction error is propagated less due to the availability of a better prediction from the anchor frames of the neighboring GOPs.

The energy-aware complexity adaptation scheme for MVC targeting mobile devices employs several quality-complexity classes (QCCs), such that each class evaluates a certain set of coding modes (thus a certain complexity and energy requirement) and provides a certain video quality. It thereby enables a run-time trade-off between complexity and video quality. To support asymmetric view quality and exploit the binocular suppression properties, views for one eye are encoded with high-quality class and views for the other eye are encoded using a low-quality class. Our complexity adaptation adapts the QCCs for different views at run time depending upon the current battery level.

Fast Motion and Disparity Estimation: Our fast ME/DE algorithm computes the confidence of predictors (motion/disparity vectors of the neighboring MBs) in the 3D-Neighborhood to completely skip the search step. The predictors are classified according to a confidence level and the search pattern is replaced by a reduced number of candidate vectors (up to 13). To exploit this knowledge, accurate motion and disparity fields must be available. Therefore, at least one frame using DE and one using ME must be encoded with a near-optimal searching algorithm. In our scheme, to avoid a significant quality loss, all anchor frames and the frames situated in the middle of the GOP are encoded using the TZ Search algorithm (Tang et al. 2010).

Once the motion and disparity fields are established, all remaining frames are encoded based on predictors available in these fields.

Hierarchical Rate Control: The HRC for MVC employs a joint solution for the multiple actuation levels of rate control. The proposed HRC employs a Model Predictive Control-based rate control that jointly considers GOP-phase and frame-level stimuli to accurately predict the bit allocation and define an optimal control action at coarse grain. This guarantees smooth bitrate and video quality variations along time and view domains while supporting any MVC hierarchical prediction structure. To further optimize the bit allocation within the frames, the HRC implements a Markov Decision Process to refine the control action at BU level taking into consideration image properties to define and prioritize Regions of Interest (RoI). The fine-grained adaptation promotes an increase in objective and subjective video qualities inside the frame. The target bitrate at each time instant is predicted based on the bitrate distribution within the 3D-Neighborhood.

Thresholds Definition Methodology: The energy-efficient algorithms, mainly those based on statistic-based heuristics, are very sensible to the thresholds. For this reason we consider the threshold definition methodology as part of this work. Our schemes employ QP-based threshold equations in order to guarantee proper reaction to changing QP values and keep the energy efficiency. The thresholds for a subset of QPs are derived from extensive correlation statistical analysis of the 3D-Neighborhood. Probability Density Functions (PDF) considering a Gaussian distribution are typically used to model the coding properties distribution. The QP-based threshold equations are then modeled and formulated using polynomial curve fitting from the set of thresholds statically defined.

3.5.3 Energy-Efficient Architectures

The overview of our architectural contribution to the energy-efficient MVC realization is presented in the following. The implementation details are given in Chap. 5.

Motion and Disparity Estimation Hardware Architecture: A pipelined hardware architecture was designed to fit the fast ME/DE algorithm introduced in Sect. 3.5.2 exploiting four levels of parallelism inherent to the MVC prediction structure which are view, frame, reference frame, and MB levels. To reduce the energy consumption related to the memory leakage a multibank on-chip memory and the dynamic window formation-based power gating control are presented. Finally, an application-aware power gating is proposed and integrated to the architectural proposal. The goal of the ME/DE architecture is to deliver the performance for real-time ME/DE for up to 4 views HD1080p while reducing the overall energy consumption. The hardware architecture is composed of five main modules (a) programmable search control unit, (b) shared SAD calculator, (c) on-chip video memory, (d) address generation unit, and (e) energy/complexity-aware control unit.

Multibank On-Chip Video Memory: Our multibank on-chip memory is designed to feed the SAD calculation by employing 16 parallel banks and provide high throughput in order to meet high definitions requirements. Each bank is partitioned into multiple sectors, such that each sector can be individually power-gated to reduced energy through leakage saving. The on-chip video memory behaves as a cache. Thus, it does not require complete reading of the entire search window. Only the required data is prefetched according to an application-aware prefetching technique such as dynamic window formation. The control of the power gating is obtained from the application-aware power gating. The size and the organization of the memory are obtained by an offline analysis of the ME/DE memory and energy requirements within the 3D-Neighborhood.

Memory Design Methodology: Based on the offline memory usage analysis, an algorithm is proposed to determine the size of the on-chip memory by evaluating the trade-off of leakage reduction (on-chip energy) and cache misses (off-chip access energy; result of reduced-sized memory). Afterwards, the organization (banks, sectors) is obtained by considering the throughput constraint. Each bank is partitioned into multiple sectors to enable a fine-grained power management control. The data for each prediction direction is stored in distinct memory sections.

Dynamic search window Formation-Based Data Reuse: Instead of prefetching the complete rectangular search window, a selected partial window is dynamically formed and prefetched for each search stage of a given fast ME/DE search pattern depending upon the search trajectory inferred within the 3D-Neighborhood. In other words, the search window is dynamically expanded depending upon the search history of neighboring MBs and the outcome of previous search stages. The search trajectories of the neighboring MBs and their spatial and temporal properties (variance, SAD, motion, and disparity vectors) are considered to predict at run time the shape of the search window for the current MB. The goals are significantly reducing energy for off-chip memory accesses and reducing the total amount of on-chip memory bits.

Application-Aware Power Gating: One key source of leakage is the big on-chip SRAM memory required to store a big rectangular search window, which is inevitable in case of DE. The unused regions of the rectangular search window indicate a waste of on-chip memory hardware. Therefore, significant leakage reduction may be obtained by reducing the size of the on-chip memory, while considering an analysis of the memory requirements of fast ME/DE schemes. Thus, an application-aware power gating scheme is employed. Depending upon the fast ME/DE search pattern, search direction, MB properties, and 3D-Neighborhood memory usage, the amount of required data is predicted. Only the sectors to store the required data are kept powered on and the remaining sectors are voltage scaled to sleep power states.

Each energy-efficient algorithmic and architectural contribution introduced in this section is detailed in Chaps. 4 and 5. They were designed and evaluated through simulations considering benchmark video sequences and recommended test conditions (Su et al. 2006; ISO/IEC 2011). The simulation setup and the energy reduction gains are presented, discussed and compared to the state of the art in Chap. 6.

3.6 Summary of Application Analysis for Energy and Quality

The computational and energy requirements demanded for optimal MVC encoding are orders of magnitude beyond the reality of current embedded systems. As demonstrated along this section, MVC optimal encoding requires up to 1000 BIPS while current processors delivers about 180 MIPS. In this scenario, state-of-the-art batteries would be able to power the MVC encoder for just a few minutes. Thus, there is a need to reduce the MVC complexity and attack the main sources of energy consumption.

As quantified along this section, mode decision and ME/DE represent more than 90 % of MVC encoder consumption. Moreover, in the ME/DE block, the memory-related energy is dominant in relation to the computation-related energy. Aware of this behavior, a series of energy-oriented contributions are presented.

Along this monograph are presented energy-efficient algorithms and hardware architectures to enable the real-world implementation of the MVC video encoder. Among the algorithms are a Multilevel Fast Mode Decision and a Fast ME/DE algorithm. These solutions employ the 3D-Neighborhood correlation to predict the full RDO-MD or to avoid unnecessary ME/DE searches. Additionally, an Energy-Aware Complexity adaptation algorithm is proposed to enable run-time adaptation in face of varying coding parameters and video inputs. To avoid eventual quality losses posed by these heuristic-based algorithms, a HRC is presented.

A motion and disparity estimation architecture is proposed in order to provide real-time performance and increased energy efficiency to the most complex MVC encoding block. The Fast ME/DE algorithm is considered along with on-chip memory design techniques to reduce energy consumption. Moreover, the on-chip memory is controlled by our Application-Aware Power Gating. The external memory accesses are reduced by the Dynamic Search Window Formation algorithm.

Chapter 4

Energy-Efficient Algorithms for Multiview Video Coding

The energy consumption in MVC encoding is directly related to the high computational effort and the intense memory access driven by the data processing. Therefore, the energy-efficient algorithms for the Multiview Video Coding proposed in this monograph are based on complexity reduction and complexity control techniques. Moreover, in addition to the energy consumption perspective, meaningful complexity reduction is also required at the performance perspective in order to make MVC real-time encoding feasible for real-world embedded devices.

This chapter presents the proposed energy-efficient algorithms targeting complexity reduction for the Multiview Video Coding through fast mode decision and fast motion and disparity estimation techniques. An energy-aware complexity adaptation algorithm designed to offer run-time adaptivity to changing scenarios (battery level, user constraints, video content) of battery-powered embedded devices is further presented. Aware of the rate-distortion losses posed by such complexity reduction techniques we also present a video-quality management technique to avoid visual degradation. The quality management employs a rate control unit able to maximize the video quality for a given target bitrate while providing smooth quality and bitrate variations at spatial, temporal, and disparity domains.

The studies of correlation within the 3D-Neighborhood build the foundation for all algorithms proposed in this chapter. These studies contemplate the analysis of coding mode, motion and disparity fields, and bitrate allocation. Additionally, the profiling of the mode distribution and motion/disparity vectors are key enablers for energy-efficient solutions able to provide high complexity reduction at a negligible cost in terms of coding efficiency.

4.1 3D-Neighborhood Correlation Analysis

4.1.1 Coding Mode Correlation Analysis

In this section is discussed the 3D-Neighborhood correlation considering the coding mode used for the different macroblocks in a video sequence. It discusses the mode distribution profiling and presents a statistical analysis considering coding modes, RDCost correlations, and video properties for multiple multiview video sequences.

4.1.1.1 Coding Mode Distribution Analysis

The graph presented in Fig. 4.1 quantifies the mode distribution in anchor and non-anchor frames of the *Ballroom* and *Vassar* sequences for various QP values (22–37). In anchor frames the mode distribution follows the typical distribution trend of H.264/AVC-based encoding at lower QPs (Huang et al. 2006), i.e., more Intra-coded MBs at lower QP values and more SKIP and large block partitions of Inter-coded MBs at higher QP values. On the contrary, for non-anchor frames, a major portion of the total MBs (50–70 %) is encoded as SKIP for QP>22. The percentage of the SKIP-coded MBs goes up to 93 % (average 63 %) in *Vassar*, a well-known test video which has slow-motion intensity. The second dominant mode is Inter-16×16. Notice that, for QP>27, the percentage distribution of the Intra-coded MBs in non-anchor frames diminishes to less than 1 %.

The uneven mode distribution for non-anchor shows there is a great potential of complexity reduction in the non-anchor frames if the SKIP or Inter-16×16 coding mode are correctly predicted for a MB. In the following analysis, we show that variance/gradient information in conjunction with the coding mode and RDCost correlation in the 3D-Neighborhood provides a good prediction of the SKIP and/or Inter-16×16 coding modes.

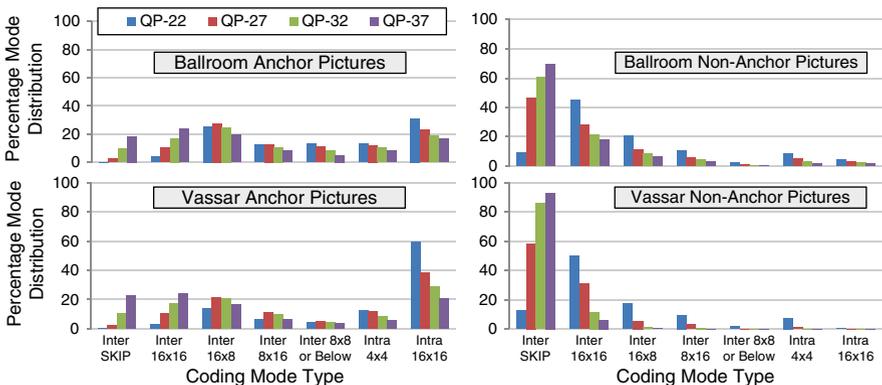


Fig. 4.1 Coding mode distribution

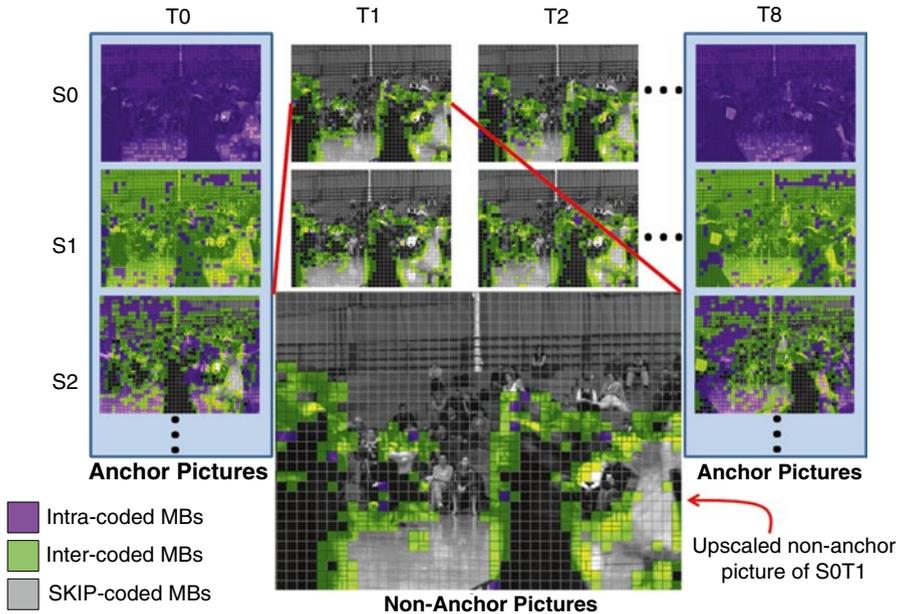


Fig. 4.2 Visual analysis of the coding mode correlation

The mode distribution analysis provides high-level information about the features of a video sequence. This analysis is required for relating the distribution of predictions modes to the video features for a given QP. An in-depth analysis is provided in section “Analyzing the Coding Mode Correlation” where Fig. 4.2 provides a subjective analysis of the optimal mode distribution in the *Ballroom* sequence encoded using the exhaustive RDO-MD.

4.1.1.2 Analyzing the Coding Mode Correlation

The first observation provided by Fig. 4.2 is the distinct mode distribution in the anchor and non-anchor frames. It is noteworthy that the number of SKIP-coded MBs is much higher in the non-anchor frames. This is due to the fact that a higher correlation space is available for non-anchor frames compared to the anchor ones, and consequently, there is higher likelihood to provide a better prediction employing the SKIP mode.

The upscaled frame (S0T1) of *Ballroom* sequence in Fig. 4.2 demonstrates that most of the MBs in the background of the scene (spectators and wall) and partially foreground objects (suits of the dancers and floor) of a non-anchor frame are encoded using the SKIP mode. The MBs at the object borders (dancers) are encoded using temporal/view-prediction modes (i.e., Inter-coded MBs) or spatial-prediction modes (i.e., Intra-coded MBs). Only a few high-textured MBs containing moving spectators in the background are encoded using spatial/temporal/view-prediction modes.

Note in Fig. 4.2 that the MBs belonging to the same region tend to use the same coding mode when considering spatial, temporal, or disparity collocated MBs. For instance, consider frame SOT1, the dancer borders share the same coding mode used by the spatial neighboring MBs that belong to this border. Also, the same coding mode tends to be shared with temporal and disparity collocated MBs in frames SOT2 and S1T1, respectively.

However, different neighboring MBs in the 3D-Neighborhood exhibit different amount of correlation to the current MB. Figure 4.3 shows the coding mode *hits* (averaged over various QPs and video sequences) using the exhaustive RDO-MD. A coding mode *hit* corresponds to the case when the optimal coding mode of a neighbor is exactly the same as that of the current MB. Otherwise it is given as a coding mode *miss*. The coordinates on the x - and y -axis correspond to the MB number in the corresponding column and row of a frame, e.g., (2,4) means the 2nd MB of the 4th row. The eight neighbor frames in the 3D domain are evaluated and named according to the cardinal points presented in Fig. 4.3. There are total 44 neighbors: 4 spatial, 18 temporal, 18 disparity, and 4 disparity-temporal. Note, the disparity and disparity-temporal neighbors consider the GDV rounded to an integer number of MBs.

Figure 4.3 illustrates that the spatial neighbors in the current frame exhibit the highest coding mode correlation to the current MB (i.e., *hits* > 70 %) followed by the disparity neighbors in the *North* and the *South* view frames (i.e., *hits* > 66 %). The coding mode *hits* of the disparity neighbors is less than that of the spatial neighbors due to the variations near the object borders and an inaccuracy in the GDV. The lower number of *hits* for the temporal and disparity-temporal neighbors is basically due to the motion properties. On overall, for non-anchor frames, in more than 98 % of the cases the optimal coding mode of an MB is present in the 3D-Neighborhood. It means that by testing the coding modes of all 44 neighbors it is highly probable to find the optimal coding mode for the current MB (more than 98 % of the MBs in the current frame). Moreover, due to the availability of a limited set of optimal coding modes in the non-anchor frames (typically much less than the number of modes tested in an exhaustive RDO-MD), a significant complexity reduction may be achieved.

As discussed above, there is a big potential of finding the optimal encoding mode in the 3D-Neighborhood. However, a big number of different coding modes may exist in this neighborhood. Thus, in order to reduce the number of probable modes, additional information is needed. In this monograph we consider video and image properties and the RDCost as additional information. The study related to these properties is presented in the following sections.

4.1.1.3 Analyzing the Video Properties

Along our studies multiple video and image properties—including variance, brightness, edges and gradient—were evaluated in order to provide useful information to build fast mode decision algorithms. Among these properties, variance and gradient

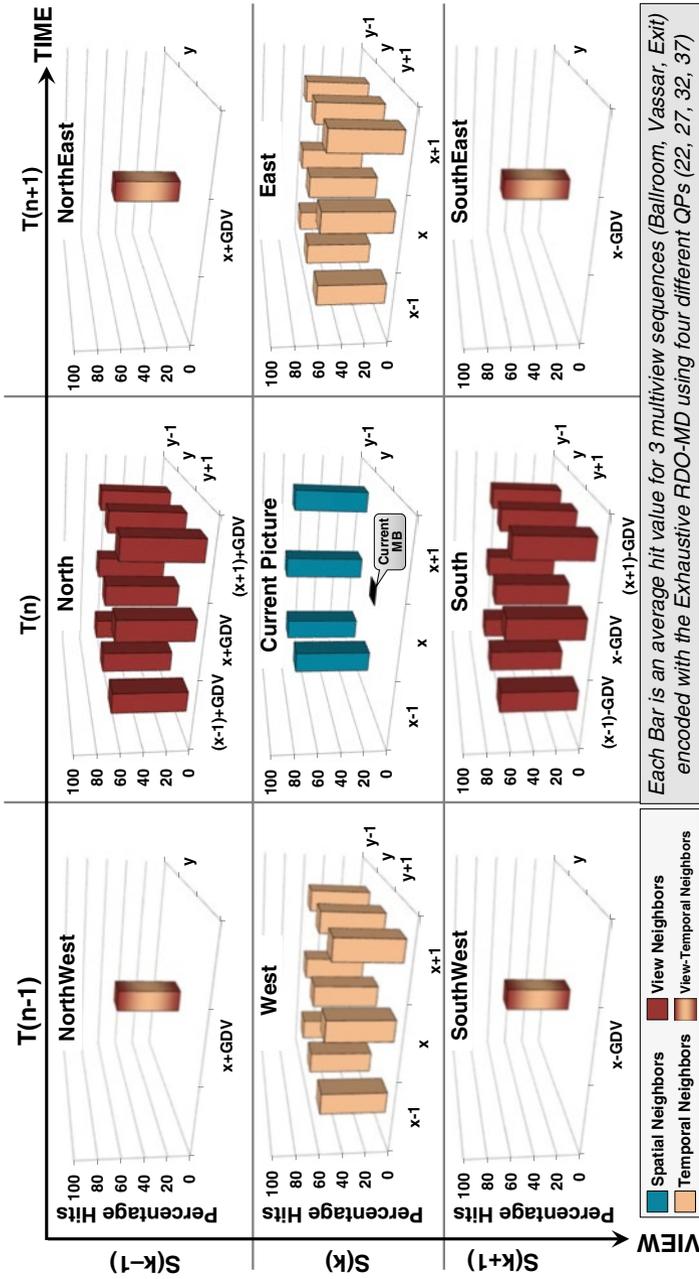


Fig. 4.3 Coding mode hits in the 3D-neighborhood

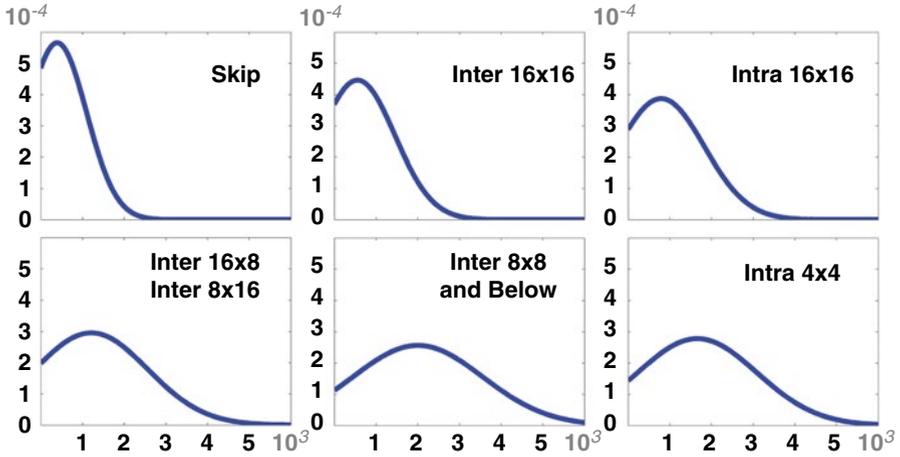


Fig. 4.4 Variance PDF for different coding modes

information showed to be the most helpful to identify highly correlated neighboring MBs and their possible coding modes. The complete evaluation based on statistical analysis for the variance is presented below. For that, multiple video sequences were considered while assuming a Gaussian distribution for the video properties. The variance is defined in Eq. (4.1) while horizontal (Δx) and vertical (Δy) gradients are determined by Eq. (4.2), where ρ_i represents the pixels of a MB:

$$Var_{MB} = \sum_{i=1}^{256} (\rho_i - \rho_{AVG})^2; \rho_{AVG} = \left(\sum_{i=1}^{256} \rho_i + 128 \right) \gg 8, \quad (4.1)$$

$$\Delta x = \left(\sum_{i=0}^{15} \sum_{j=0}^{15} \left| \frac{\partial f}{\partial x} \right| + 128 \right) / 256, \quad \Delta y = \left(\sum_{i=0}^{15} \sum_{j=0}^{15} \left| \frac{\partial f}{\partial y} \right| + 128 \right) / 256; \quad (4.2)$$

$$\frac{\partial f}{\partial x} = \rho(i, j) - \rho(i-1, j), \quad \frac{\partial f}{\partial y} = \rho(i, j) - \rho(i, j-1)$$

Figure 4.4 shows different PDF (Probability Density Function) plots for the variance related to various coding modes. It is noticeable that the peaks for the SKIP and Inter/Intra- 16×16 modes are at 400 and 700, respectively. Therefore, MBs with low variance are more likely to be encoded as SKIP than Inter/Intra- 16×16 . On the contrary, MBs with high variance (1,500–2,500) are more likely to be encoded using smaller block partitions. The PDFs for gradient are omitted; however, they have a similar distribution to that of the variance.

Since there is a considerable overlap between the PDFs of 16×16 and smaller block partitions, in order to obtain a more robust/accurate prediction about the coding modes, RDCost (see section “Analyzing the RDCost”) and coding mode correlation in the 3D-Neighborhood are considered along with the variance and gradient information.

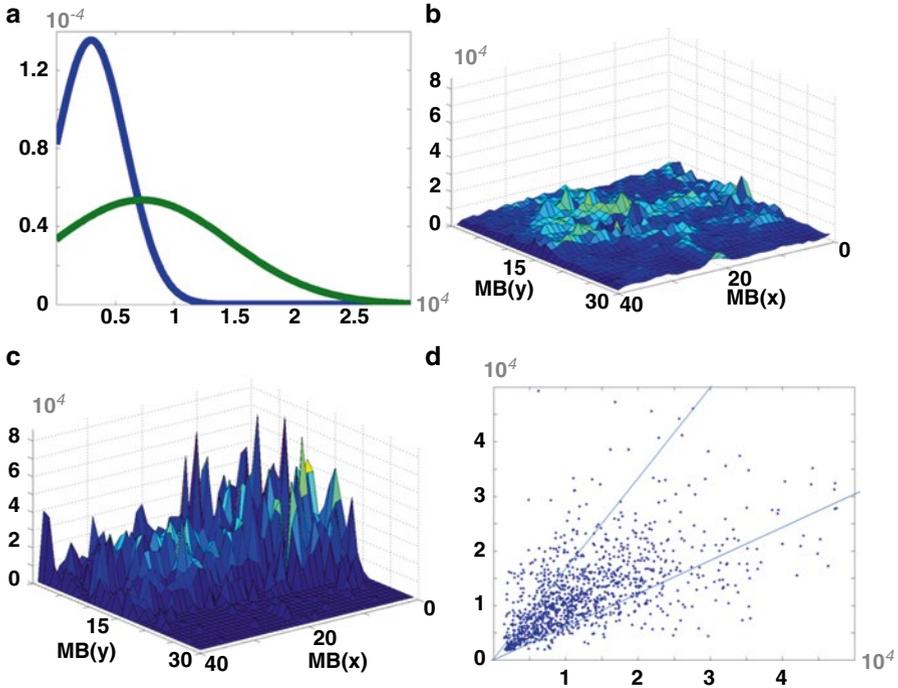


Fig. 4.5 (a) PDF for RDCost difference (between the current and the neighboring MBs) for SKIP *hit* and *miss*; (b, c) Surface plots of RDCost difference for the SKIP coding mode hit and miss; (d) RDCost prediction error for spatial neighbors

4.1.1.4 Analyzing the RDCost

In order to determine which neighbor has a probable coding mode *hit* or *miss*, we compute the difference between RDCost of a neighbor and the predicted RDCost of the current MB (as the actual RDCost is not available before the RDO-MD process). In the following we analyze the relationship between the RDCost difference and the above discussed coding mode *hit/miss*. Figure 4.5a presents the PDF for the RDCost difference for coding mode *hits* and *misses* in case of the SKIP mode. The PDF shows that a SKIP coding mode can be predicted with a high probability of a *hit* when the variance of an MB is low. Figure 4.5b, c shows MB-wise surface plot of the RDCost difference (averaged over all frames of the *Ballroom* sequence) for *hits* and *misses*, respectively. These plots demonstrate that most of the *hits* occur when the RDCost difference is below 10 k, while the number of *miss* increases when the value of RDCost difference goes above 70 k. This behavior also conforms to the PDFs in Fig. 4.5a. This analysis shows that the value of RDCost difference provides a good hint for a *hit* in case of a SKIP coding mode. Similar behavior was observed in the *hit* and *miss* PDFs for other coding modes. Here, we discuss the PDFs for the SKIP mode as an example since it is the dominant coding mode in non-anchor frames, especially for higher QP values.

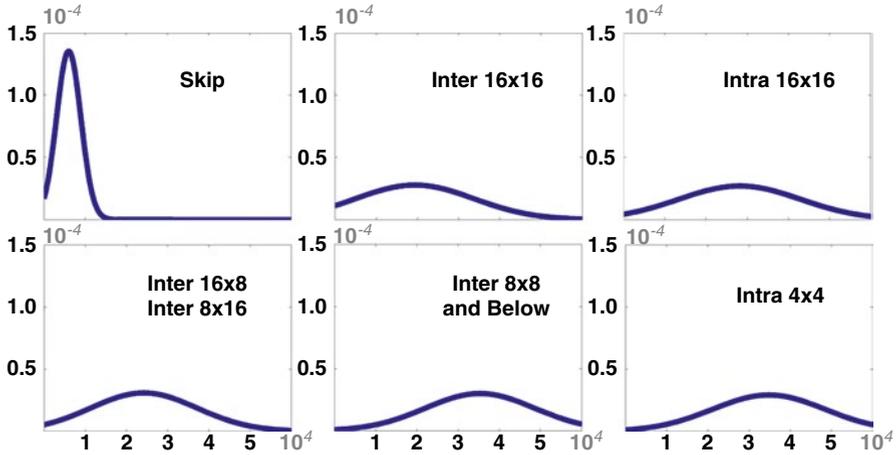


Fig. 4.6 PDF of RDCost for different prediction modes in Ballroom sequence

Figure 4.6 presents the PDFs of predicted RDCost for different coding modes. Variable shapes of the PDFs already hint towards the exclusion of improbable mode for a given value of the predicted RDCost. Since a good prediction is important to determine a near-optimal coding mode, we have evaluated the accuracy of the predicted RDCost and optimal RDCost to analyze the risk of misprediction.

Once the RDCost is not available without the exhaustive RDO-MD we tested different predictors for the current MB RDCost in the 3D-Neighborhood. After analyzing the mean and median RDCosts predictors, we have determined that the median RDCost of the spatial neighbors [see Eq. (4.3)] provides the closest match to the optimal RDCost. In Eq. (4.3) S_L , S_T , and S_{TL} represent left, top, and top/left spatial neighbors, respectively:

$$RDCost_{PredCurr} = Median(S_L, S_T, S_{TL}). \quad (4.3)$$

Figure 4.5d shows the optimal RDCost vs. predicted RDCost for each MB in 36 frames of the *Vassar* third view. It illustrates a high correlation between the two values (approximately 0.88). Figure 4.7 shows the error surface for the predicted RDCost compared to the optimal RDCost highlighting the regions of misprediction (i.e., borders of the moving objects).

4.1.1.5 Coding Mode Analysis Summary

Our detailed analysis illustrates that it is possible to accurately predict the optimal coding modes, mainly for non-anchor frames, if the coding mode distribution, video

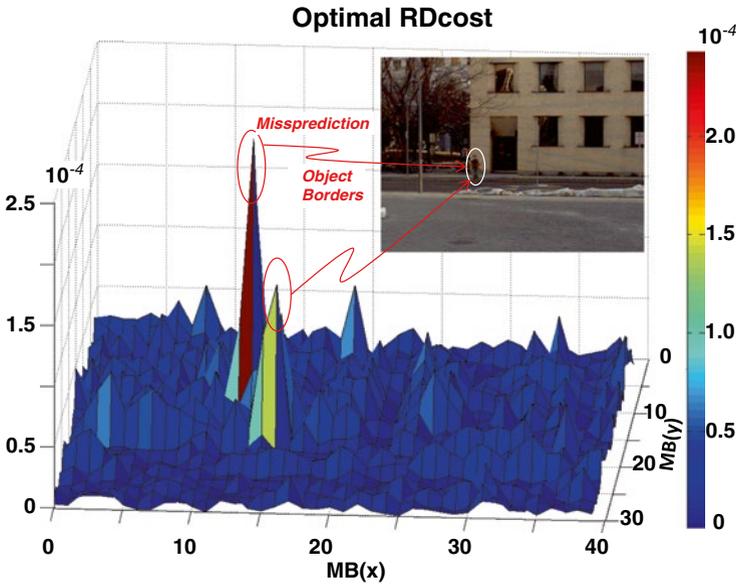


Fig. 4.7 Average RDCost prediction error for spatial neighbors in Vassar Sequence

statistics, and RDCost correlation in the 3D-Neighborhood are considered. It leads to a high potential of complexity and energy consumption reduction during the MVC encoding process. The main conclusions that enable our fast algorithms are summarized below.

- SKIP MB is the dominant prediction mode (47–97 %) in non-anchor frames for $QP > 27$.
- The inter-coded MBs with big partitions are dominant over the smaller partitions and intra-coded MBs in non-anchor frames.
- Different prediction modes exhibit different variance, gradient, and RDCost properties which may be used to identify more- and less-probable coding modes for fast mode decision.
- The spatial, temporal, and view neighborhood exhibit up to 77 %, 62 %, and 69 % coding mode hits; thus there is a high probability to find a correct prediction of the coding mode in the 3D-Neighborhood.
- RDCost provides means to identify neighbors with relatively high *hit* probability at run time.
- The median RDCost of the spatial neighbors provides an accurate RDcost prediction for the current MB.

Mispredictions may occur at the object borders, objects with high motion, and in the foreground objects where the displacement is different from GDV.

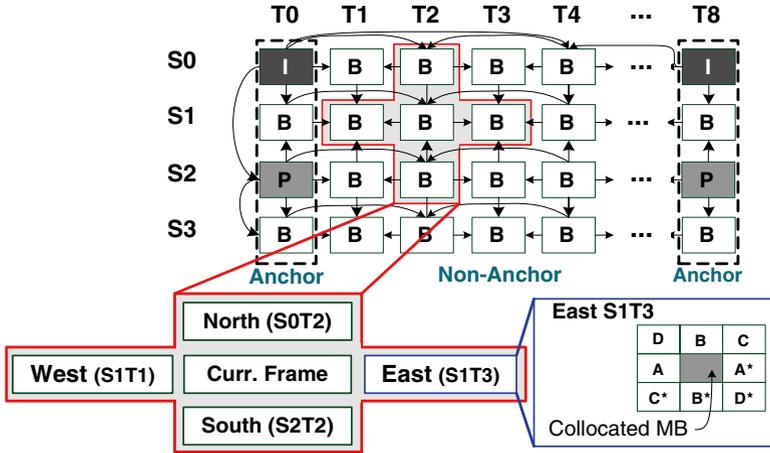


Fig. 4.8 MVC prediction structure and 3D-neighborhood details

4.1.2 Motion Correlation Analysis

Before proceeding to the motion/disparity vectors correlation analysis, we briefly recall the basic prediction structure of MVC to a level of detail necessary to understand our analysis. MVC uses the motion and disparity estimation tools to eliminate the temporal and view redundancies between frames, respectively. The prediction structure used in this work is presented in Fig. 4.8. *I* squares represent intra-predicted frames (i.e., no ME/DE is used), *P* are frames using unidirectional prediction or estimation (in this example the *P* frames use only DE in one direction), and *B* frames use bidirectional prediction having reference frames in at least two directions. The arrows represent the prediction directions: frames at the tail side act as reference frames to the frames pointed by the arrowheads. Note that some frames have up to four prediction directions. In order to provide random access points, the video sequence is segmented in Groups of Pictures (GOPs) where the frames located at the GOP borders are known as anchor frames and are encoded with no reference to the previous GOP. All other frames are called non-anchor frames.

In our observations we noticed that the same objects in a 3D scene are typically present in different views (except for occlusions). Consequently, the motion perceived in one view is directly related to the motion perceived in the neighboring views (Deng et al. 2009). Moreover, considering parallel cameras, the motion field is similar in these views (Kim et al. 2007). Analogously, the disparity of one given object perceived in two cameras remains the same for different time instances when just translational motion occurs. Even for other kinds of motion the disparity is highly correlated.

Based on these observations an analysis of the motion and disparity vectors is performed to quantify the MV/DV correlation (here the term correlation is subjectively used, it is defined as the difference between the predictors and the optimal

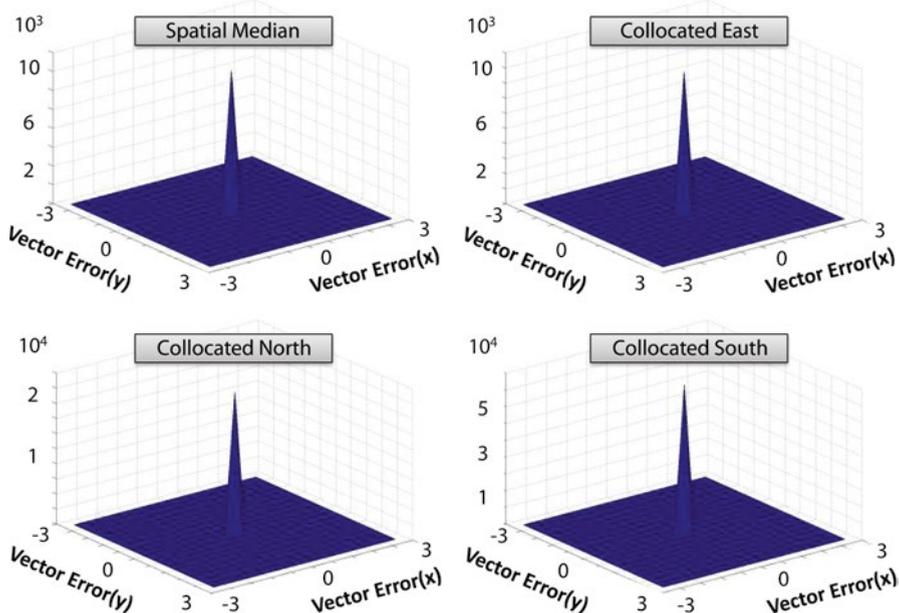


Fig. 4.9 MV/DV error distribution between predictors and optimal vector (Ballroom, Vassar)

vector, i.e., MV/DV error) in the 3D-Neighborhood (i.e., spatial, temporal, and view domains). A set composed of one spatial median predictor, six temporal predictors, and six disparity predictors is analyzed. The temporal predictors are selected from the previous and next frames (in the displaying order) called West and East neighbor frames, respectively. For each neighboring frame, three predictors are calculated. They are (a) the collocated MB (MB in the reference frame with the same relative position of the current MB), (b) median up (using the median formula specified by the MVC standard (JVT 2009b) to calculate the spatial predictor), and median down (median of A^* , B^* , C^* , and D^* as shown in Fig. 4.8). The disparity predictors from the North and South neighboring view frames are obtained by considering the GDV.

Figure 4.9 illustrates the MV/DV error distribution for *Vassar* (low motion) and *Ballroom* (high motion) test video sequences in the 3D-Neighborhood. Each plot represents the difference between a given predictor (in this case for the spatial predictor and three collocated predictors in different neighboring frames) and the optimal vector of the current MB. It shows that for the majority of the cases, the predictor vectors have similar values in comparison to the optimal vector. Even, most of the predictors have exactly the same value of the optimal vector. Our analysis shows that this observation is valid for the other nine predictors in all directions of the 3D-Neighborhood as depicted in Fig. 4.9 (only few error plots are shown here).

To quantify the MV/DV error distribution in the 3D-Neighborhood, several experiments were carried out to measure the frequency in which a given predictor is equal to the optimal vector (i.e., $MV_{Pred} = MV_{Curr}$). When this condition is satisfied,

Table 4.1 Predictors hit rate and availability

Predictor	Neighbor frame	Hit [%]	Available [%]
Spatial	n.a.	94.12	99.90
All	<i>West</i>	96.94	51.52
	<i>East</i>	97.93	60.30
	<i>North</i>	97.94	65.40
	<i>South</i>	98.67	21.29
Collocated	<i>West</i>	58.43	99.90
	<i>East</i>	66.79	99.90
	<i>North</i>	95.39	72.39
	<i>South</i>	96.75	23.48
Median Up	<i>West</i>	54.74	99.90
	<i>East</i>	63.78	99.90
	<i>North</i>	93.17	73.99
	<i>South</i>	94.61	23.98
Median Down	<i>West</i>	54.99	99.89
	<i>East</i>	63.92	99.89
	<i>North</i>	93.21	74.13
	<i>South</i>	94.70	23.93

it is denoted as a so-called *hit*. A set of different conditions was defined including, for example, the case when all predictors of a given neighbor frame (collocated, median up, and median down) are *hits*. Table 4.1 presents the detailed information on the vector *hits* where the *Availability* is the percentage of cases when that predictor is available. The disparity predictors present the higher number of *hits* followed by the spatial and temporal predictors. Considering the quality of predictors in the same neighboring frame, the collocated predictors present better results in relation to median up and median down. The latter two present similar number of *hits*. In the case where all predictors of a given neighboring frame are available, the predictor is highly accurate providing up to 98 % *hits*.

In conclusion, there is a high vector correlation available in the 3D-Neighborhood that can be exploited during the ME/DE processing. Once the predictors point to the same region as the optimal vector, there is no need for search patterns exploiting a large search range. Moreover, for most of the cases the predictors' accuracy is enough to completely avoid the ME/DE search and refinement stages.

4.1.3 Bitrate Correlation Analysis

In this section we present a detailed bitrate distribution analysis to provide a better understanding towards the bitrate distribution during the MVC encoding process and its correlation with spatial, temporal, and disparity neighborhood. The analysis is presented in a top-down approach starting with the view-level-related discussion,

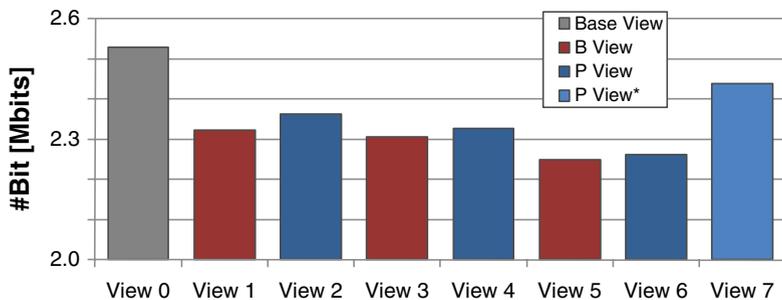


Fig. 4.10 View-level bitrate distribution (Flamenco2, QP=32)

following to frame-level and concluding with BU-level considerations. For that, we used eight views of the *Flamenco2* VGA video sequence encoded at a fixed QP, that is, without rate control, for an IBP view coding (0-2-1-4-3-6-5-7) order and hierarchical bi-prediction at temporal domain. One basic unit is defined as one macroblock.

Figure 4.10 shows the uneven bitrate distribution along different views. This distribution is highly related to the prediction hierarchy inside a GOP. The View 0 or Base View is encoded independently with no inter-view prediction. It leads to reduced possibilities of prediction, and consequently, worse prediction, more residues, and higher bitrate. B-Views (Views 1, 3, and 5) fully exploit the inter-view correlation by performing disparity estimation (in addition to spatial and temporal predictions) to upper and bottom neighboring views. This increased prediction decision space results in improved prediction quality and tends to lead to reduced bitrates. P-Views (Views 1, 3, 5, and 7) represent the intermediate case performing disparity estimation in relation to a single neighboring view. P-Views typically present bitrate in the range between Base View and B-Views bitrates. Note, in Fig. 4.10 the View 7 is a P-View, but its reference view is closer if compared to other P-Views. While View 2 is two views distant to its reference view (View 0), View 7 is just one view distant to View 6. It usually results in a reduced bitrate for View 7 due better disparity estimation prediction.

The bitrate relations associated with prediction hierarchy, however, are not always true and vary with the video/image properties of each view. For instance, in the example provided in Fig. 4.10, View 6 (P-View) presents reduced bitrate in relation to View 1 and View 3 (both B-Views). Thus, we may conclude that even employing bi-prediction at disparity domain the Views 1 and 3 are harder to predict in relation to View 6 and produce higher bitrate. Similar observation is the increased bitrate generated by View 7 if compared to other P-Views. Reduced bitrate is expected, but for View 7 an increased bitrate is measured. These observations show that besides the relation to the prediction structure (as discussed above), the bitrate distribution has a high dependence on the video content of each view. Hard-to-predict views typically present high texture and/or high motion/disparity objects and require more bits to reach a given video quality.

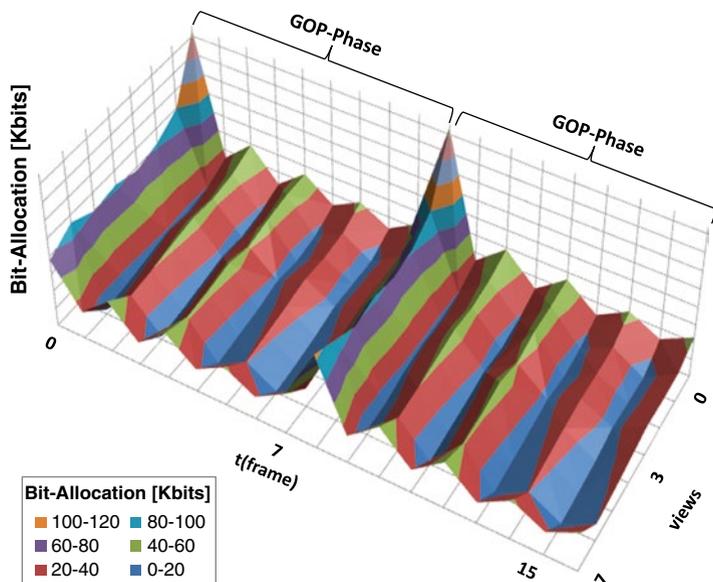


Fig. 4.11 Frame-level bitrate distribution for two GGOPs (Flamenco2, QP=32)



Fig. 4.12 Basic unit-level bitrate distribution (Flamenco2, QP=32)

The bitrate distribution at frame level presented in Fig. 4.11 shows that inside each GOP the frames that present higher bitrate are located at lower hierarchical prediction levels. This is related to the distance of temporal references; the farther the reference the more difficult is to find a good prediction. Therefore, more error is inserted resulting in higher bitrates. In B-Views this effect is attenuated once this view is less dependent on temporal references due to the higher availability of disparity references. Figure 4.11 illustrates that for neighboring GGOPs the frames at same relative position exhibit similar and periodic rate distribution pattern, the GOP-Phase.

Inside each frame the number of bits generated for each BU is also related to the video content. Figure 4.12 shows that the homogeneous and low motion/disparity

background requires lower bitrate if compared to the dancers’ region and to the textured floor for a similar quality. However, the Human Visual System (HVS) requires a higher level of details for textured and border regions to perceive good quality, and consequently, these regions deserve higher objective quality. Therefore, textured regions must be detected and receive further increased number of bits during the encoding process through QP reduction.

Summary: The frame-level bitrate distribution depends on the prediction hierarchy and the video content of each frame. Due to correlation of video content, an effective rate control must consider the neighboring frames at temporal, disparity, and GOP-phase domains. The video properties have to be considered at BU level in order to locate and prioritize regions that require higher quality.

4.2 Thresholds

To determine the thresholds for our energy-efficient algorithms a statistical analysis is performed. We analyze the Probability Density Function (PDF) of a set of MBs sharing the same image/coding property (for instance, same MB type, similar motion vectors, bitrate distribution range, etc.) in order to infer relations between this property and the encoder decisions such as coding mode and motion vectors. From this statistical distribution a region of interest is defined according to the level of confidence designed for the energy-efficient algorithm. Consider the PDF for Gaussian (or Normal) distribution presented in Fig. 4.13 where μ and σ denote the average and standard deviation of an image property, respectively. Consider also that a region of confidence including $P\%$ of all MBs is required. In this case, a pair of threshold points that cover $P\%$ of the area under the bell curve must be defined. The percentage area under the curve $P\%$, represented by the gray filled are in Fig. 4.13, is calculated using Eqs. (4.4) and (4.5):

$$F(\mu + n\sigma; \mu, \sigma^2) - F(\mu - n\sigma; \mu, \sigma^2) = \Phi(n) - \Phi(-n), \tag{4.4}$$

$$P\% = [\Phi(n) - \Phi(-n)] \times 100. \tag{4.5}$$

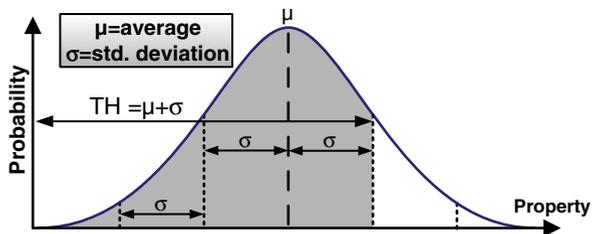
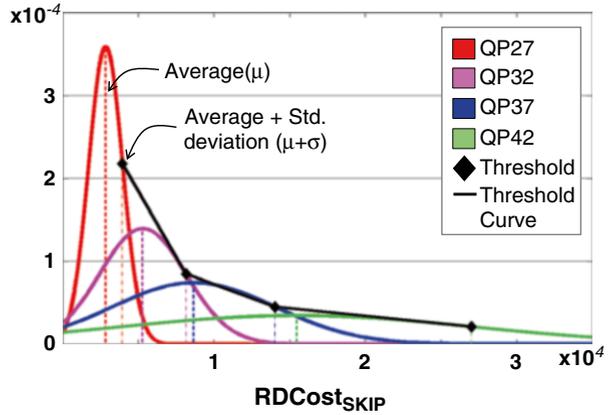


Fig. 4.13 PDF showing the area of high probability as the shaded region

Fig. 4.14 PDF of RDCost for SKIP MBs



Note that the relation between an image property and the coding property to be inferred from it varies with the changing quantization parameter (QP). This comes from the fact that the QP parameter influences on the decisions taking by the encoder. For example, considering the mode decision process, the QP changes the quantization itself but also changes the λ parameter [see Eq. (2.3)] used for balancing the quality vs. bitrate trade-off. Therefore, the presented statistical study needs to be replicated for different QPs. To avoid the complete data extraction and statistical analysis for every single QP value (that goes from 0 to 51 while the practical use typically goes from 22 to 37) we analyze a set of QPs and derive a generic equation for every QP. The generic equation is derived using polynomial curve fitting.

To provide a practical example, we present the threshold derivation for the RDCost property of SKIP MBs encoded using the exhaustive RDO-MD. Figure 4.14 shows the PDFs for the *Vassar* sequence encoded using various QP values. Notice that the PDF for QP 27 shows a concentrated distribution centered in a relatively low RDCost range, i.e., a small average (μ) and standard deviation (σ). Contrarily, the PDFs for relatively high QPs (32–42) exhibit a low peak centered in a relatively high RDCost range.

Recall that the goal is to define thresholds for detecting SKIP MBs with a given confidence considering a Gaussian distribution for the RDCost. In this problem, the confidence region is defined from zero to a threshold point defined in terms of average (μ) and standard deviation (σ)— $TH_{RDCost} = \mu + n\sigma$ —where n represents a multiplier factor for the standard deviation. Higher n represents that more MBs will attend this condition. For the example presented in Fig. 4.13, the threshold was defined as $TH_{RDCost} = \mu + \sigma$ ($n=1$) and all MBs with $RDCost < \mu + \sigma$ belong to the high-confidence zone represented by the gray-filled area in Fig. 4.13. The region of confidence is defined as follows in Eq. (4.6):

$$F(\mu + \sigma; \mu, \sigma^2) - F(0; \mu, \sigma^2) \approx 0.84. \quad (4.6)$$

Equation (4.6) shows that up to 84 % MBs belong to the high-confidence region. We define these points of high probability as the RDCost thresholds for a set of QP

Fig. 4.15 Threshold curves for RDCost

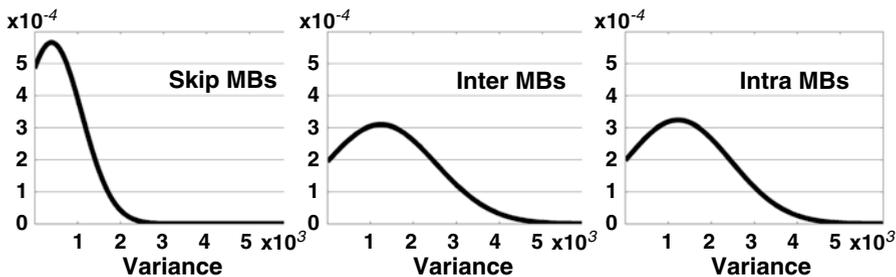
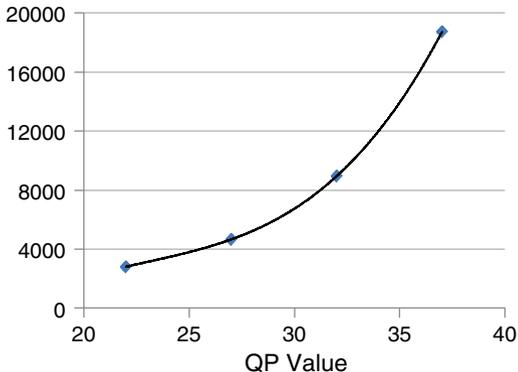


Fig. 4.16 PDF of variance for different prediction modes

represented as diamond points in the Fig. 4.14. Different points for different QPs are used to derive the QP-based threshold Eq. (4.7) using polynomial curve fitting in order to extend to any QP value, as depicted in Fig. 4.15:

$$TH_{RD} = 4.06QP^3 - 279.89QP^2 + 6755.90QP - 53541. \tag{4.7}$$

Figure 4.16 demonstrates a similar statistical analysis using the variance property of a SKIP MB. It shows that the SKIP MBs have a PDF peak in the low variance range compared to other inter/intra-coded MBs. The variance thresholds (TH_{Var}) for predicting a SKIP coding mode are computed in the same way as of TH_{RD} , i.e., $Var_{MB} < \mu + \sigma$. The QP-based threshold is given by Eq. (4.8):

$$TH_{Var} = 0.02QP^3 - 2.02QP^2 + 72.30QP + 196.04. \tag{4.8}$$

The presented statistical threshold derivation strategy was widely applied in our algorithms for different problems, video properties, and values of n . However, to avoid repetition, the thresholding strategy is not discussed again but referred to this section. Only threshold equations are presented for the proposed algorithms.

4.3 Multilevel Mode Decision-based Complexity Adaptation

The Multilevel Mode Decision-based complexity adaptation is presented along this section encapsulating the concepts of Early SKIP prediction, 3D-Neighborhood mode ranking, and image/video properties. The complexity adaptation technique considers the battery status to define the quality vs. complexity operation point.

4.3.1 Multilevel Fast Mode Decision

In this section we propose a complete multilevel mode decision scheme based on the 3D-Neighborhood correlation and exploitation of additional statistical and video information. It incorporates an early SKIP prediction scheme in order to reduce the number of coding modes evaluated during the encoding process.

The detailed flowchart of our multilevel fast mode decision for non-anchor in MVC is presented in Fig. 4.17. The scheme operates in six phases (1) RDCost-based confidence-level ranking, (2) early SKIP prediction, (3) evaluating high-confidence modes, (4) evaluating low-confidence modes, (5) video properties-based mode decision, and (6) size/direction-based mode decision. At the end of each phase (except for phase 1), a condition is evaluated for early termination of the scheme. We explain these phases in the subsequent sections. Similar description is presented in (Zatt et al. 2011b).

4.3.1.1 RDCost Confidence-Level Ranking

Firstly, the 3D-Neighborhood information is fetched and the RDCost for the current MB is predicted using the spatial neighbors considering their high ratio of coding mode *hits* [Eq. (4.9)]. A list of candidate prediction modes (*CandidateList*) is formed from the 3D-Neighborhood. Each candidate mode is associated with a rank value (R_{MODE}). This value is calculated as the accumulated confidence level of the neighbors with the similar coding mode ($CL_{NBi}(Mode)$; Eqs. (4.10) and (4.11)). This confidence level of a neighbor is computed by evaluating the normalized difference (*NDiff*) between its actual RDCost and the predictive RDCost for the current MB [Eq. (4.12)]. Note that the confidence-level calculation depends upon the quality of RDCost prediction (section “Analyzing the RDCost”). The candidate list is then sorted according to the rank value [Eq. (4.13)]:

$$RDCost_{PredCurr} = Median(S_L, S_T, S_{TL}), \quad (4.9)$$

$$R_{MODE} = \sum_{i=1}^{44} CL_{NBi}(Mode), \quad (4.10)$$

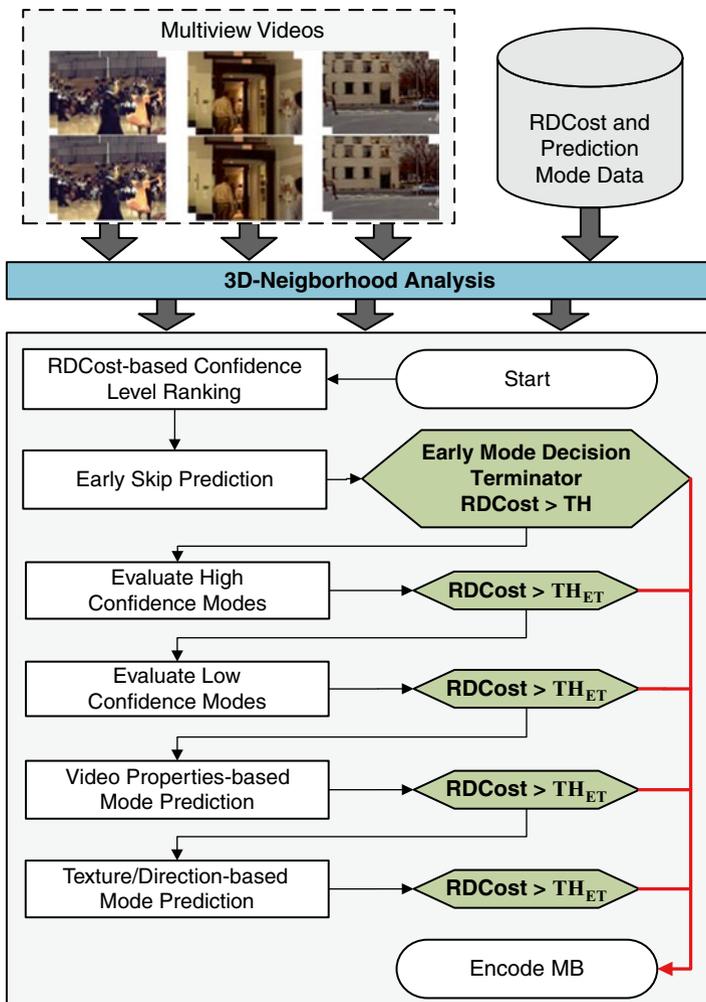


Fig. 4.17 Overview of the multilevel fast mode decision

$$CL_{NB_i}(Mode) = (Clip(NDiff(NB_i), 0, 1)), \quad (4.11)$$

$$NDiff(NB_i) = 1 - Abs(RDCost_{PredCurr} - RDCost_N) / Diff_{MAX} \quad (4.12)$$

$$CandidateList = Sort(R_{SKIP}, R_{INTER16 \times 16}, \dots, R_{INTRA4 \times 4}), \quad (4.13)$$

where NB_i is the i th neighbor, and $Diff_{MAX}$ is the maximum RDCost difference. The Sort function in Eq. (4.13) sorts the values in a descending order.

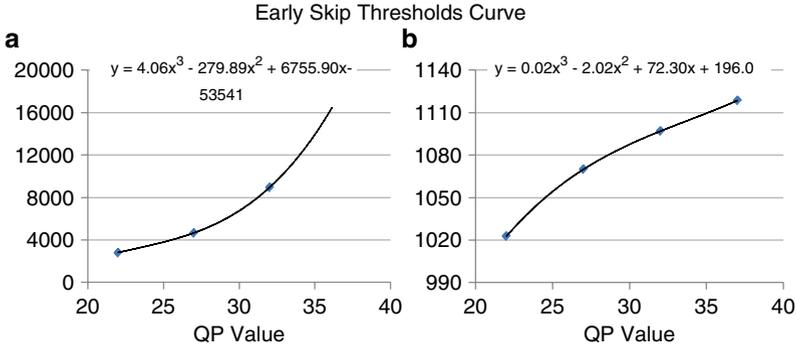


Fig. 4.18 Early SKIP threshold curves for (a) RDCost and (b) Variance

4.3.1.2 Early SKIP Prediction

Based on the analysis of high SKIP MBs distribution in non-anchor frames (section “Coding Mode Distribution Analysis”), our scheme employs an early SKIP prediction. In case a SKIP mode is correctly predicted, significant complexity reduction is obtained as the ME and DE are entirely skipped. This early SKIP mode prediction is only performed if sufficient correlation is available in the 3D-Neighborhood. To avoid a misprediction (that may result in significant PSNR loss) the early SKIP prediction depends upon three conditions considering the mode rank, variance, and RDCost, as presented in Eq. (4.14):

$$\begin{aligned}
 & ((R_{SKIP} > TH_{Rank}) \&\& \\
 \text{EarlySKIP} = & (\text{Variance} < TH_{Var}) \&\& \\
 & (RDCost_{PredCurr} < TH_{RD})) \quad (4.14)
 \end{aligned}$$

The QP-based thresholds for RDCost ($TH_{RDCost_{ES}}$) and variance ($TH_{Var_{ES}}$) were obtained using the corresponding PDF analysis. The area of high probability is considered as the average plus one standard deviation. A threshold is thereby given as $TH = \mu + \sigma$. The PDFs for four different QP values are used to determine four thresholds at different QPs. A QP-based threshold formulation is obtained using the polynomial curve fitting. Figure 4.18 presents thresholds ($TH_{RDCost_{ES}}$ and $TH_{Var_{ES}}$) for four QPs and the corresponding curve fitting. The threshold for ranks ($TH_{R_{ES}}$) was obtained (using an exhaustive analysis) as 15 % of the total confidence level accumulated on the entire *CandidateList* (i.e., sum the ranks of all modes). The Early SKIP is also discussed in (Zatt et al. 2010).

4.3.1.3 Early Mode Decision Terminator

After the early SKIP mode prediction, the tested mode is evaluated for the early mode decision termination. If the tested RDCost is bigger than the threshold TH_{ET} ,

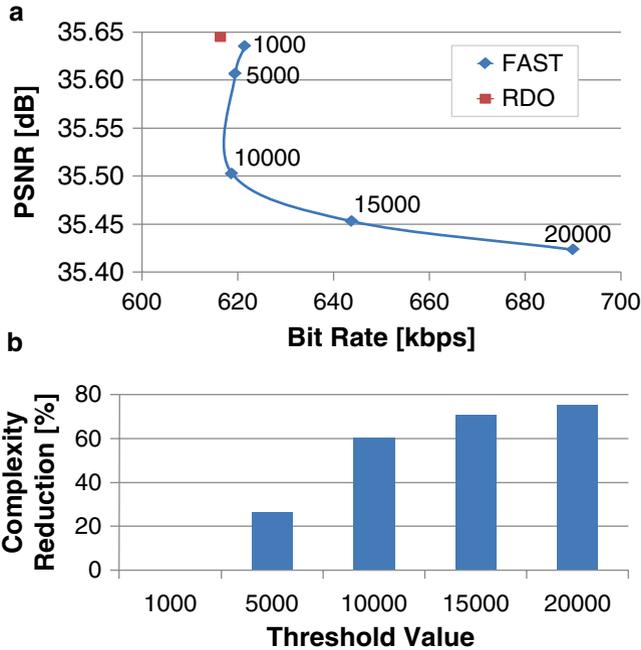


Fig. 4.19 Evaluation of thresholds for early termination (Ballroom, QP=32)

the mode decision proceeds to the next phase. Otherwise, the mode decision is terminated and the best tested mode is used for encoding the current MB.

The threshold for early mode decision termination controls the achieved complexity reduction and the resulting PSNR loss. An excessively high value threshold provides high complexity reduction at the cost of severe PSNR loss. We have performed an exhaustive analysis to determine these thresholds. Figure 4.19 shows the RD curve for five different test threshold values and their corresponding complexity reduction (bars) for QP=32. It is noted that $TH_{ET}=5,000$ provides minimal PSNR loss and low complexity reduction, while $TH_{ET} \geq 10,000$ provides a high complexity reduction at the cost of considerable PSNR loss (i.e., >0.15 dB). In order to provide a trade-off between achieved complexity and the resulting quality loss, we propose two complexity reduction levels or complexity reduction strengths:

- *Relax* complexity reduction: it provides a reasonable complexity reduction while considering a low PSNR loss.
- *Aggressive* complexity reduction: it provides a high complexity reduction at the cost of a slightly higher PSNR loss (but still visually unnoticeable in many cases, as we will show in results chapter).

From an exhaustive analysis of various multiview sequences (encoded using exhaustive RDO-MD), we obtained the plots and QP-based equations for *Relax* (blue) and *Aggressive* (red) complexity reduction (see Fig. 4.20).

This early termination is employed after each phase of our dynamic complexity reduction scheme as explained in the subsequent sections.

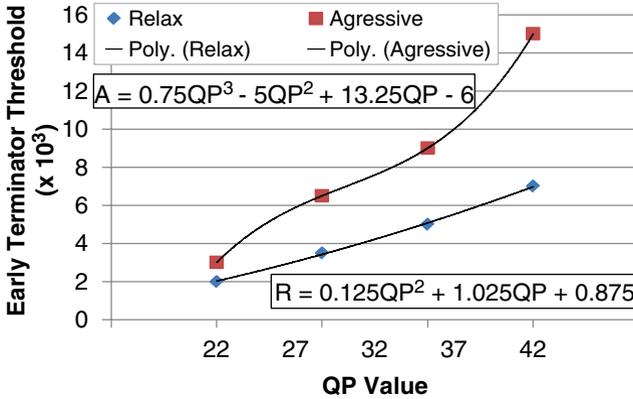


Fig. 4.20 Early termination threshold plots for Relax (*blue*) and Aggressive (*red*) complexity reduction

4.3.1.4 High-Confidence Modes and Low-Confidence Modes

The modes in the sorted *CandidateList* are partitioned into *high-confidence* and *low-confidence* modes using TH_{HighCL} . The threshold TH_{HighCL} is determined (using an exhaustive analysis) as 25 % of the total confidence level accumulated on the entire *CandidateList*. First, all of the *high-confidence* modes (i.e., $R_{MODE} \geq TH_{HighCL}$) are evaluated. Afterwards, the condition for early termination is evaluated. If the condition is not satisfied, all of the *low-confidence* modes (i.e., $R_{MODE} < TH_{HighCL}$) are evaluated. If the termination condition is not satisfied after evaluating the *low-confidence* modes, the mode decision proceeds to the next phase.

4.3.1.5 Video Properties-Based Mode Prediction

As discussed in Fig. 4.1, SKIP and Inter- 16×16 are the two most occurring modes in the non-anchor frames. In case sufficient correlation is not available in the 3D-Neighborhood, the variance property of a frame is considered to evaluate SKIP and Inter- 16×16 coding modes (in case these were not evaluated in the previous phases). The thresholds used in the conditions of this phase are derived using the PDFs presented in section “Analyzing the Video Properties” considering the region of high probability as discussed in Fig. 4.13.

4.3.1.6 Texture Direction-Based Mode Prediction

In the last phase a texture direction based prediction is employed to evaluate modes other than SKIP and Inter- 16×16 (if they were not tested in the previous phases). The direction of the gradient is considered to exclude improbable modes.

The RDCosts of Inter- 16×16 and Inter- 8×8 modes are compared to determine whether to evaluate bigger or smaller partitions for the current MB. If the RDCost of Inter- 16×16 is less than that of Inter- 8×8 , larger partitions (Inter- 16×8 or Inter- 8×16) are evaluated depending upon the dominant direction of the gradient. Otherwise, smaller partitions (i.e., Inter- 8×4 and Inter- 4×8) are evaluated accordingly Fig. 4.17.

In case of larger partition, SKIP mode is always tested (if not tested in the earlier phases). Similarly, Inter- 4×4 is always tested in case of smaller partitions (if not tested in the earlier phases). Finally, after all the processed phases, the best mode (i.e., the mode with the minimum RDCost) is used for coding the current MB.

4.3.2 Energy-Aware Complexity Adaptation

Besides the algorithms able to perform the fast mode decision, a complexity adaptation algorithm is required to adapt the mode decision at run time according to the changing application scenarios. Targeting MVC encoding systems where battery level, user constraints, and video content may vary widely along the time, we propose in this section an energy-aware complexity adaptation for MVC targeting mobile devices. Our algorithm, also presented in (Shafique et al. 2010b), employs several *Quality-Complexity Classes* (QCCs), such that each class evaluates a certain set of coding modes (thus a certain complexity requirement) and provides a certain video quality. To support asymmetric view quality, views for one eye are encoded with high-quality class and views for the other eye are encoded using a low-quality class. Our algorithm adapts the QCCs for different views at run time depending upon the current battery level.

4.3.2.1 Employing Asymmetric View Quality

A reduction in the complexity and energy consumption can be obtained by exploiting the binocular suppression theory which is based on the psycho-visual studies of stereoscopic vision (Stelmach and Tam 1999). According to this study, if the video quality of left and right eye views differs, the overall perceived quality is close to the high-quality sharper view (Stelmach and Tam 1999). However, for a blocky image, the perceived quality is the average of left and right eye views. Based on the binocular suppression theory (Stelmach and Tam 1999) and considering the in-loop deblocking filter of MVC (JVT 2008) (that reduces the blocking artifacts), views for two eyes can be encoded at different qualities (i.e., exploiting asymmetric view quality), thus requiring different computational efforts.

Figure 4.21 shows the MVC prediction structure for a four-view scenario employing asymmetric view quality. Assuming that the viewer is always exposed to adjacent views (S_n and S_{n+1}), the even views (S_0 and S_2) are encoded in higher quality while odd views (S_1 and S_3) are encoded in lower quality. In such way, the

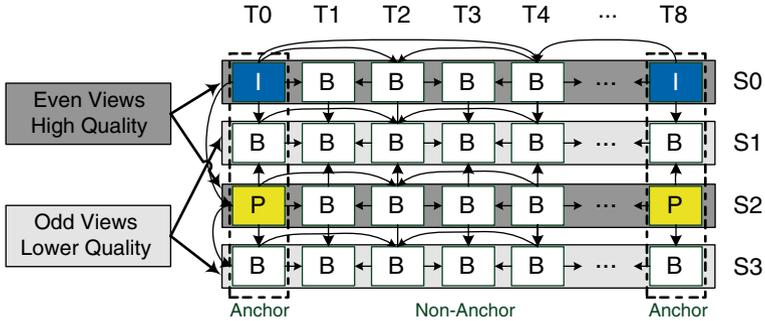


Fig. 4.21 MVC coding structure for asymmetric coding

viewer sees one high-quality and one low-quality view resulting in a perception near to the high-quality view. The use of high quality in even views is explained by the fact they are used as reference to odd views.

Although in (Stelmach and Tam 1999) the low-quality frames were synthetically blurred for analysis, this knowledge can be extended to a real scenario and applied in techniques to reduce the MVC coding complexity. In our scheme, the odd views will be submitted to more aggressive mode decision resulting in a lower quality in relation to their neighboring views.

In the following section is presented the energy-aware complexity adaptation algorithm besides the *Quality-Complexity Classes* (QCCs) and *Quality States* (QS) description.

4.3.2.2 QCCs: Quality-Complexity Classes

In order to employ the asymmetric view quality and the battery-level sensitivity to our scheme we define three *Quality-Complexity Classes* (QCCs).

QCC 1: MBs of *QCC 1* are exposed to the more aggressive mode decision of our scheme including SKIP and Inter 16×16 modes. Therefore, they have the lowest video quality and the higher complexity reduction.

QCC 2: This class presents the intermediate video quality and complexity reduction. Modes of *QCC 1* plus Intra 16×16 , Inter 16×8 , 8×16 , and 8×8 are evaluated.

QCC 3: More computationally intense class, and consequently, the one that provides better video quality. It includes the coding modes available in *QCC 1* and *QCC 2* plus small blocks such as Intra 4×4 , Inter 8×4 , 4×8 , and 4×4 .

Figure 4.22 presents the high-level diagram of our scheme showing the mode decision flow. The QCCs are related to three different prediction phases according to the dashed blocks in Fig. 4.22. *QCC 1* is subject to Phase 1; *QCC 2* to Phase 1 and Phase 2; and *QCC 3* is subject to Phase 1, Phase 2, and Phase 3. However, even

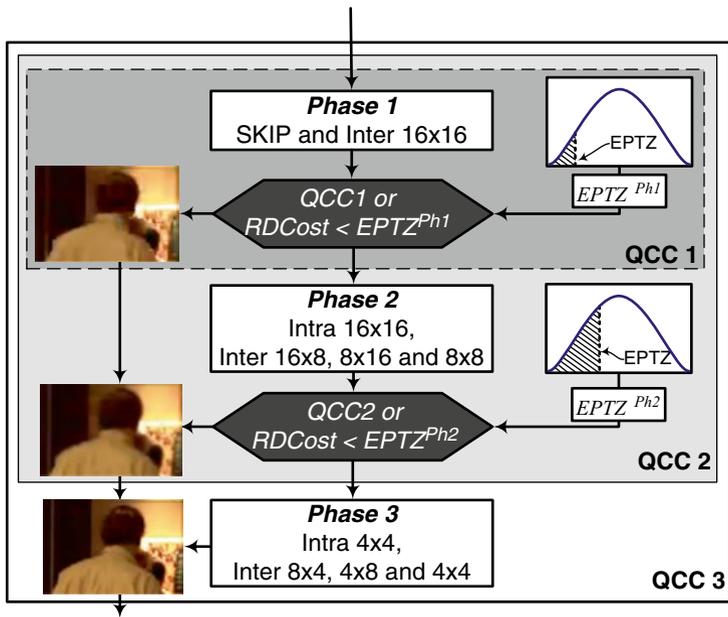


Fig. 4.22 Energy-aware MVC complexity adaptation scheme

for *QCC 2* and *QCC 3* big part of MBs are SKIP or Inter 16×16 and there is no need to test small block sizes. For this reason, the early prediction terminator zone (EPTZ) was defined.

4.3.2.3 Mode Decision Algorithm for Different QCCs

The proposed algorithm is presented in Fig. 4.23. It is composed of three phases. In Phase 1 the RDCost of SKIP and Inter 16×16 are calculated. If the current MB is *QCC 1* or the RDCost of one of the predicted modes is under the $EPTZ^{Ph1}$ limit the MD is terminated. Phase 2 calculates the RDCost for Intra 16×16 and one out of three inter modes, 16×8, 8×16, and 8×8, depending upon the gradient direction. The MD is terminated if the MB is *QCC 2* or the best RDCost is lower than $EPTZ^{Ph2}$ limit. For MBs *QCC 3* with RDCost higher than $EPTZ^{Ph2}$ one of four small block modes (Intra 4×4, Inter 8×4, 4×8, and 4×4) is tested. Finally, the best prediction mode, i.e., mode of lowest RDCost, is used to encode the current MB.

With the RDCost characterization and using the strategy described in Sect. 4.2, we defined the EPTZ as being $RDCost_{SKIP} < \mu_{RD} - \sigma_{RD}$ for *QCC 2*. In other words, if the best RDCost for previously tested modes is within EPTZ (see PDFs in Figs. 4.22 and 4.24) the mode decision is terminated. For *QCC 3* there are two early termination

```

MB Mode Decision (Current MB)
  //Phase 1
  01. Calculate RDCost (SKIP, Inter16x16);
  02. If ( Class 1 or RDCost < EPTZPh1 )
  03.   Exit;
  04. End If
  //Phase 2
  05. Calculate RDCost (SKIP, Intra16x16);
  06. If ( Gradient_Horiz > 1.25 * Gradient_Vert )
  07.   Calculate RDCost (Inter8x16);
  08. Else If ( Gradient_Horiz < 0.8 * Gradient_Vert )
  09.   Calculate RDCost (Inter16x8);
  10. Else
  11.   Calculate RDCost (Inter8x8);
  12. End If
  13. If ( Class 2 or RDCost < EPTZPh2 )
  14.   Exit;
  15. End If
  //Phase 3
  16. For ( all 8x8 partitions )
  17.   If ( Gradient_Horiz > 1.25 * Gradient_Vert )
  18.     Calculate RDCost (Inter4x8);
  19.   Else If ( Gradient_Horiz < 0.8 * Gradient_Vert )
  20.     Calculate RDCost (Inter 8x4);
  21.   Else
  22.     If ( Inter16x16 < Intra16x16 )
  23.       Calculate RDCost (Inter4x4);
  24.     Else
  25.       Calculate RDCost (Intra4x4);
  26.     End If
  27.   End If
  28. End For
  29. Encode MB (Mode with lowest RDCost);
  
```

Fig. 4.23 Pseudo-code of mode decision for different QCCs

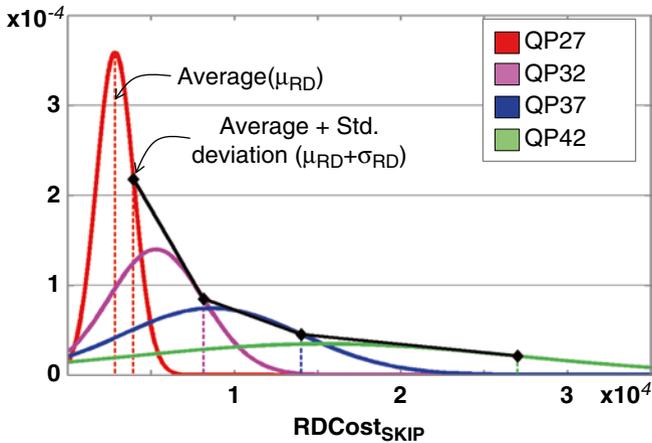
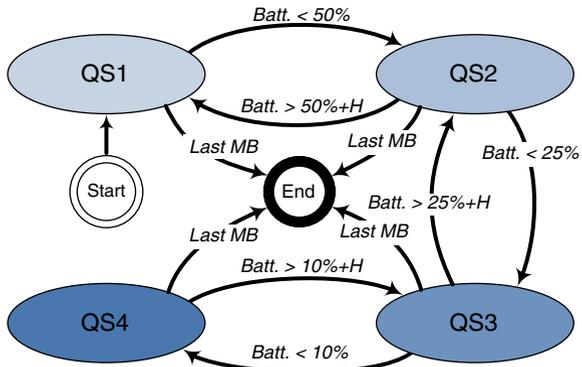


Fig. 4.24 Probability density function for RDCost

Table 4.2 Quality states

Quality state	Video quality	Even views	Odd views
<i>Quality State 1(QS1)</i>	High quality	<i>QCC3</i>	<i>QCC3</i>
<i>Quality State 2(QS2)</i>	Medium quality	<i>QCC3</i>	<i>QCC2</i>
<i>Quality State 3(QS3)</i>	Low quality	<i>QCC2</i>	<i>QCC1</i>
<i>Quality State 4(QS4)</i>	Lowest quality for battery saving	<i>QCC1</i>	<i>QCC1</i>

Fig. 4.25 Run-time complexity adaptation state machine

points, one at Phase 1 defined as $RDCost_{SKIP} < \mu_{RD} - 2 * \sigma_{RD}$ and other in Phase 2 $RDCost_{SKIP} < \mu_{RD} - \sigma_{RD}$. Once the distribution is different for each QP, the EPTZ limit approximated by polynomial curve fitting is given by the following QP-based Eqs. (4.15), (4.16), and (4.17):

$$\left(EMT^{Ph1}\right)_{QCC2} = \left(EMT^{Ph2}\right)_{QCC3} = RDCost_{AVG} - 1.5RDCost_{SD}, \quad (4.15)$$

$$\left(EMT^{Ph1}\right)_{QCC3} = RDCost_{AVG} - 0.5RDCost_{SD}, \quad (4.16)$$

$$\left(EMT^{Ph1}\right)_{QCC2} = \left(EMT^{Ph2}\right)_{QCC3} = 29.663QP^2 - 1409.1QP + 18766. \quad (4.17)$$

4.3.2.4 Energy-Aware Complexity Adaptation Algorithm

Associated with the QCCs, our scheme employs four different *Quality States* (QS). The QSs consider the binocular suppression theory (Stelmach and Tam 1999) using asymmetric view quality and react, at run time, to the changing battery level. As summarized in Table 4.2, *QS1* presents the highest quality and encodes all views as *QCC3*. In turn, *QS2* and *QS3* use the view quality asymmetry encoding odd view in lower quality than even views. *QS4* provides the lowest quality and highest complexity reduction coding all views as *QCC1*.

The QS control is performed by one state machine that receives an indication of the battery level as input. Figure 4.25 presents the transitions between the four

possible states. The quality states just change to the immediately superior or inferior quality in order to have smooth video quality variation. The hysteresis (H) is fixed as 5 % in order to avoid quick oscillations between different states, and consequently, video quality fluctuations. This state machine can be easily adapted to consider other external parameters such as user presets and time constraints.

4.3.3 Multilevel Fast Mode Results

The detailed results of our multilevel fast mode decision algorithm compared to the RDO-MD solution implemented in the JMVC are presented along this section. Table 4.3 presents the results for Δ PSNR, Δ Bitrate, and complexity reduction (i.e., time saving, TS). For JMVC using the exhaustive RDO-MD the results are presented in coding time (column T, in seconds), PSNR (dB), and Bitrate (column BR, in kbps). The values for a certain QP value are obtained by averaging over all eight views. The last row named *Average* presents the average results over all sequence. The experiments were performed for eight views considering IPB view coding order. For more details on the experimental setup refer to Sect. 6.1.

Figure 4.26 illustrates the PSNR (lines) and time savings (bars) comparison of *Relax* and *Aggressive* levels averaged over all views and QPs for *Ballroom* and *Exit* sequences. It is noted that the difference between the RD curves of *Relax* and *Aggressive* is more significant at low bitrates and this difference diminishes at higher bitrates. The time savings of the *Aggressive* level are significantly higher compared to *Relax* level at higher bitrates while providing slight RD difference. *Relax* scheme was developed to keep video quality in all QP ranges and for this reason is forced to reduce TS for big and small QP ranges presenting higher TS for intermediate QPs. In *Aggressive* scheme the higher TS is prioritized for the whole QP range.

4.3.3.1 View-Level Time Saving Evaluation

A view-wise Δ PSNR and time savings comparison of *Relax* and *Aggressive* levels is provided in Fig. 4.27 for the *Exit* sequence encoded using QP=32. Odd views—with north and south views (i.e., Views 1, 3, 5) available in the neighborhood—present higher time savings compared to the views with just one (i.e., Views 2/4/6/7) or none available neighboring views (i.e., View 0). Additionally, Views 1, 3, and 5 also present a smaller PSNR loss. This higher complexity reduction and reduced PSNR loss is due to the larger correlation space in the 3D-Neighborhood. It implies that more neighboring MBs are available for the prediction. Consequently, a more accurate *CandidateList* is generated.

4.3.3.2 Tested Modes Evaluation

The high time savings provided by the multilevel mode decision comes from the reduced number of coding modes tested. Figure 4.28 provides the result of the

Table 4.3 Detailed results for Δ PSNR, Δ Bitrate, and time savings compared to the exhaustive RDO-MD

Video sequence	QP	JMVC			Proposed relax			Proposed aggressive		
		T [s]	PSNR [dB]	BR [kbps]	TS [%]	Δ PSNR [dB]	Δ BR [%]	TS [%]	Δ PSNR [dB]	Δ BR [%]
Ballroom	22	2,682.53	41.111	3,176.849	54.77	0.005	2.923	59.03	0.002	5.820
	27	2,490.47	38.415	1,319.736	61.23	0.039	3.015	70.12	0.038	10.150
	32	2,315.22	35.667	654.338	57.04	0.039	0.934	65.71	0.084	3.060
	37	2,121.62	32.884	360.162	52.67	0.025	0.453	63.07	0.065	1.500
Exit	22	2,671.07	41.601	2,114.453	60.29	0.006	3.937	67.68	0.045	7.510
	27	2,268.79	39.456	652.491	71.36	0.016	2.402	80.08	0.099	12.090
	32	2,065.18	37.508	292.673	70.19	0.030	1.101	78.10	0.109	5.910
	37	1,900.21	35.293	163.436	67.92	0.043	0.357	78.37	0.123	2.060
Vassar	22	2,963.66	40.743	3,007.434	55.26	0.001	2.837	69.35	0.001	5.470
	27	2,519.33	37.828	850.324	77.18	0.021	1.198	81.24	0.001	11.250
	32	2,114.88	35.490	259.826	76.59	0.020	-0.028	82.44	0.055	3.450
	37	1,827.57	33.294	117.428	74.69	0.008	-0.196	82.23	0.034	3.660
Race1	22	4,908.26	42.340	2,549.767	64.55	0.006	8.349	78.10	0.005	11.550
	27	4,631.98	39.422	1,182.855	71.70	0.036	6.514	80.09	0.028	17.800
	32	4,269.55	36.501	552.949	70.49	0.036	3.443	78.13	0.064	9.800
	37	3,806.22	33.795	294.763	69.05	0.028	1.543	74.92	0.074	6.650
Rena	22	2,238.63	46.555	1,347.801	68.54	-0.205	12.283	67.09	-0.212	22.960
	27	1,960.61	43.846	587.514	70.55	-0.215	14.289	70.44	-0.310	36.550
	32	1,685.64	40.535	293.333	71.79	0.028	6.971	70.87	-0.220	33.740
	37	1,452.22	37.396	163.581	66.58	0.043	3.038	73.03	-0.154	26.880
Akko&Kayo	22	2,644.20	43.53	1,743.05	65.67	-0.056	10.152	66.81	-0.050	14.770
	27	2,560.85	40.79	808.48	69.66	-0.015	9.395	74.27	-0.020	21.920
	32	2,466.77	37.59	433.65	65.66	0.036	3.130	71.05	0.000	15.280
	37	2,320.73	34.45	254.24	59.72	0.023	1.597	70.02	0.020	8.970
Breakdancers	22	5,893.46	41.449	4,899.089	53.81	0.002	5.393	62.39	0.004	7.150
	27	4,817.77	39.841	1,454.553	63.14	0.038	7.492	76.10	0.061	12.910
	32	4,116.45	38.432	667.955	62.98	0.054	7.861	74.95	0.111	11.700
	37	3,487.59	36.629	378.932	59.53	0.094	3.615	76.11	0.237	7.150
Uli	22	4,826.40	40.476	8,152.865	44.35	-0.001	2.024	47.53	0.007	3.545
	27	4,638.51	38.591	3,801.339	61.34	0.013	3.245	66.62	0.048	5.488
	32	4,326.23	36.238	2,056.013	57.22	0.002	2.115	61.41	0.037	4.280
	37	3,945.21	33.554	1,162.239	61.54	0.078	3.676	68.11	0.137	7.150
Poznan_Hall2	22	10,737.88	42.9111	5,149.584	63.48	-0.003	5.649	67.68	0.045	7.510
	27	8,911.15	41.693	1,417.539	69.02	0.016	5.124	80.08	0.099	12.090
	32	7,778.938	40.303	693.789	65.46	0.020	3.913	78.10	0.109	5.910
	37	6,702.69	38.606	420.110	61.56	0.047	1.269	78.37	0.123	2.060
GT_Fly	22	6,334.83	41.247	6,437.935	55.16	0.017	9.547	64.68	0.028	12.975
	27	5,168.38	39.705	2,029.539	59.98	0.042	8.832	72.68	0.052	10.685
	32	4,511.77	38.280	946.270	62.40	0.044	6.562	67.60	0.079	8.655
	37	2,185.58	36.819	610.563	59.74	0.051	3.217	67.00	0.116	9.984
Average	22	3,603.53	42.226	3,373.914	58.59	-0.023	6.310	65.03	-0.013	9.926
	27	3,236.04	39.774	1,332.162	67.52	-0.001	6.151	75.17	0.010	15.093
	32	2,919.99	37.245	651.342	65.98	0.031	3.600	72.84	0.043	10.178
	37	2,607.67	34.661	361.847	63.30	0.044	1.857	73.12	0.077	7.606
	AVG	3,091.81	37.183	1,172.794	63.85	0.013	4.479	71.54	0.029	10.701

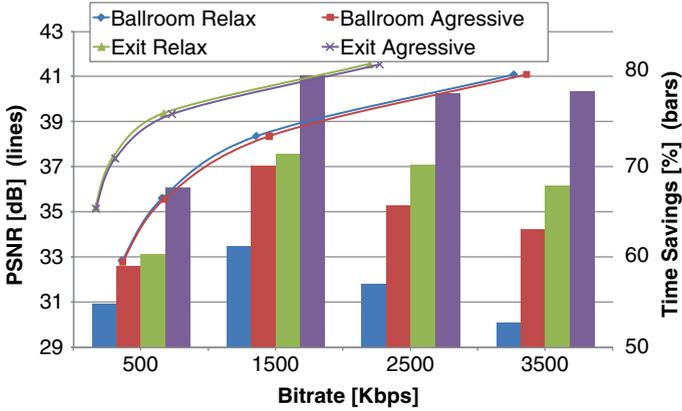


Fig. 4.26 Average tested modes (QP = {22,27,32,37,42}, GOP = 8, Views = 8)

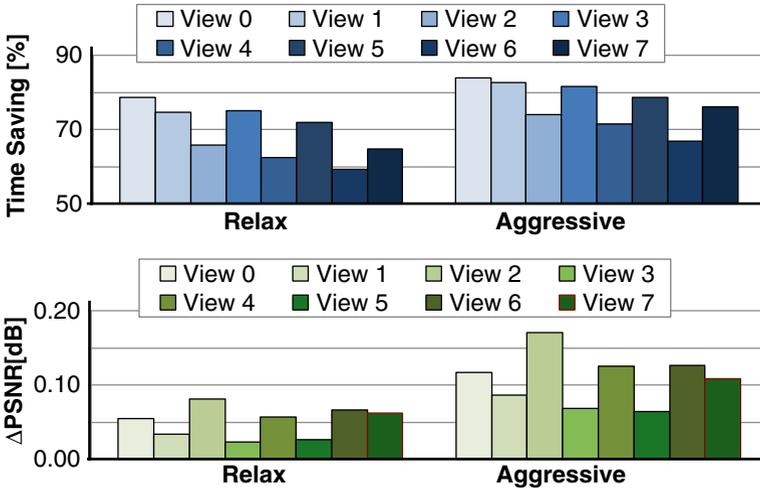
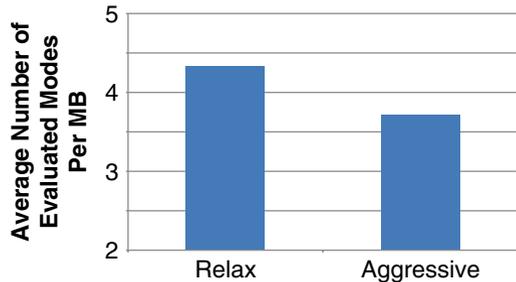


Fig. 4.27 View-level time savings and Δ PSNR comparison of Relax and Aggressive levels (Exit sequence, QP = 32)

Fig. 4.28 Average tested modes for all sequences



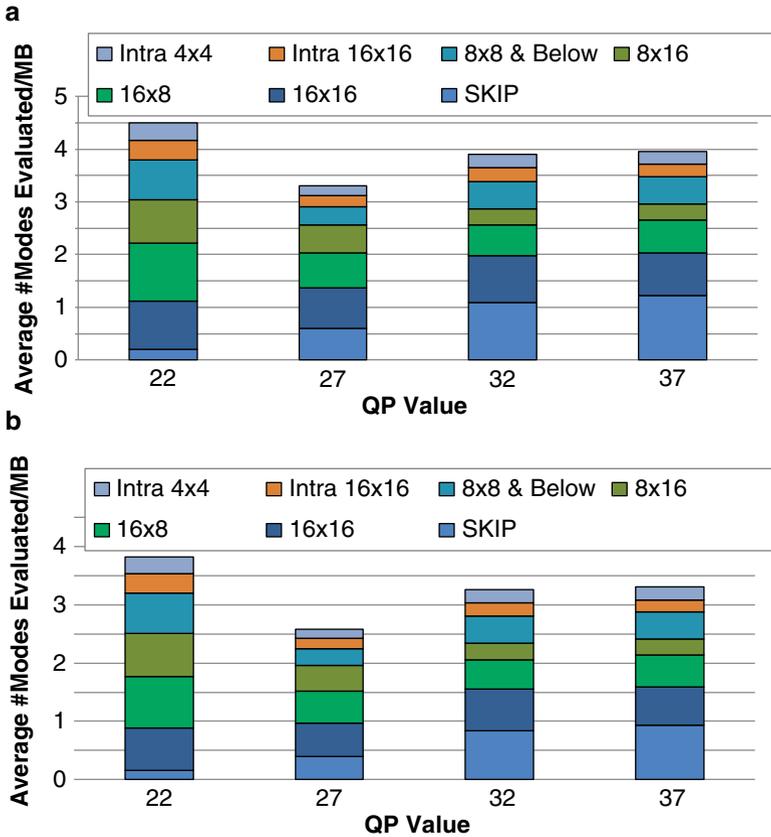


Fig. 4.29 Detailed number of evaluated modes for (a) *Relax* and (b) *Aggressive* (Exit Sequence)

average number of modes evaluated per MB considering the different operation modes. The *Relax* and *Aggressive* levels of our complexity reduction scheme process only 3.7 and 2.3 modes per MB, respectively.

The distribution of the evaluated modes for *Relax* and *Aggressive* complexity reduction levels is presented in Fig. 4.29. It is noted that the number of SKIP mode increases for higher QPs while the number of other modes decreases accordingly. This behavior confirms the analysis of optimal mode distribution discussed in section “Coding Mode Distribution Analysis”. For QP 32 and above, the number of evaluated modes increases to maintain a high video quality.

4.3.3.3 Frame-Level Time Saving Evaluation

To analyze the frame-wise comparison of *Relax* and *Aggressive* levels, we have plotted the PSNR and time savings for View 0, 1, and 2 of *Exit* test sequence encoded using QP=32, as shown in Figs. 4.30 and 4.31. Please note that the plots only

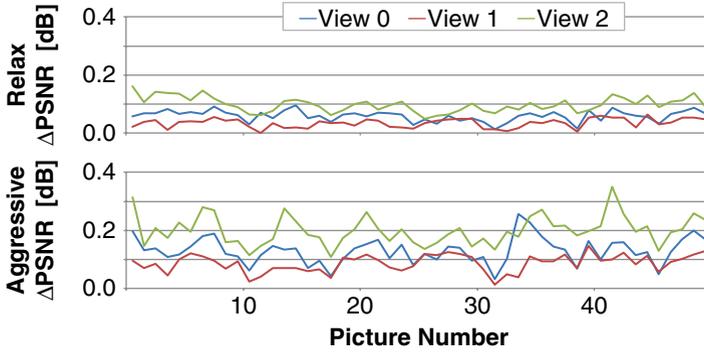


Fig. 4.30 Frame-wise PSNR loss comparison of Relax and Aggressive levels (Exit, QP=32)

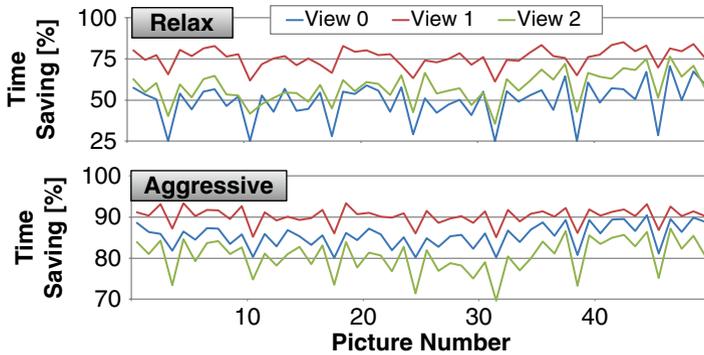


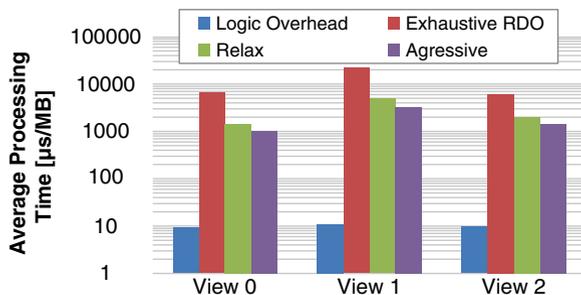
Fig. 4.31 Frame-wise time saving comparison of Relax and Aggressive levels (Exit, QP=32)

contain results for non-anchor frames, which are the primary focus for complexity reduction in this algorithm.

There are 0, 1, and 2 available neighboring views available for the View 0, View 1, and View 2, respectively (representing all possible cases). View 2 exhibits a higher Δ PSNR and lower time savings, while View 1 exhibits higher time savings and a lower Δ PSNR for most of the frames when compared to the other plotted views. This ratifies the view-level results from Fig. 4.27.

The sudden variations (i.e., valleys) in Fig. 4.31 correspond to the frames in the middle of the GOPs, i.e., frames that have temporal-neighbors from the anchor frames. In this case, more intra modes are evaluated (in phases 3 and 4) in addition to the inter modes leading to a lower complexity reduction. View 1 due to the availability of all view-neighbors suffers less with such variations.

Fig. 4.32 Overhead of our scheme



4.3.3.4 Multilevel Mode Decision Algorithm Overhead

The overhead of our complexity reduction scheme is already computed in the total processing time and time savings. Figure 4.32 compares the average overhead of the computational logic of our scheme with the average processing time of one MB encoded using different schemes. It is noted that the overhead is 0.15 % of the average MB encoding time using the exhaustive RDO-MD. Figure 4.32 shows that the overhead of our scheme is insignificant compared to its time savings.

In this section was presented the multilevel fast mode decision algorithm focusing on complexity reduction MVC that exploits the image properties, RDCost, and the correlation in the 3D-Neighborhood to provide complexity reduction with insignificant PSNR loss. Our detailed analysis provides the foundation for the proposed scheme. In order to react to the changing QP values, QP-based threshold equations are deployed.

For a trade-off between the desired complexity reduction and the resulting quality loss, two different operational levels are proposed for our scheme: the *Relax* and *Aggressive* modes. However, to better exploit the complexity reduction vs. RD performance, a control algorithm able to select at run time the most appropriate complexity reduction level is desirable. In the following section an energy-aware complexity adaptation based on fast mode decision is proposed.

4.3.4 Energy-Aware Complexity Adaptation Results

This section presents the detailed experimental results for each *Quality State* of the proposed energy-aware complexity adaptation algorithm compared to the RDO-MD. The overall results for the complexity adaptation algorithm are presented in section “Comparing the Energy-Aware Complexity Adaptation to the State-of-the-Art” in Chap. 6. The experiments used the experimental setup described in Sect. 6.1.

Table 4.4 presents the detailed PSNR, bitrate (BR), and time saving (TS) results of the four *Quality States* of our scheme compared to the exhaustive mode decision (RDO-MD). For the *QS1* state our scheme provides a TS of up to 77 % with negligible PSNR loss (avg. 0.089 dB). The TS goes up to 87 % for *QS4* with an average PSNR loss of 0.195 dB.

Table 4.4 Comparison between the Quality States (QS)

QP	Ballroom			Exit			Vassar			Rena			QP average			Total average			
	TS	Δ PSNR	Δ BR	TS	Δ PSNR	Δ BR	TS	Δ PSNR	Δ BR	TS	Δ PSNR	Δ BR	TS	Δ PSNR	Δ BR	TS	Δ PSNR	Δ BR	
<i>QS1</i>	22	75.71	0.095	-3.45	77.28	0.093	-4.13	75.21	0.084	-7.143	82.26	0.292	-3.221	77.62	0.14	-4.49	75.29	0.089	-5.48
	27	75.13	0.074	-5.33	77.28	0.072	-4.75	73.84	0.031	-8.527	79.64	0.207	-2.186	76.47	0.10	-5.20			
	32	73.15	0.075	-2.36	76.45	0.070	-5.70	71.97	0.015	-12.22	76.18	0.094	-1.493	74.44	0.06	-5.44			
	37	71.62	0.031	-1.24	75.95	0.003	-6.05	70.80	0.115	-15.49	72.24	0.073	-4.429	72.65	0.06	-6.80			
<i>QS2</i>	22	77.96	0.096	0.70	79.29	0.094	0.07	77.68	0.086	-2.344	84.18	0.291	0.305	79.78	0.14	-0.32	76.96	0.093	-0.53
	27	76.97	0.079	1.84	78.60	0.072	0.58	75.18	0.036	-2.381	81.21	0.210	1.193	77.99	0.10	0.31			
	32	75.22	0.089	2.01	78.08	0.077	0.69	73.49	0.023	-2.294	77.26	0.099	0.707	76.01	0.07	0.28			
	37	73.27	0.053	1.54	77.22	0.027	0.07	72.51	0.091	-8.604	73.30	0.063	-2.573	74.07	0.06	-2.39			
<i>QS3</i>	22	84.86	0.116	4.38	84.52	0.112	4.32	85.81	0.097	1.676	85.67	0.280	4.260	85.22	0.15	3.66	82.64	0.123	4.76
	27	83.98	0.088	6.42	83.43	0.084	7.45	83.61	0.050	2.660	83.31	0.228	3.919	83.58	0.11	5.11			
	32	82.86	0.124	6.14	82.73	0.112	8.20	81.57	0.056	3.974	80.23	0.186	3.536	81.85	0.12	5.46			
	37	81.43	0.131	5.48	81.53	0.113	7.83	80.03	0.059	3.790	76.66	0.127	2.116	79.91	0.11	4.80			
<i>QS4</i>	22	87.96	0.148	6.42	86.93	0.140	6.75	87.86	0.123	2.456	87.54	0.318	6.106	87.57	0.18	5.43	85.26	0.195	7.40
	27	87.29	0.112	8.95	85.40	0.117	11.39	85.59	0.065	3.690	85.83	0.320	6.175	86.03	0.15	7.55			
	32	86.42	0.173	9.24	84.57	0.205	13.08	83.58	0.086	5.536	83.53	0.364	6.530	84.53	0.21	8.60			
	37	85.44	0.239	8.96	83.32	0.264	13.02	82.20	0.090	5.277	80.74	0.356	4.782	82.93	0.24	8.01			

To calculate the objective quality of a sequence we consider the average PSNR between all possible stereo view points (VP) of a sequence. For example, a sequence with four views has three stereo VPs (View 0 and View 1, View 1 and View 2, View 2 and View 3). To calculate the PSNR of a given VP considering the binocular suppression we use the Eq. (4.18), as proposed in (Ozbek et al. 2007):

$$PSNR^{VP} = (1 - \alpha) \cdot PSNR^{HighQuality} + \alpha \cdot PSNR^{LowQuality}; \alpha = 1/3. \quad (4.18)$$

4.4 Fast Motion and Disparity Estimation

According to the motivational analysis presented in Sect. 3.1.1 and challenges discussed on Sect. 3.2, the two main sources of complexity and energy consumption in the MVC encoder are the mode decision and the motion and disparity estimation units. Along Sect. 4.3 distinct solutions for reducing the complexity and energy for the MD were proposed. Moreover, an energy-aware complexity adaptation based on mode decision was presented in order to enable run-time adaptivity to changing system and content scenarios. In this section the target is to present solutions to reduce the complexity and energy consumption associated with the second main complexity source, the ME/DE unit.

In this section is presented a correlation analysis related to motion and disparity vectors (MV, DV) followed by a Fast ME/DE algorithm. Our Fast ME/DE algorithm was designed taking into account a future hardware implementation.

4.4.1 Fast Motion and Disparity Estimation Algorithm

Our Fast ME/DE scheme (Zatt et al. 2011c) is based on the previously presented 3D-Neighborhood analysis. However, to exploit this correlation the motion and disparity fields must be available. In order to establish these fields at least one frame using DE and one using ME must be encoded with the optimal or a near-optimal searching algorithm. In our scheme, to avoid major quality loss, all anchor frames and the frames situated in the middle of the GOP are encoded using the TZ search algorithm [the fast ME/DE algorithm used in JMVC (JVT 2009a)]. The anchor frames are encoded using DE, while the frames in the middle of a GOP use ME or ME and DE according to the view they belong. These frames encoded with high effort are herein referred as Key Frames (KF), while the others are the Non/key Frames (NKF). Once the motion and disparity fields are available all NKF can be encoded based on these fields. The complete ME/DE search pattern is skipped for all NKF. It only uses the predictors inferred from the 3D-Neighborhood.

Figure 4.33 presents the flow diagram of our proposed fast ME/DE scheme based on the 3D-Neighborhood vectors correlation. It employs two different prediction classes: Ultra Fast Prediction and Fast Prediction. The scheme is composed of three

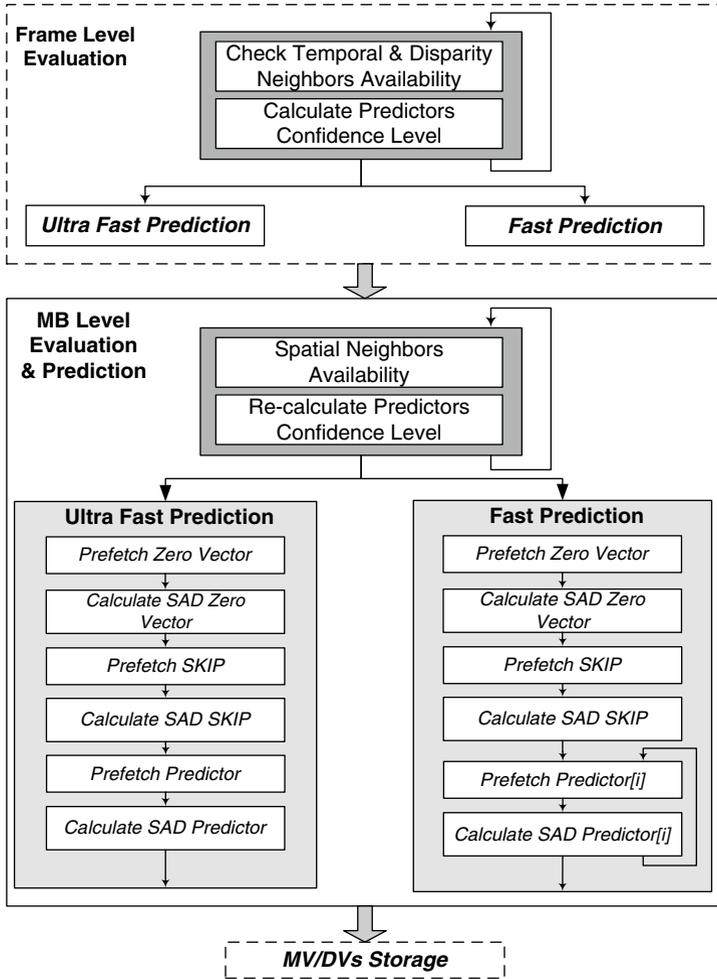


Fig. 4.33 Flow diagram of the adaptive fast ME/DE

main phases (1) Frame-Level MV/DVs evaluation, (2) MB-Level MV/DVs Evaluation and Prediction, and (3) MV/DV Storage. Figure 4.33 considers only the NKF coding. For KF the TZ search algorithm is used (as discussed earlier) to provide a good motion/disparity field.

In Frame-Level MV/DVs evaluation, the presence of all temporal-disparity predictors is checked. If available, they are read from a MV/DV memory. The spatial predictor is not loaded in this phase since it is not available due to the spatial dependencies, if a hardware architecture scenario is considered. With the available data, the current MB is pre-classified in one out of the two prediction classes according to the predictors' Confidence Level. Note, the spatial predictor is required for the SKIP vector calculation. So, this predictor is also considered in our algorithm to classify the current MB in the MB-Level evaluation and prediction phases.

Table 4.5 Comparison of our fast ME/DE algorithm to TZ Search

Video	QP	<i>TZ Search</i>			Fast ME/DE		
		Time [s]	BR [kbps]	PSNR [dB]	TS [%]	Δ BR [%]	Δ PSNR [dB]
Ballroom	22	215.1	3,298.026	40.709	85.9	8.4	0.011
	32	175.7	651.640	35.119	86.2	11.7	0.060
	42	127.1	188.178	29.318	84.9	19.8	0.190
Vassar	22	171.1	3,415.744	40.456	83.0	1.1	0.010
	32	110.0	315.142	35.226	82.4	6.7	0.013
	42	72.1	64.888	30.563	79.7	6.7	0.043
Breakdancers	22	583.5	5,680.204	41.172	86.0	15.7	0.087
	32	384.6	788.800	38.010	85.1	14.6	0.275
	42	262.9	277.970	33.803	82.8	8.9	0.304
Uli	22	600.9	12,245.676	39.819	82.9	9.3	0.053
	32	516.8	2,949.870	34.960	83.0	13.7	0.134
	42	406.5	849.462	28.944	82.4	8.5	0.191
Poznan_Hall2	22	1,206.4	6,245.748	42.654	86.8	8.2	0.012
	32	811.2	1,002.354	40.138	86.1	10.3	0.017
	42	726.8	521.816	35.465	84.5	7.7	0.027
GT_Fly	22	1,321.8	7,123.971	40.980	81.3	10.0	0.103
	32	954.8	1,113.591	38.057	81.9	17.3	0.265
	42	801.0	613.548	33.849	79.8	12.3	0.298
Average	22	683.1	6,334.895	40.965	84.3	8.8	0.046
	32	492.2	1,136.900	36.918	84.1	12.4	0.127
	42	399.4	419.310	31.990	82.4	10.7	0.176
	Avg.	524.9	2,630.368	36.625	83.6	10.6	0.116

The predictors Confidence Level is calculated based on the offline *hit* value, as presented in Table 4.1. Each predictor is associated with a Confidence Level (*hit* value). If one predictor has a Confidence level higher than a threshold ($CL_{pred} > CL_{TH}$), the current MB is classified to be encoded as Ultra Fast Prediction. Otherwise, the MB is classified to be encoded with the Fast Prediction. In case of Ultra Fast Prediction MBs, only three vectors are tested: the predictor with highest Confidence Level (also referred as Common Vector), the Zero vector, and the SKIP vector. The Zero and SKIP vectors are tested because of their high occurrence. Fast Prediction MBs test all available predictors in addition to the Zero and SKIP vectors. It is important to mention that even if all predictors are available and different (this worst case rarely occurs), only 13 predictors are tested.

4.4.2 Fast ME/DE Algorithm Results

In Table 4.5 the fast ME/DE results are detailed for the four evaluated sequences considering three different QPs (22,32,42). The TZ Search with a search range of $[\pm 64, \pm 64]$ is used for comparison as it is used for the Key frames and performs 23 \times faster compared to the Full Search (not used for performance comparison),

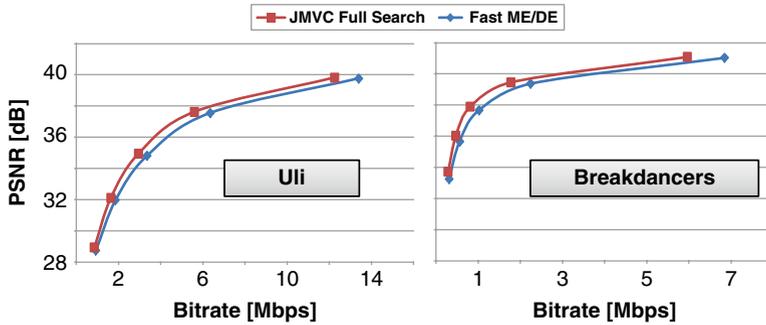


Fig. 4.34 Rate-distortion comparison with full search

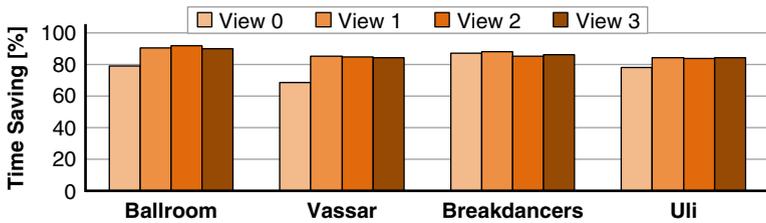


Fig. 4.35 View-level execution time savings compared to TZ Search

while providing the similar rate-distortion results (Yang 2009). Compared to the TZ Search, our fast ME/DE provides 83 % execution time saving at the cost of 11 % increase in bitrate and 0.114 dB of PSNR loss. In the best case, the execution time savings go up to 86 %, which represents a significant complexity reduction. Moreover, the reduced number of candidate blocks leads to a lower number of external memory accesses.

Figure 4.34 presents the RD curves for two XGA (1024×768) video sequences *Uli* and *Breakdancers*. It is noted that the RD curves of our fast ME/DE algorithm are close to that of the Full Search (used for quality comparison only). Compared to the Full Search, our fast ME/DE algorithm suffers from an insignificant loss of 0.116 dB (on average).

The view-level execution time savings are presented in Fig. 4.35. Note that for all views (except for View 0 of *Vassar* sequence) the time saving is $\geq 80\%$. The execution time savings for the high-motion sequences are slightly more than that in the low-motion sequences. Figure 4.36 presents the average number of SAD operations for ME/DE of one MB using Full Search, TZ Search, and our fast ME/DE algorithm. Averagely, the proposed scheme reduces more than 99.9 % in comparison to Full Search and 86 % to TZ Search. Note, the detailed results in Figs. 4.35 and 4.36 are for QP32.

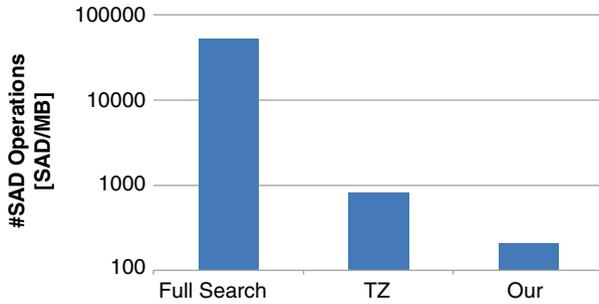


Fig. 4.36 Comparison of the number of SAD operation

4.5 Video-Quality Management for Energy-Efficient Algorithms

Although the energy and complexity reduction algorithms presented along this chapter were carefully designed to reduce undesirable effects to the coding efficiency, they lead to some level of quality drop due to the heuristics and simplifications inserted in the encoding process. For this reason, an algorithm to manage the coding process and compensate eventual video-quality losses is required. This management algorithm, however, must also consider and optimize the video quality and bitrate trade-off in order to increase the coding efficiency while respecting bandwidth constraints as discussed in the following.

Despite the high coding efficiency provided by MVC, the transmission and storage of 3D videos remain a big challenge, especially for services operating over bandwidth/buffer-constrained infrastructures. It becomes even more challenging due to changing input video properties, run-time variations on video encoder state, battery level, and user preferences. Thus, to provide high video quality while meeting channel bandwidth/buffering constraints it is necessary to further optimize the bandwidth usage by intelligently regulating the bits allocation. Therefore, a rate control algorithm is implemented to dynamically find a good compromise between the coding efficiency and video quality by adapting the QP.

In this section is presented the Hierarchical Rate Control (HRC) (Vizzotto et al. 2012) for MVC that employs coupled Model Predictive Control-based frame-level RC and Markov Decision Process-based BU-level RC. Before presenting the HRC, however, a bitrate allocation study within the 3D-Neighborhood is detailed.

4.5.1 Hierarchical Rate Control for MVC

In this section is presented the proposed Hierarchical Rate Control (HRC) for MVC, depicted in Fig. 4.37. The HRC is responsible for controlling the encoder output bitrate, in accordance with the user preferences and/or channel limitations, by

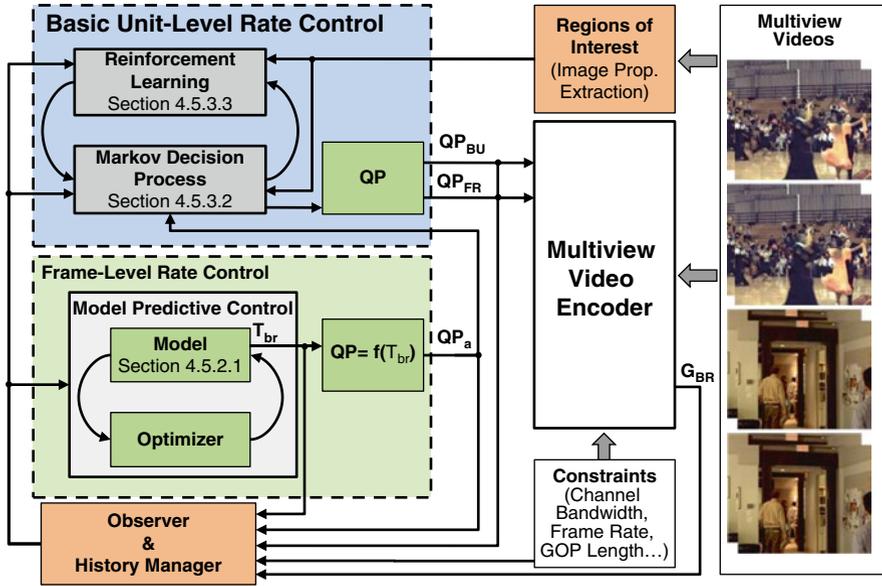


Fig. 4.37 Hierarchical rate control system diagram

monitoring the MVC encoder and actuating through QP adaptation. It can be conceptually divided in two actuation levels (1) frame level (that encapsulates GOP and frame levels) at coarse grain and the (2) basic unit level at fine grain. The MVC encoder receives the video sequences as input along with all user preferences and configurations to start the encoding process. The Model Predictive Control-based frame-level RC models the system behavior considering the encoding hierarchy and predicts the bitrate allocation at frame level considering temporal, view, and GOP-phase (inter-GOP) correlation. It defines the optimal QP for the predicted frames, the base QP, and forward it to the Markov Decision Process-based basic unit-level RC. At BU level, a fine-grained decision is taken to define the QP variation considering the image properties in terms of regions of interest. The fine-grained adaptation promotes an increase in objective and subjective video qualities inside the frame by allocating more bits to the RoI (in our case the hard-to-predict regions; see section “Region of Interest” in Chap. 2 and Sect. 4.1.3). The decision maker considers the previous knowledge, by implementing the Reinforcement Learning (RL) method, to increase or decrease the QP in relation to the base QP. To couple the frame and BU level in HRC, the RL unit feedbacks both the MPC and the MDP to keep system consistency and avoid mismatches. The HRC employs an observer unit able to read, store, and manage the MVC encoder feedback (generated bitrate) and variables that define the encoder system state (target bitrate, QP, input constraints, etc.) in order to support the bitrate prediction and actions/decision taking. Also, an image properties extractor is employed to build the RoI map used for BU-level RC. This integration allows HRC to properly exploit the influence of spatial, temporal, view, and GOP-phase inputs to define a global optimal control action.

MPC-based frame-level Rate Control: It is responsible for predicting the bitrate allocation and defining an optimal QP value for the current frame while minimizing a performance cost function. Our MPC-based RC deals with multiple stimuli superposition building the input horizon using previously encoded frames from temporal and view neighborhood. The proposed scheme also incorporates the GOP-phase for accurate bitrate prediction.

MDP-based Basic Unit-level Rate Control: The BU-level RC receives the QP defined at frame level and adjusts the QP for each BU. The proposed Markov Decision Process-based RC takes the decisions over a map of states based on a set of possible actions (QP adaptations) and the associated rewards. The texture-based map of states is linked to the map of RoI and provides the structure to make decision.

Coupled Reinforcement Learning: It is responsible for adapting MPC and MDP models to the dynamic system behavior. After an action is taken at BU level, the RL reads the system response, and updates the transition probabilities and the associated rewards in the MDP model. Once the frame is totally encoded, the resulting map of states is used to update the frame-level MPC. This strategy integrates frame level and BU level guaranteeing consistency and avoiding modeling mismatches.

On the following subsections the Hierarchical Rate Control will be presented in details along with the equations that describe the whole controller behavior. For simplicity we provide, in Table 4.6, the definitions of the main variables used in the HRC description.

4.5.2 Frame-Level Rate Control

The frame-level MVC Rate Control problem matches the control-theory superposition principle (Tatjewski 2010) defined as the response at a given place and time of the linear system caused by multiple stimuli. Model Predictive Control (MPC) techniques (García Carlos et al. 1989; Morari and Lee 1999) have demonstrated to accurately predict the response of multiple stimuli dynamic systems such as MVC encoder while incorporating the phase concept (periodic behavior) present in GGOP-level RC (see Sect. 4.1.3). MPC outperforms traditional feedback controllers by efficiently integrating input stimuli to state space constraints while providing flexibility by employing rolling input and output horizons (see section “Model Predictive Control” in Chap. 2).

As discussed in section “Model Predictive Control” in Chap. 2, the main goal of a Model Predictive Controller is to predict the future behavior of a system state and/or output over a finite time horizon as well as compute the future input signals at each step. These actions occur by minimizing a cost function under inequality constraints on the manipulated control or the controlled variables. In this work the MPC operates at frame level predicting the bitrate and providing the QP for each frame to be encoded. The rate controller tries to define a sequence of actions and then induce

Table 4.6 Variables definitions

Variable	Description
Frame-level rate control	
T_{BR}	Target bitrate for one frame (bits per frame)
BW	Channel bandwidth (bits per second)
FR	Frame rate (frames per second)
BA	Bit allocation (absolute)
w_I, w_P, w_B	I, P and B weight respectively (absolute)
\bar{w}_{GOP}	Average w for the current GOP (absolute)
L_{GOP}	GOP length (# of frames)
ω	Frame weight (absolute)
N_A	Number of anchor frames (# of frames)
BR	Bitrate (#bits)
H_{QP}	QP history (absolute)
QP_{FL}	Quantization parameter at Frame-level RC (discrete)
QP_{CLP}	Quantization parameter in last process (discrete)
QP_{st}	Initial quantization parameter (discrete)
Q	Quantization parameter in the optimization loop (discrete)
N_{FR}	Number of frames
Basic unit-level rate control	
M_S	RoI-Normalized variance matrix (absolute 0–1)
$M(\delta)$	MDP reward matrix (matrix of absolute RD)
BU	BU variance
μ	Average of BU_i
N_{BU}	Number of BUs
QP_{BU}	Quantization Parameter at Frame-level RC (discrete)
T_{BR}	Target bitrate for one frame (bits per frame)
R_S	BU Reward “Shared” (absolute)
R_L	Reinforcement learning Value (vector of H_R)
$f(s, \delta)$	Probability of state transition
P_R	Probability results from R_L vector of “phase” actions. Actions of R_L in a range of at least 2 horizons
$\Delta\delta$	Variation between actual BU δ and the δ of anchor frame
Mf	Variation of variance matrix values
HR	History of R_L
G_{BR}	Generated bitrate (bits per frame)
$U(s, s')$	Function to update the matrix from s to s'

the system to a desired state while the negative effects of this action are reduced respecting restrictions and taking constraints into account. In other words, the RC defines a QP that optimizes the bandwidth or bit allocation while maximizing the visual quality and reducing bitrate/quality sudden variations.

The bitrate prediction is performed considering the neighborhood correlation at temporal, view, and inter-GOP domains. As discussed in Sect. 4.1.3, there is a high correlation in the temporal and view neighboring frames inside the same GOP. Moreover, there is also a periodic pattern that repeats at GOP level, the GOP-Phase. Our MPC-based RC is able to exploit this correlation in order to accurately predict

Fig. 4.38 MPC-based RC horizons

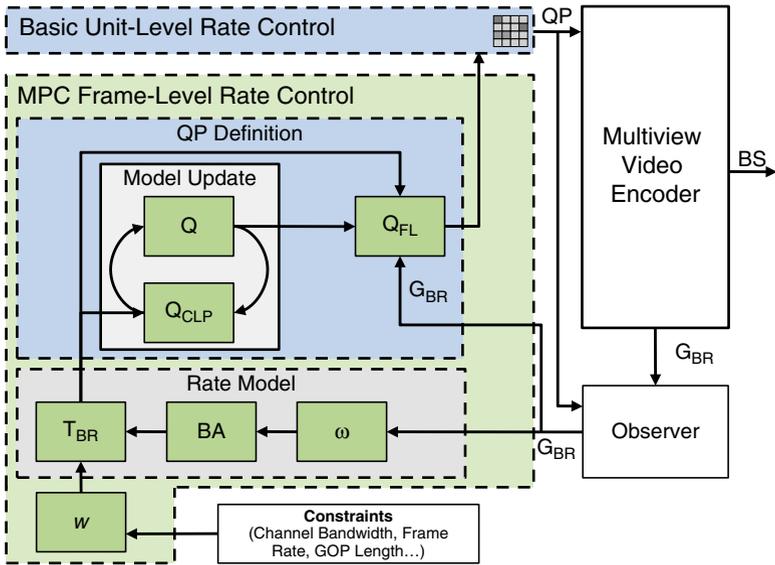
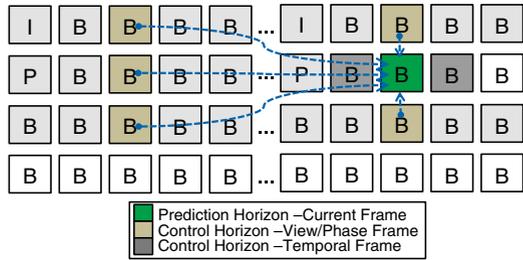


Fig. 4.39 Frame-level rate control diagram

the future bitrate. Figure 4.38 represents the previously encoded frames used for prediction (control horizon) and the current frame to be predicted (prediction horizon) for a given MVC prediction structure. Our method employs a variable weighting factor for frames considering their positions in relation to the current frame. The variable weighting factor is calculated considering the number of references and their distance to the current frame. With this extension our frame-level RC may be directly implemented in any hierarchical bi-prediction structure (HBP) while still catching the GOP-phase correlation.

Figure 4.39 shows in details the MPC optimization process and how the component functions interact to each other. The Rate Model generates, based on the neighborhood correlation, a bitrate prediction for the current frame, the target bitrate. Based on the prediction an optimal QP is defined and the internal model is updated. The system feedback and the actually used QP defined in the BU-level RC are received through the observer.

4.5.2.1 Rate Model

The MPC-based Rate Control defines the target bitrate ($T_{BR(f)}$) considering the bandwidth (BW) and frame rate (FR) constraints along with the neighboring frames weights (w) and their frames bit allocation (BA), as shown in Eq. (4.19):

$$T_{BR(f)} = \frac{BW}{FR} \pm w(BA). \quad (4.19)$$

The feedback and the correlation between frames vary with the type of each frame. The bitrate range of distinct frame types (I, P, and B) lie in different ranges; see Fig. 4.11. Thus, the weighting factors for each frame type must be different. A static weight (w_I) is statically predefined for I frames (Li et al. 2003), while P- and B-frame weights (w_P and w_B) are calculated dynamically considering the weights of temporal neighboring frames. Equation (4.20) shows how the weights are calculated considering the HBP in order to respect the local linearity inside the current GGOP, where \bar{w}_{GOP} is the average of w in the current GOP, f represents the f th frame of a given type (I, P, or B) in the processing order, $u = 1/(L_{GOP-1})$, and L_{GOP} denotes the GOP length. For a smooth weighting propagation, w is limited according to a statistically defined range:

$$\begin{aligned} w_I &= 0.75 \\ w_P &= \max \left\{ w_{f-1} - 2u, \min \left\{ \bar{w}_{GOP} - .25, w_I - 2u \right\} \right\} \\ w_B &= \max \left\{ w_{f-1} - 4u, \min \left\{ \bar{w}_{GOP} - .25, w_P - 2u \right\} \right\} \end{aligned} \quad (4.20)$$

The target bit allocation (BA) is given by a history-based weighted model to optimize MPC for best target bit allocation, as shown in Eq. (4.21). The proposed MPC-based RC was designed to differentiate the frames according to their number of references (0.2 temporal + 0.2 disparity reference frames) as it is an important data to understand how the bit allocation propagates within the 3D-Neighborhood. It allows HRC to respond to variations inside the GGOP and to become more flexible by adapting, without further extensions, to any HBP structure.

The weights $\omega_{i,j}^{m,n}$ (where i and j are the frame time instant and view; m and n denotes the number of references in the temporal and view domains, respectively); calculation is presented in Eq. (4.22):

$$BA_{(f)} = \left(BA_{(f-1)} - \frac{BA_{(f-1)}}{N_A - 1} + \frac{\omega_{i,j}^{m,n}}{\sum_0^m \sum_0^n \omega^{m,n}} - 1 \right) \times \frac{BW}{FR} \times L_{GOP}, \quad (4.21)$$

$$\omega_{i,j}^{m,n} = \frac{(BR_{i,j}^{m,n} \times QP_{i,j}^{m,n} (f-1)) + (L_{GOP} - 2)\omega_{i,j}^{m,n} (f-1)}{L_{GOP} - 1} \quad (4.22)$$

4.5.2.2 Quantization Parameter Definition

Once the prediction is performed, the RC must define a proper action in terms of QP. The QP is determined by summation of all target bitrate ($T_{BR(f)}$) in the prediction horizon, the summation of all generated bitstream in the control horizon (BR), and the history of QPs (H_{QP}), as shown in Eq. (4.23). Note, the QP defined in the frame-level (QP_{FL}) RC is not directly used by the MVC encoder but forwarded to the BU-level RC to refine the QP selection:

$$QP_{FL} = H_{QP} \times \frac{\sum_{i=1}^p T_{BR}}{\sum_{i=1}^m BR}. \quad (4.23)$$

To maintain the performance of our MPC-based controller there is a need to update the QP model. For that, the HRC implements an optimization loop with non-discrete steps (k) where Q_{CLP} denotes the quantization parameter for the frame coded in the last process. Equations (4.24) and (4.25) describe the update process where the QP value is constrained to a variation range of ± 2 QP points for smooth update. In Eq. (4.25) M is the transposed matrix of ω multiplied by target bitrate variation ($\Delta T_{BR(f)}$) for the frames belonging to the control horizon. Q_{st} is the initial QP defined by the user:

$$Q_k = \min \left\{ Q_{(k-1)} + 2, \max \left\{ Q_{(k-1)} - 2, Q \right\} \right\} \quad (4.24)$$

$$Q_{(k-1)} = \min \left\{ QP_{\max}, \max \left\{ QP_{\min}, Q_{CLP} \right\} \right\}$$

$$Q_{CLP} = \sum_{i=L} Q_k \times \det(M (\omega \times \Delta T_{br})^T) \times Q_{st} \times \frac{\sum \Delta \tilde{Q}_k^j}{N_{Fr}}. \quad (4.25)$$

4.5.2.3 Frame-Level Rate Control Evaluation

In the following are presented the detailed results of the frame-level only HRC. Table 4.7 presents the bitrate results generated using SMRC (Single-View Mode Rate Control) extrapolated from the H.264 reference software (JM) using the quadratic MAD prediction (Li et al. 2003). To measure the target bitrate accuracy, we use the Mean Bit Estimation Error (MBEE) metric presented in Eq. (4.26). On average, the proposed frame-level RC provides 1.13 % (up to 1.58 %) of bitrate error while the SMRC provides 2.46 % (up to 2.91 %). The results show that the frame-level HRC predicts more accurately the bitrate behavior and is able to adapt the QP in order to reduce the output error:

$$MBEE = \left\{ \sum_{i=0}^{GOP_{size} \times N_v} \frac{|R_t - R_a|}{R_t} \times 100 \right\} / N_{Fr}. \quad (4.26)$$

Table 4.7 Comparison of frame-level HRC Bitrate accuracy

Video	Target [kbps]	Bitrate [kbps]		Error (MBEE) [%]	
		SMRC	Frame-Level HRC	SMRC	Frame-Level HRC
Ballroom	256	263	259	2.63	1.17
	392	402	396	2.61	1.07
	512	523	518	2.16	1.13
	1,024	1,048	1,032	2.35	0.81
Exit	256	261	258	2.10	0.88
	392	402	397	2.55	1.29
	512	523	519	2.25	1.36
	1,024	1,048	1,038	2.34	1.38
Flamenco2	256	262	258	2.30	0.81
	392	402	396	2.50	1.00
	512	525	517	2.54	1.07
	1,024	1,049	1,035	2.46	1.10
Vassar	256	263	258	2.91	0.84
	392	402	397	2.56	1.25
	512	526	519	2.68	1.36
	1,024	1,049	1,040	2.44	1.58
Average	256	262	258	2.49	0.93
	392	402	397	2.55	1.15
	512	524	518	2.41	1.23
	1,024	1,049	1,036	2.40	1.22
Total average				2.46	1.13

Table 4.8 Comparison of BD-PSNR

Video	SMRC		MPRC	
	BD-BR [%]	BD-PSNR [dB]	BD-BR [%]	BD-PSNR [dB]
Ballroom	10.902	-0.328	28.603	-0.939
Exit	11.542	-0.368	36.920	-1.089
Flamenco	9.630	-0.217	29.852	-0.880
Vassar	6.514	-0.183	20.333	-0.596
Average	9.647	-0.274	28.927	-0.876

The proposed frame-level RC also provides rate-distortion (RD) results that outperform SMRC and the fixed-QP solution (non-RC). Table 4.8 summarizes the quality and bitrate outputs in terms of BD-PSNR (Bjontegaard Delta PSNR) and BD-BR (Bjontegaard Delta BR) (Tan et al. 2005) in relation to the non-RC solution. Compared to SMRC the proposal provides 0.6 dB BD-PSNR increase. The BD-BR reduction is 19.28 % in relation to SMRC.

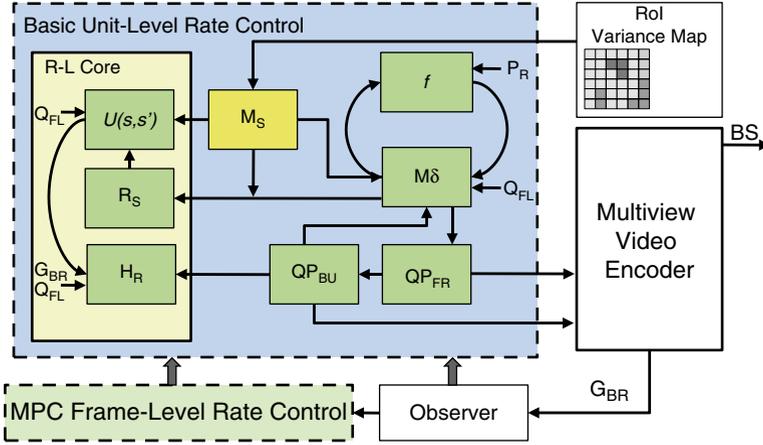


Fig. 4.40 Basic unit-level rate control diagram

4.5.3 Basic Unit-Level Rate Control

Markov Decision Process (MDP) is a mathematically based optimization model of discrete state, sequential decision making in a stochastic environment that depends only on the current state and not in previous states. However, if a controlled MDP is considered, the transition probabilities are affected by previous actions. According to this definition, the controlled MDP perfectly fits to the BU-level rate control where a decision among a set of discrete QP values has to be made considering the neighborhood history. However, in MVC, the transition probabilities between the possible states are not known a priori and vary for distinct time instants and video content. Reinforcement learning can solve MDP with no explicit probabilities definition. It calculates the probabilities of transition based on the Law of Effect theory that states: in case an action is followed by satisfactory state, the probability taking the same action again is increased. It is also possible to incorporate additional variables such as image properties into the reinforcement learning definition.

As part of the HRC we propose a BU-level Rate Control employing Markov Decision Process along with RL able to consider the image properties through a texture-based RoI map, as detailed along this section.

Figure 4.40 depicts the diagram of the proposed BU-level RC that works as a refinement of the frame-level RC. In order to refine the accuracy of bit allocation and provide smooth visual quality, our BU-level RC includes the concept of region of interest (RoI) into a Markov Decision Process that employs reinforcement learning for adapt to dynamic encoder and input variations. At each decision step, the RC monitors the state of the system and determines the next action to take based on constraints observations and the control policy. Firstly, the HRC defines the RoIs for anchor frames generating a map of weights M_S that will determine the importance of each BU inside the frame. Secondly, the weights map is linked to a map of states $M(\delta)$ in the MDP that corresponds to the QP for each BU. The MDP fits to the MVC encoder behavior by providing the structure to make decisions partly random and

partly under a control. Finally, to dynamically adjust the matrix of states for next decision, the RL is responsible to feedback the system response to the current state for both BU-level and frame-level control.

4.5.3.1 Regions of Interest

As discussed in section “Region of Interest” in Chap. 2, frames are composed of regions with distinct image properties requiring a variable number of bits to be encoded. Regular video encoders use the same QP to encode all basic units within a frame leading to inefficient bitrate distribution and undesirable quality variations inside the frame. However, it is possible to define regions to receive special treatment, the regions of interest. The BUs belonging to RoIs may be prioritized by the rate control unit in order to protect the quality of those regions. In this work, the whole frame is considered to have the same semantic relevance (this leave space for further application specific extensions) but regions that present a hard-to-predict content must be allowed to use more bits through QP reduction. According to our analysis (see Sect. 4.1.3), textured regions tend to generate more residue, and consequently, require higher bitrate.

In our solution, the RoI is determined by a normalized variance map—given by M_s in Eq. (4.27)—for all anchor frames. Additionally, HRC also keeps a second matrix of states where each value represents a bitrate of a frame inside a GOP encoding history to incorporate temporal and view neighborhood information to the MDP process. The matrixes data are used by the MDP and RL to define the rewards associated with each state and the actions taken by the control. For non-anchor frames are used the statistics given by anchor frames considering the reinforcement learning R_L :

$$M_{s(i,j)} = \frac{(\text{BU}_i - \mu)^2}{N - 1}. \quad (4.27)$$

4.5.3.2 Markov Decision Process

The HRC implements the BU-level RC by employing the Markov Decision Process. The MDP works over a matrix of independent states $M_f(s)$ representing the QPs of each BU within a frame. Each BU has a set of possible actions A with associated rewards R_s and transition probabilities $f(s,\delta)$. In our model the possible actions are increase, decrease or maintain the QP value defined at frame level, as shown in Eqs. (4.31) and (4.32). A matrix of coefficients $M(\delta)$ is used to define the reward for each action according to Eq. (4.28). The rewards R_s are calculated based on the RoI map M_s , matrix of coefficients $M(\delta)$, and the reinforcement learning R_L (see section “Reinforcement Learning” in Chap. 2), as shown in Eq. (4.29). For each action there is a probability of transition $f(s,\delta)$ defined by Eq. (4.30):

$$M(\delta) = \sum \frac{\text{QP}_{\text{BU}} \times \text{BS}}{\text{Max}_{\text{QP}} \times (T_{\text{BR}} / N_{\text{BU}})}, \quad (4.28)$$

$$R_s = R_L \times |M(\delta) - M_s|, \quad (4.29)$$

$$f(s, \delta) = P_R \mp \Delta\delta, \quad (4.30)$$

$$QP_{BU} = \begin{cases} QP_{FR} + 1 \quad \forall f(s, \delta) > +1 \\ QP_{FR} - 1 \quad \forall f(s, \delta) < -1 \\ QP_{FR} \quad \forall -1 < f(s, \delta) < +1, \end{cases} \quad (4.31)$$

$$QP_{FR} = \text{trunc} \left(\frac{\sum M_f(s)}{N_{BU}} \right). \quad (4.32)$$

4.5.3.3 Coupled Reinforcement Learning

The RL agent incorporates the knowledge of previous events in the decision-making process through monitoring the MVC system response and updating state transitions probabilities and rewards at both frame and BU levels. The BU-level feedback happens by updating the history of reinforcement learning h_R ; see Eq. (4.33). Equation (4.34) gives the final MDP state matrix that is used as obtained knowledge for the upcoming frames. The QP of the frame updated using Eq. (4.34) and calculated according to Eq. (4.30) QP_{FR} provides feedback to the MPC at frame level:

$$H_R = \frac{\Delta T_{BR} \times \sum QP_{BUL}^k}{\sum G_{BR L}^k \times \Delta QP_{FL}}, \quad (4.33)$$

$$U(s, s') = QP_{FL} \begin{cases} M_f(s, s') \quad \forall -1 > f(s, \delta) > +1 \\ M_f(s, s) \quad \forall -1 < f(s, \delta) < +1 \end{cases} \quad (4.34)$$

4.5.4 Hierarchical Rate Control Results

In this section are presented the detailed results of the proposed HRC compared to the baseline solution (i.e., the JMVC without RC) and to the SMRC (Single-View Mode Rate Control). The comparison with the state of the art is presented in Chap. 6. Table 4.9 presents the accuracy in terms of MBEE (less is better) for our HRC compared to baseline RC solutions. The test conditions are detailed in Sect. 6.1. On average, our Hierarchical Rate Control provides 1.6 % MBEE decrease while ranging from 0.7 % to 1.37 %. The superior accuracy is a result of the ability to adapt the QP jointly at frame and BU levels considering the neighborhood correlation and the video content properties.

Table 4.10 presents the objective rate distortion in BD-PSNR (Bjøntegaard Delta PSNR) and BD-BR (Bjøntegaard Delta Bitrate) (Tan et al. 2005) in relation to JMVC without RC. The HRC provides 1.86 dB BD-PSNR increase along with

Table 4.9 Comparison of proposed HRC Bitrate accuracy

Sequence		Bit-Rate [kbps]				MBEE [%]			
		Target	JMVC	SMRC	HRC	JMVC	SMRC	HRC	
VGA	<i>Ballroom</i>	256	268	263	258	4.64	2.63	0.75	
		392	408	402	395	4.06	2.61	0.78	
		512	529	523	516	3.33	2.16	0.78	
		1024	1058	1048	1032	3.30	2.35	0.78	
	<i>Exit</i>	256	267	261	258	4.29	2.10	0.94	
		392	408	402	396	3.99	2.55	0.92	
		512	528	523	516	3.21	2.25	0.83	
		1024	1056	1048	1031	3.14	2.34	0.72	
	<i>Flamenco 2</i>	256	268	263	258	4.79	2.91	0.71	
		392	409	402	395	4.34	2.56	0.71	
		512	530	526	516	3.56	2.68	0.84	
		1024	1059	1049	1031	3.41	2.44	0.70	
	<i>Vassar</i>	256	267	262	258	4.27	2.30	0.75	
		392	407	402	395	3.73	2.50	0.72	
		512	528	525	516	3.13	2.54	0.86	
		1024	1056	1049	1033	3.15	2.46	0.86	
Average		256	268	262	258	4.50	2.49	0.79	
		392	408	402	395	4.03	2.55	0.78	
		512	529	524	516	3.31	2.41	0.83	
		1024	1057	1049	1032	3.25	2.40	0.76	
XGA	<i>Break dancers</i>	512	525	524	518	2.47	2.41	1.23	
		768	801	788	776	4.33	2.54	1.08	
		1024	1052	1050	1034	2.72	2.56	1.00	
		2048	2101	2109	2070	2.58	2.99	1.06	
	<i>Uli</i>	512	525	525	519	2.46	2.54	1.37	
		768	801	789	776	4.28	2.72	1.08	
		1024	1052	1052	1034	2.74	2.72	0.95	
		2048	2101	2101	2069	2.59	2.60	1.05	
	Average		512	525	525	519	2.46	2.48	1.30
			768	801	788	776	4.30	2.63	1.08
		1024	1052	1051	1034	2.73	2.64	0.97	
		2048	2101	2105	2070	2.58	2.80	1.05	
HD	<i>GT Fly</i>	1024	1050	1049	1037	2.54	2.44	1.27	
		1536	1581	1575	1553	2.93	2.54	1.11	
		2048	2104	2101	2069	2.73	2.59	1.03	
		4096	4202	4219	4140	2.59	3.00	1.07	
	<i>Poznan Hall2</i>	1024	1049	1050	1038	2.44	2.54	1.37	
		1536	1582	1578	1553	2.99	2.73	1.11	
		2048	2104	2104	2068	2.73	2.73	0.98	
		4096	4202	4203	4139	2.59	2.61	1.05	
Total Average						3.40	2.55	0.95	

Table 4.10 BD-PSNR and BD-BR comparison

JMVC 8.5 vs.	VGA					XGA			HD1080p		AVG
	Ballroom	Exit	Flamenco2	Vassar	Bdancer	Uli	Poznan	GTGly			
SMRC	0.328	0.368	0.217	0.183	0.215	0.208	0.253	0.012	0.223		
BD-BR	-9.831	-10.348	-8.784	-6.116	-8.963	-9.805	-12.180	-6.711	-9.092		
HRC	1.585	2.375	2.103	1.176	2.060	1.870	2.085	2.055	1.914		
BD-BR	-31.588	-47.458	-38.199	-27.335	-46.112	-49.660	-48.760	-47.250	-42.045		

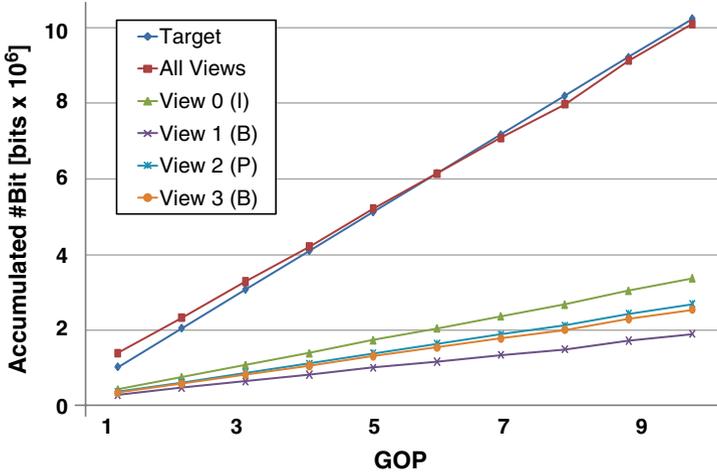


Fig. 4.41 View-level bitrate distribution (Flamenco2)

BD-BR reduction of 40.05 %, on average. If compared to SMRC, the HRC delivers 1.6 dB increased BD-PSNR and 31.08 % reduced BD-BR. Remember, besides superior RD performance the HRC also outperforms SMRC in terms of accuracy.

In the following we present HRC detailed results for *Flamenco2* sequence encoded at 1024 kbps. For simplicity, we analyze only the first four views. Figure 4.41 shows the target bitrate, the total accumulated bitrate and the accumulated bitrate for each view. The presented bitrate distribution is smooth also at view level without abrupt oscillations. As expected from the discussion in Sect. 4.1.3, the base view (View 0, I-view) is more bitrate hungry followed by P-views (View 2) and B-views (View 1, 3).

The frame-level bitrate distribution is further detailed for the GOP #8 in Fig. 4.42. It shows, graphically, the smooth bitrate and PSNR variations delivered by our solution considering frame level. Note, the HRC surface presents no sudden variations for both bitrate and PSNR. Compared to the other solutions, it is clear that the bitrate and quality provided by our HRC are significantly smoother even when compared to SMRC.

Analogous analysis was performed to demonstrate the behavior of our RC at BU level. Figure 4.43 shows the bitrate distribution for a frame region (zoomed image) in sequence *Flamenco2*. Observe that for HRC the bitrate varies with the texture complexity due to our RoI-aware MDP implementation. For the homogeneous background, reduced number of bits is spent while for textured objects and borders (dancer) more bits are allocated. Note that, in Fig. 4.43, the HRC bitrate distribution surface plot accurately fits the object shapes. This behavior prioritizes the regions where the HVS requires a higher level of details tending to lead to a superior overall perceived quality. SMRC is unable to accurately react to the image content. In addition, the HRC also results in smoother variations within the same region

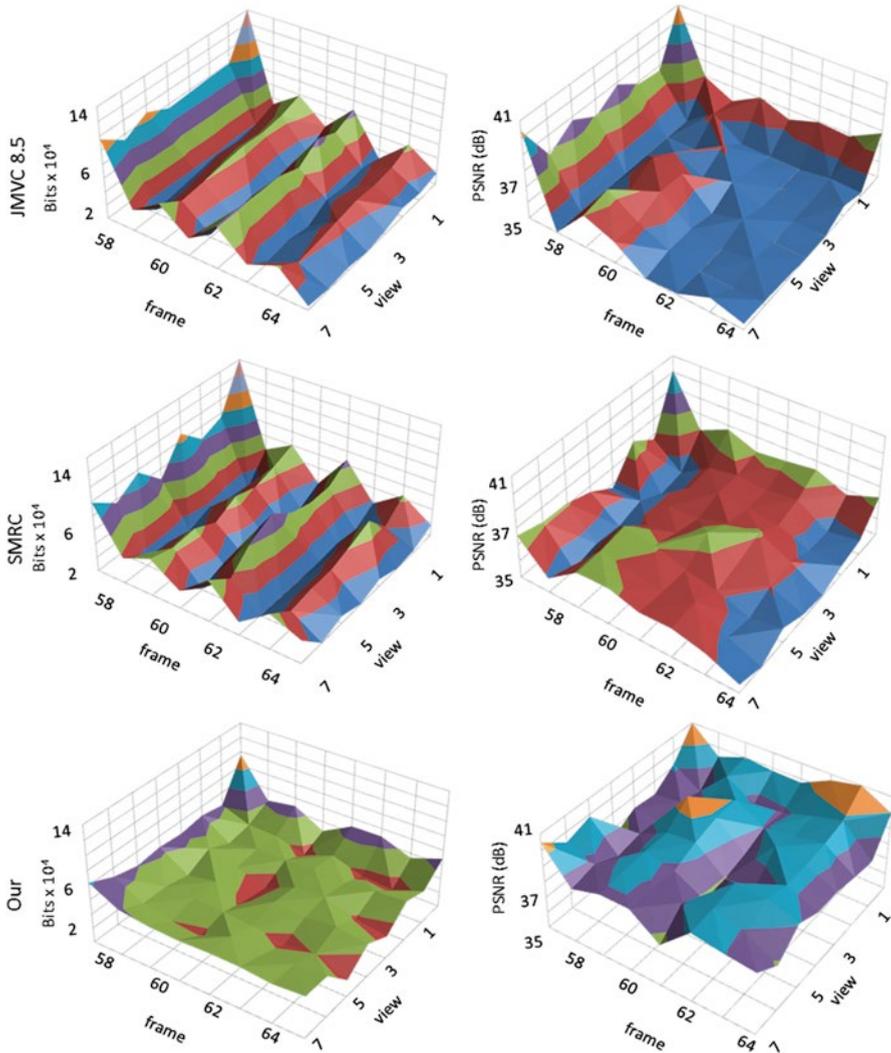


Fig. 4.42 Bitrate and PSNR distribution at frame level (GOP #8)

(dancer’s body or background), as shown in Fig. 4.43. It avoids sudden quality variations and the resulting coding artifacts inside those regions.

The evaluation presented demonstrates that it is possible to maximize the video quality while obeying to bandwidth constraints by implementing an efficient RC algorithm. The proposed HRC is a powerful tool in order to protect the MVC encoder from quality losses typically posed by fast MD and fast ME/DE heuristics such as those shown in Sects. 4.3 and 4.4.

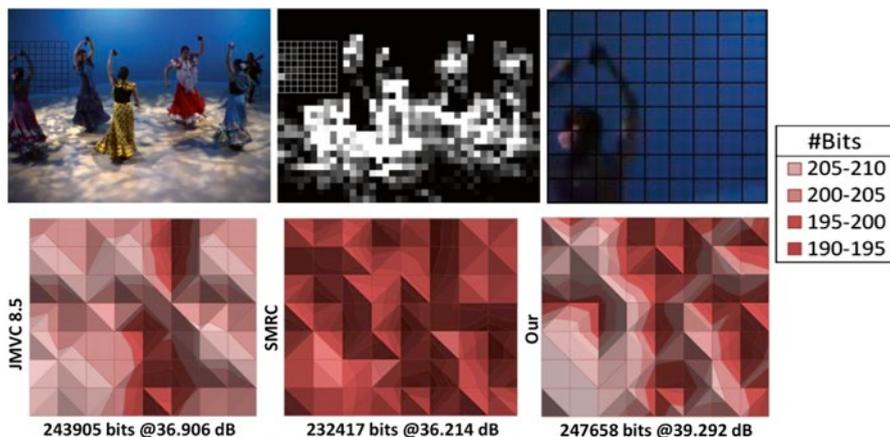


Fig. 4.43 Bitrate distribution at BU level (GOP #8)

4.6 Summary of Energy-Efficient Algorithms for Multiview Video Coding

To provide basis to the energy-efficient algorithms a complete 3D-Neighborhood correlation analysis was presented focusing on coding modes, motion/disparity vectors, and bitrate distribution at both frame and Basic Unit levels.

After that, the Multilevel Mode Decision-based complexity adaptation scheme that incorporates the Early SKIP concept describes a 6-step sophisticated algorithm for complexity reduction. It employs two complexity reduction operation modes while exploiting the 3D-neighborhood correlation along with video properties. To handle the energy versus quality trade-off an energy-aware complexity adaptation algorithm is presented.

Targeting the ME/DE complexity reduction the Fast ME/DE is presented in Sect. 4.4. This algorithm defines two classes of frames: the key and non/key frames. Depending on the prediction mode inferred from 3D-Neighborhood information, the MBs belonging to non/key frames are submitted to the evaluation of 3 or 13 candidate blocks. It represents a meaningful overall complexity reduction.

Aware of the video-quality drawback posed by the energy-efficient MD and ME/DE algorithms, a Hierarchical Rate Control was developed in order to manage and compensate eventual quality drawbacks. The goal is to improve the video quality by optimizing the bit allocation. For that, the HRC operates in frame-level and BU-level rate control actuation levels. At frame level a Model Predictive Controller is used while the BU-level RC exploits a Markov Decision Process along with Reinforcement Learning.

Chapter 5

Energy-Efficient Architectures for Multiview Video Coding

Although the fast ME/DE provides significant complexity reduction, a high-throughput hardware architecture is required for real-time ME/DE in MVC. Without a dedicated hardware, ME/DE for MVC in real-time mobile applications is unfeasible. Therefore, in addition to our fast ME/DE algorithm we propose a novel motion and disparity estimation hardware architecture designed to provide real-time MVC encoding for up to four views HD1080p (1920×1080) based on the proposed fast ME/DE algorithm. Firstly, we are going to start presenting a high-level architectural template description. The architectural template, presented in Sect. 5.1, will give the required basis for a better understanding of the proposed architecture.

In Sect. 5.2 the ME/DE parallelism and hardware scheduling are discussed together. Two scheduling schemes designed for generic search patterns and for the fast ME/DE algorithm proposed in Sect. 4.4 are presented, respectively. The scheduling exploits the multiple levels of parallelism available in the MVC encoding structure. A novel on-chip video memory is presented in Sect. 5.4 to reduce the number of on-chip bits and the external memory communication. The dynamic search window formation strategy (Sect. 5.3) that accurately predicts the memory access pattern from the 3D-Neighborhood is employed to manage the on-chip video memory power state. Finally, in Sect. 5.4.2 an application-aware power-gating algorithm is presented. Also, the memory sizing, partitioning, and management techniques are detailed along this chapter.

5.1 Motion and Disparity Estimation Hardware Architecture

Aware of the dominant memory-related energy consumption we present in this section a hardware architecture featuring novel data prefetching scheme and on-chip memory management solution. Considering the previous MBs memory access our architecture is able to build a search map and predict the memory usage.

It enables to read only the required data from external memory, avoiding performance drawback. Additionally, the mentioned memory access prediction is forwarded to the on-chip memory run-time management in order to adapt the power states of each memory sector, resulting in reduced energy consumption. The features of this architecture are summarized below.

Based on in-depth memory access correlation analysis (see Sect. 5.3.1), it was possible to conclude that the memory access prediction can be further improved in relation to the solution presented in Sect. 5.4.2. As a result, novel memory power-gating control schemes may be proposed. In this section we present a ME/DE hardware architecture featuring an application-aware power-gating scheme able to consider the 3D-Neighborhood correlation and reduce the energy consumption related to memory. The memory hierarchy and power management are carefully explained along this section. Below the main features implemented in this architectural solution are summarized.

Hardware architecture with Multibank On-Chip Memory: A hardware architecture with parallel SAD modules is proposed. A pipelined schedule is proposed to enable high throughput. Moreover, the hardware is equipped with a multibank on-chip memory to provide high throughput in order to meet high-definition requirements. The size and the organization of the memory are obtained by an analysis of the fast ME/DE scheme. Each bank is partitioned into multiple sectors, such that each sector can be individually power-gated to save leakage. The control of the power gates is obtained from the application layer.

An On-Chip Multibanked Video Memory: based on the offline memory usage analysis, an algorithm is proposed to determine the size of the on-chip memory by evaluating the trade-off of leakage reduction and misses (as a result of reduced-sized memory). Afterwards, the organization (banks, sectors) is obtained by considering the throughput constraint. Each bank is partitioned into multiple sectors to enable a fine-grained power management control. The data for each prediction direction is stored in distinct sections.

Dynamically Expanding Search Window Formation Algorithm: Instead of prefetching the complete rectangular search window, a selected partial window is formed and prefetched for each search stage of a given fast ME/DE scheme depending upon the search trajectory, i.e., the search window is dynamically expanded depending upon the outcome of each search stage. An analysis is performed to highlight the issues related to the expansion of the partial window at each search stage. The search trajectories of the neighboring MBs and their spatial and temporal properties (variance, SAD, motion, and disparity vectors) are considered to predict at run time the form of the search window for the current MB. This results in a significantly reduced energy for off-chip memory accesses.

Application-Aware Power-Gating Scheme for the On-Chip Memory: Depending upon the fast ME/DE scheme and the macroblock properties, the amount of required data is predicted. Only the sectors to store the required data are kept powered on and the remaining sectors are power-gated. A power-gating scheme is employed. Depending

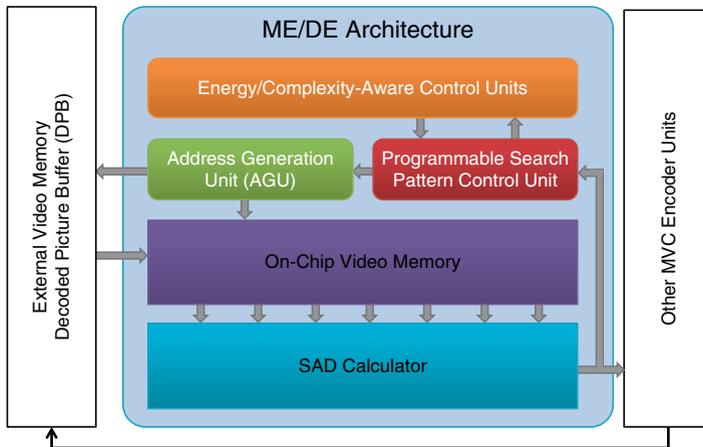


Fig. 5.1 ME/DE hardware architecture template

upon the previous macroblock (i.e., using the knowledge from the application that determines a prediction direction), different sectors can be completely power-gated.

The architectural template overview is presented in Fig. 5.1. It is composed of five main blocks named (a) Energy/Complexity-aware Control Units; (b) Programmable Search Pattern Control Unit; (c) Address Generation Unit (AGU); (d) On-chip Video Memory; and (e) SAD Calculator. Energy/Complexity-aware Control Units box is not detailed in this section since this block represents the implementation of all energy/complexity-aware control techniques presented in Sects. 5.1.1, 5.1.2, 5.1.3, and 5.1.4. Our ME/DE architectures communicate with the remaining MVC encoder blocks by providing SAD values and motion/disparity vectors for the mode decision unit. The reference frames data is read from the external memory that stores the decoded picture buffer (DPB). The MVC encoder writes the DPB after the encoded frames are reconstructed and filtered by the deblocking filter.

The proposed Programmable Search Pattern Control Unit was designed employing a microprogrammed style in order to facilitate the implementation and experimentation of multiple search patterns. It communicates with the Energy/Complexity-aware Control Units in order to provide search pattern information such as search pattern used, memory regions accessed, and number of candidate blocks tested. Energy/Complexity-aware Control Units feedback the Programmable search pattern control unit with energy/complexity budget, search pattern to be employed for future MBs, vector predictors, search directions to be exploited, active on-chip memory sectors, etc. This communication and the hardware actually implemented inside the Energy/Complexity-aware Control Units depend on which energy-efficient techniques are designed for the specific architectural solution.

Once the search pattern is defined, the candidate blocks are forwarded to the AGU as a set of points inside the search window. The AGU is responsible for translating these points into a sequence of actual memory addresses. As the on-chip

video memory is implemented in a cache fashion, the cache tags are generated using the address provided by the AGU according to a predefined tag format defined in Sect. 5.1.3. The on-chip video memory is implemented using SRAM memory to locally store samples belonging to the search window. The samples are brought from the external memory in block-based read operations. To check if the samples required by the search control are available on-chip, the above-mentioned cache tags are tested employing a fully associative approach.

The sum of absolute differences (SAD) Calculator is composed of an array of 4-sample SAD Processing Elements (PEs), an array of adder trees, and comparator trees. The number of PEs depends on the throughput require. The PEs connectivity depends on the block sizes supported by the architecture and the number of candidate blocks processed in parallel. The SAD Calculator is fed in parallel by the on-chip video memory. The number of PEs and the memory width must be jointly defined in order to maximize the hardware usage and processing throughput.

In the following sections the ME/DE hardware modules are presented in details.

5.1.1 SAD Calculator

All the data processing is performed in the SAD Calculator unit. It receives the current MB samples, which are stored in a small local buffer (omitted in Fig. 5.1), and the reference samples to determine the SAD between the original block and the reference block according to Eq. (5.1):

$$\text{SAD} = \sum_{i=1}^n |\text{Orig}(i) - \text{Ref}(i)|. \quad (5.1)$$

Each Processing Element, as depicted in Fig. 5.2, calculates the SAD for four samples in parallel. PEs are composed of four subtractors, one absolute operator, and three adders. Although the hardware description supports multiple sample bit-depth the implementation was limited to 8-bit sample inputs. The PEs are associated using adder trees to generate the SAD for a whole block of $N \times N$. In the example presented in Fig. 5.2, the SAD Calculator is designed to process a 4×4 block in parallel by associating four PEs (PE0...PE3 process one 4×4 block). In this scenario, each adder tree requires further three adders in two logic levels. The larger the block to be processed, the bigger the adder tree. For 16×16 blocks, 63 adders are required in six logic levels. Therefore, pipelining is required for bigger block sizes in order to avoid operation frequency reduction. For simplicity, Fig. 5.2 omits pipeline barriers.

After the SADs are calculated for the multiple block processed in parallel, the SAD Comparators Tree is used to select the smallest SAD values. Along with the SAD value the SAD Calculator feedbacks the Programmable Search Pattern Control with the position where the smallest SAD was found. This information is used to

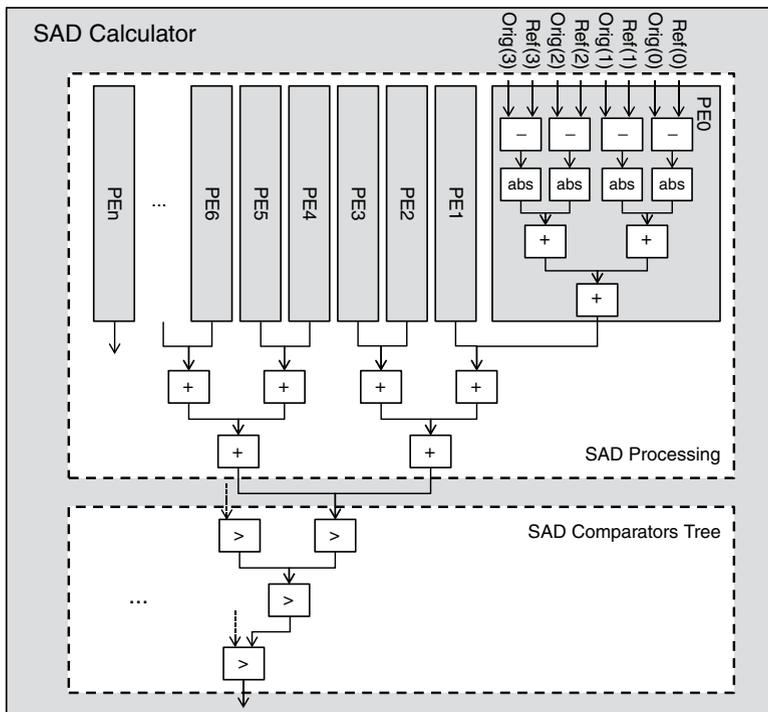


Fig. 5.2 SAD Calculator architecture

decide the following steps of the search process. The SAD Comparators Tree size and logic depth depend on the number of blocks processed in parallel.

5.1.2 Programmable Search Control Unit

Hardware implementations for the search control unit are typically limited to one single search pattern. In the architecture proposed in this monograph, we implement a Programmable Search Control Unit able to support multiple search patterns without hardware redesign by employing the microprogramming concept. By simply reprogramming the Search Pattern Memory (SPM) it is possible to change the search pattern (or shape). It allows fast hardware ME/DE algorithms design and verification.

The Programmable Search Control Unit is composed of a finite state machine (FSM) and an SPM, both presented in Fig. 5.3. Firstly, the FSM identifies the current MB position and reads, from the SPM, the first search pattern. By adding the current MB position and the coordinates of each search point defined in the SPM,

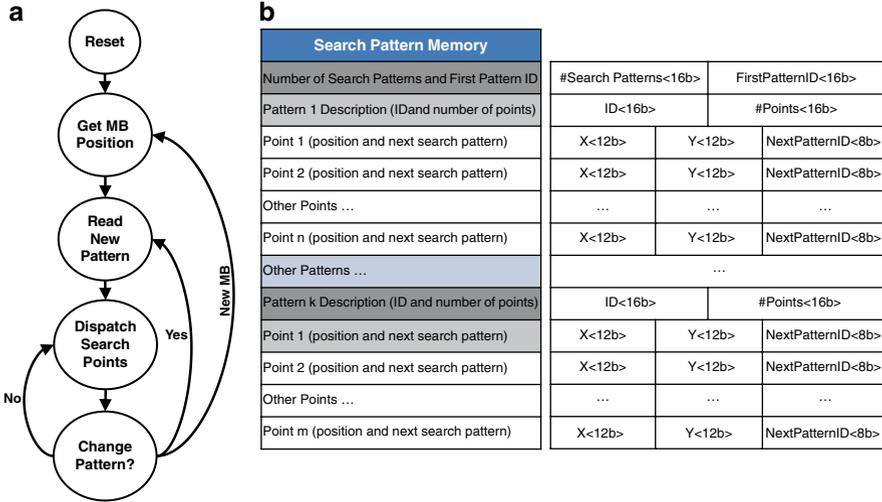


Fig. 5.3 Programmable Search Control Unit (a) FSM and (b) program memory

the Programmable Search Control Unit determines and dispatches, in parallel, all the search points for that specific search pattern step. Depending on the feedback from the SAD Calculator, the next search step is selected among three options: repeat the same pattern, use another pattern described in the memory, or process next MB.

The SPM program memory organization is presented in Fig. 5.3. The left table shows the description of each line while the right table specifies the fields and bit-depth of each field (the number of x -bits is represented as $\langle x \rangle$). A 32-bits program memory is used. The first SPM line brings the total number of patterns programmed and the ID of the first pattern (*FirstPatternID*) to be used, where each search pattern has a unique ID sequentially defined. The search pattern is described starting by a line containing a 16-bit ID (actually only the eight LSB are considered) and the number of points belonging to that specific search pattern. In the following, each point is described using a (X, Y) coordinates pair and the *NextPatternID*. X and Y are 12-bits integer numbers representing the displacement between the search point and the search pattern central reference point. The search pattern central reference is initially defined as the MB position and evolves according to the algorithm interaction assuming the best SAD point as center. The 12-bit coordinates enables a search range of up to $[\pm 2048, \pm 2048]$ in relation to the search pattern central reference. The *NextPatternID* specifies the next pattern to be used in case this point presents the lowest SAD among all points of the current pattern. In case the search point represents a terminating point (in case it is the lowest SAD the search ends) the *NextPatternID* is defined as the reserved value 0xFF.

Table 5.1 provides a simple example using a two-step Log Search (Marlow et al. 1997) with window size $W=16$ (search range $[\pm 8, \pm 8]$) and finishes with a local Cross-Search (Ghanbari 1990) refinement. The first Log Search step (ID 0x0000)

Table 5.1 Search Pattern Memory example

Addr	Instruction		
0	0x000003	0x000001	
1	0x0000	0x0009	
2	0	0	0x02
3	8-	0	0x01
4	8-	8-	0x01
5	0	8-	0x01
6	8	8-	0x01
7	8	0	0x01
8	8	8	0x01
9	0	8	0x01
10	8-	8	0x01
11	0x0001	0x0009	
12	0	0	0x02
13	4-	0	0x02
14	4-	4-	0x02
15	0	4-	0x02
16	4	4-	0x02
17	4	0	0x02
19	4	4	0x02
20	0	4	0x02
21	4-	4	0x02
22	0x0003	0x0004	
23	1-	1	0xFF
24	1	1	0xFF
25	1	1-	0xFF
26	1-	1	0xFF

leads to the second Log Search step (ID 0x0001) except for the central position (line 2) that leads to the Cross-Search refinement (ID 0x0002). After the second Log Search step all points lead to the Cross-Search refinement (ID 0x0002). The terminating step Cross-Search points to the reserved terminating pattern ID 0xFF.

Although it is possible to easily extend the Programmable Search Control Unit, the current implementation requires modifications in the FSM to support features such as early termination and thresholds adaptation.

5.1.3 On-Chip Video Memory

The on-chip video memory used in this monograph works as a cache memory composed of an address comparator block and the on-chip SRAM memory itself, as represented in Fig. 5.4. The address requested by the Search Control and forwarded by the AGU (still using video representation) is compared to the Tags of each

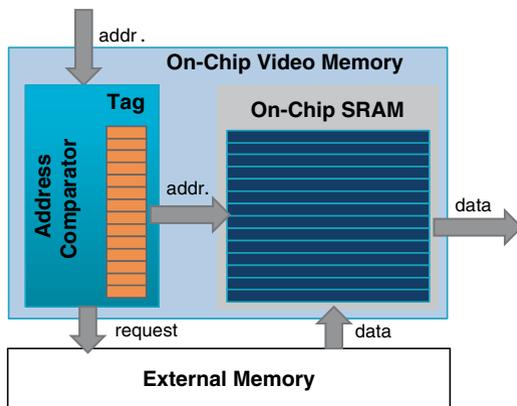


Fig. 5.4 On-chip video memory organization

Cache Tag:

ViewID<4b>	FramePOC<6b>	MBPosX<8b>	MBPosY<8b>
------------	--------------	------------	------------

Fig. 5.5 On-chip video memory cache tag

memory entry. Each entry represents a full MB of the reference frame and the Tags comparison is performed in parallel. In case the reference is available on-chip the requested data is transmitted to the SAD Calculator unit. Otherwise, the read request is sent to external memory. The addresses are provided to the AGU that translates the address from video representation to a burst of addresses mapped to the actual memory address space (see Sect. 5.1.4). After updating the data, the samples are sent to the SAD Calculator.

The tag in the on-chip video memory is defined as shown in Fig. 5.5 where the $\langle nb \rangle$ represents a value with n bits wide. The tag is composed by an unique view identifier, six LSBs of the Picture Order Counter (POC) within that specific view, and the X and Y coordinates of the reference MB. By using this tag it is possible to support up to 16 views, access reference frames within a 64-frames temporal window, and support up to $2k \times 4k$ (QDH) video resolutions. This definition, however, can be easily extended by increasing the bitdepth of each field in order to handle increased demands.

Remember this is just a template description, in the following sections we are going to present, in detail, the SRAM organization, sizing, and energy management of the on-chip video memory for different scenarios.

5.1.4 Address Generation Unit

The Address Generation Unit (AGU) is used to convert the addresses defined in video representation provided by the Search Control to a linear memory

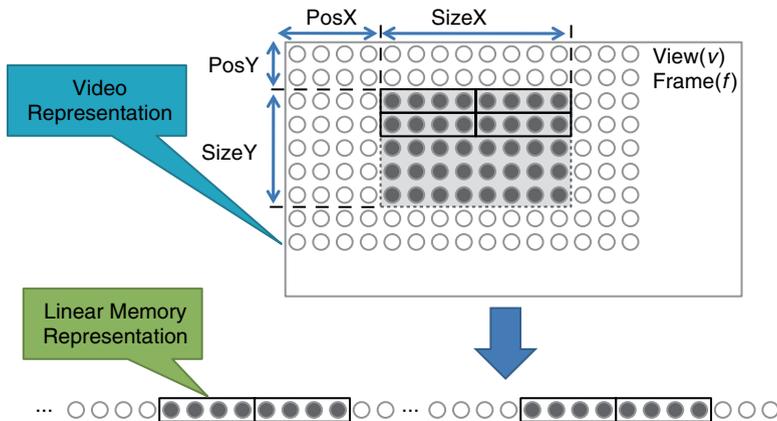


Fig. 5.6 Address Generation Unit (AGU) for dedicated video memory

representation. Video content is represented using 2D arrays; however, when mapped to the external memory these 2D arrays must be translated to a sequence of 1D addresses. This process is depicted in Fig. 5.6, where the circles represent video samples. The math for generating the memory mapped addresses is detailed in the following paragraph. Note, this AGU is used for reading-only purpose since other blocks of the MVC encoder are responsible for writing the external video memory.

The view (v) and frame (f) identifiers associated with the frame resolution ($FrameWidth$ and $FrameHeight$) are used to define the frame base address ($FrameBaseAdd$), as presented in Eq. (5.2). Also, the line $LineStride$ is defined by the frame width. To read a block pointed by positions ($PosX$, $PosY$) and size ($SizeX$, $SizeY$) a sequence of linear memory addresses (Add_0 , Add_1 , etc.) are generated by the AGU. If the block requires more than one access per sample line (as the example Fig. 5.6) two sequential addresses are generated by the AGU. This process is repeated $SizeY$ times, always considering the $LineStride$ displacement, to complete the block reading. The general linearization definition is provided in Eq. (5.2). Note that the video representation addresses refer to sample positions, and to map it to memory positions the memory word size ($MemWordSize$) in number of samples is taken into consideration.

$$\begin{aligned}
 Line\ Stride &= \lceil FrameWidth / MemWordSize \rceil \\
 FrameStride &= LineStride * FrameHeight \\
 FrameBaseAdd &= (v * FramesperView + f) * FrameStride \\
 Add_0 &= FrameBaseAdd + (PosY * LineStride) + \lfloor PosX / MemWordSize \rfloor \\
 Add_1 &= Add_0 + 1 \\
 Add_2 &= Add_0 + LineStride \\
 Add_3 &= Add_2 + 1 \\
 &\dots
 \end{aligned}
 \tag{5.2}$$

5.2 Parallelism in the MVC Encoder and ME/DE Scheduling

Although significant complexity reduction was achieved through the Fast ME/DE algorithm presented in Sect. 4.4, real-time motion and disparity estimation feasibility for mobile devices depends on energy-efficient dedicated hardware architectures. The dedicated hardware architectures must exploit the parallelism available in the MVC prediction structure and feature an optimized scheduling scheme in order to optimize the hardware usage and energy consumption. Therefore, a deep understanding of the distinct parallelism levels and search algorithms behavior is required. Along this section four levels of parallelism available in the MVC encoder are discussed in Sect. 5.2.1. Possible scheduling schemes are proposed and detailed in Sect. 5.2.2.

5.2.1 Parallelism in the MVC Encoder

Due to the prediction structure used in the MVC as depicted by the arrows in Fig. 5.7, four levels of parallelism can be exploited to achieve high throughput. For an easy understanding, the frames in Fig. 5.7 are ordered according to the coding sequence using numbers for the key frames (KF) and the alphabet order for nonkey frames (NKF). The I frames are not processed by ME/DE and are considered available. Frames 2', 4', and 6' belong to the previous GOP.

View-Level Parallelism: Although MVC defines the *Time First* decoding order (i.e., all frames inside a GOP of a given view are processed and then the next view is processed), this order is not mandatory (i.e., not forced by the standard) during the encoding process, as far as the bitstream respects it. For instance, views S1 and S3 can be encoded completely in parallel after S0 and S2 reference views are available.

Frame-Level Parallelism: Within a view there are frames with no dependencies between them. For example, using one reference frame per prediction direction (1 west, 1 east, 1 north, and 1 south) frames A and B can be processed in parallel. Analogously, it is possible to process the frames C, D, E, and F in parallel.

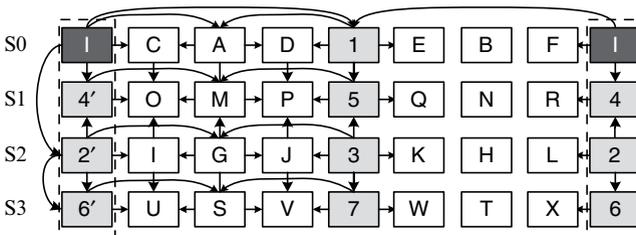


Fig. 5.7 MVC prediction structure in our fast ME/DE algorithm

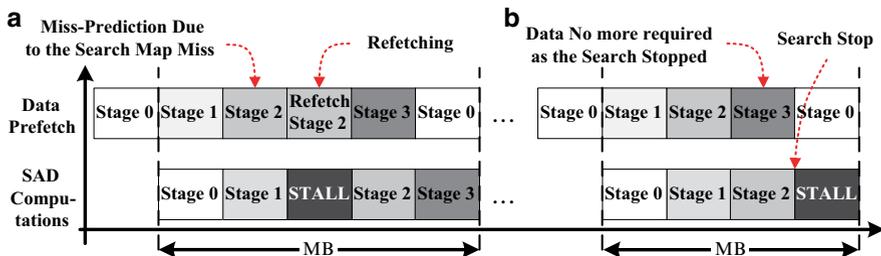


Fig. 5.8 Pipeline processing schedule of our ME/DE architecture

Reference Frame-Level Parallelism: Every single frame can be predicted using up to four reference frames to be encoded. The search in different reference frames has no data dependencies allowing the parallel processing. For instance, while encoding Frame P the search in reference frames 5, D, M, and J may be performed in parallel.

MB-Level Parallelism: Each MB has data dependencies to the previous encoded MB due to motion vector prediction process and SKIP vector prediction. However, using the Fast ME/DE scheme (see Sect. 4.4.1) it is possible to start the predictors evaluation before the spatial neighboring MB. Additionally, the prefetching and SAD calculation may be started before obtaining the previous MB results (also possible for Zero Vector).

5.2.2 ME/DE Hardware Architecture Pipeline Scheduling

Along this section we present two scheduling strategies implemented in our architectural solution. Firstly, the scheduling designed for any search pattern algorithm is described in Sect. 5.2.2.1. In the following, a scheduling scheme specific for our Fast ME/DE algorithm is presented in Sect. 5.2.2.2. The scheduling for our Fast ME/DE considers two logical search control units, one for the KF (using TZ Search) and other for the NKF (using the fast search). The hardware implementation, however, can be simplified and implemented as a single control unit, as presented in Sect. 5.5.

5.2.2.1 Scheduling for Generic Search Pattern

In Fig. 5.8 is presented the MB-level ME/DE processing pipeline scheduling for the selected search algorithm including the data prefetching and SAD computation for different search stages. During the SAD computations of the preceding search stage, the partial search window data is prefetched for the succeeding search stage. However, in case of a Search Map miss (see Sect. 5.3.2) where the data needs to be

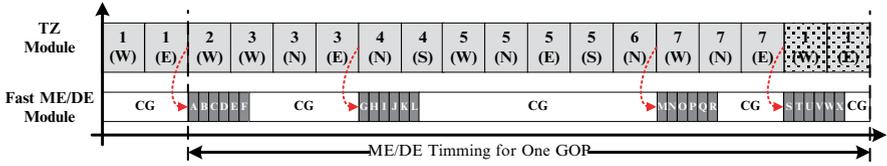


Fig. 5.9 GOP-level pipeline schedule

fetched again, stall for one candidate data prefetch happens (see Fig. 5.8a). In case the search algorithm stops due to early termination criteria, the prefetch data in the search window is wasted (see Fig. 5.8b). This scheduling scheme was employed in our hardware architecture and experimented for TZ Search and Log Search algorithms in order to evaluate the performance while considering the Search Map prediction hits/misses.

5.2.2.2 Scheduling for the Fast ME/DE Algorithm

A novel scheduling scheme (Zatt et al. 2011c) designed for the Fast ME/DE search algorithm (Sect. 4.4) is proposed and depicted in Fig. 5.9. The numbers and letters are consistent to Fig. 5.7. The letters between parentheses represent the prediction direction E (East), W (West), N (North), and S (South). The dotted boxes represent a frame that belongs to the next GOP. This scheduling assumes the existence of two logic control units operating in parallel: one processing the TZ search for KF and the other processing the fast ME/DE for NKF. However, it can be easily mapped to a serial architecture or extended to a more parallel hardware to further exploit the multiple levels of parallelism inherent to MVC.

Each time slot of the TZ Module is dedicated to perform the search for a complete KF in one reference frame. It is noticeable that the coding time for a given GOP is the time to perform 16 TZ searches. This number represents a reduction of 81 % in the number of complete TZ searches if compared to a system without our fast ME/DE search (that performs 88 complete TZ searches).

For NKF encoding there is a Fast ME/DE module. After the required reference frames are processed by the TZ module (solving the data dependencies) all NKF in the same view are processed following the predefined coding order (as shown by the alphabetic order). The data dependencies between the KF and NKF are represented in Fig. 5.9 by dashed arrows. To avoid pipeline stalls, the GOP-level pipeline starts the TZ search in KF of next GOP while the fast ME/DE Module concludes the current GOP processing. Since the Fast ME/DE represents less than 1 % of the processing effort of TZ Search the fast ME/DE module processes the NKF in a burst, and in the following, it is clock-gated (CG) until the next usage. For simplification, in Fig. 5.9, the encoding of NKF does not present the details showing which prediction direction is tested. However, in the slot of given frame A, all required prediction

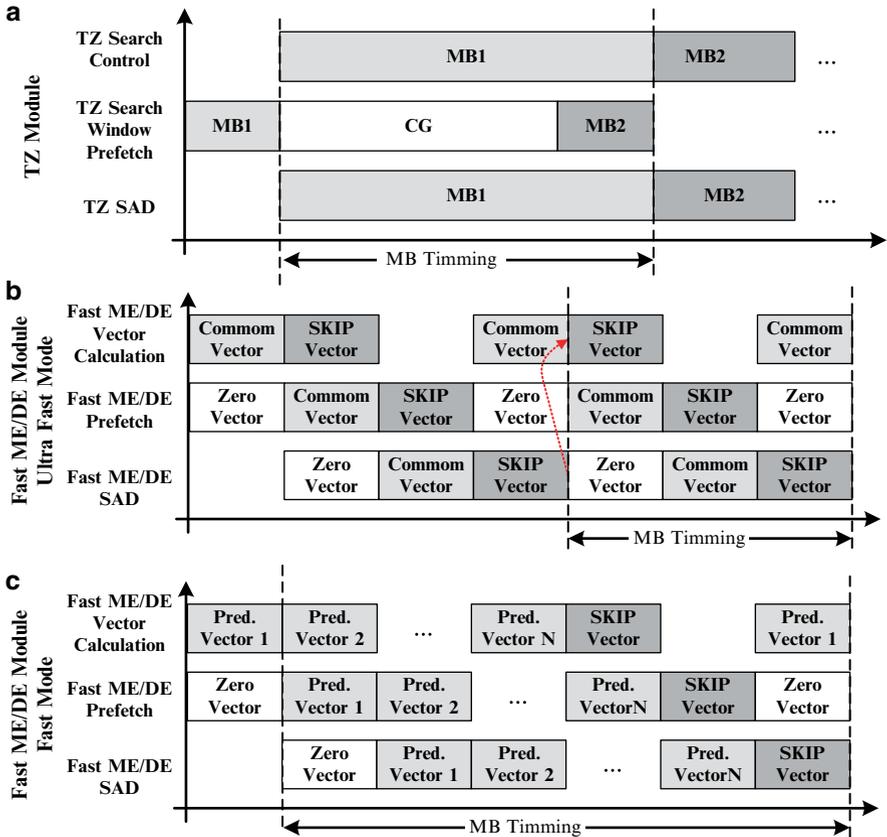


Fig. 5.10 MB-level pipeline schedule for (a) TZ module and fast ME/DE module in (b) Ultra fast and (c) Fast operation modes

directions are serially tested (in the specific case of frame A, West and East directions).

The internal scheduling of the TZ Module operates at MB level, as presented in Fig. 5.10a. The three main tasks are the algorithm control which is always active, the TZ search window prefetch logic which can be clock-gated after bringing the search windows from the external reference memory (DPB), and the TZ SAD computation that starts processing as soon the first useful reference block is available.

As the Fast ME/DE scheme has two prediction modes (the Fast and Ultra Fast prediction modes) two distinct pipeline scheduling are required for the Fast ME/DE Module. The Ultra Fast scheduling is presented in Fig. 5.10b and the Fast scheduling in Fig. 5.10c. Three tasks are considered: Fast ME/DE Vector Calculation, Data Prefetching, and SAD calculation. First, the Zero Vector is tested because it has no data dependencies with the spatial neighbors. Afterwards, the predictors are evaluated and the Common Vector (if it exists, for algorithm details check Sect. 4.4) or

Predictor Vector 1 are processed (represented by the gray blocks in Fig. 5.10b, c). This is the second vector evaluation step (MB-Level Evaluation in Fig. 4.33) once the vectors can be calculated based on the Frame-Level Evaluation (Sect. 4.4.1). If the spatial vector points to a different position, additional data is then fetched and processed (Predictor Vector N). The last vector to be tested is the SKIP predictor that depend upon the previous MB. In this pipeline stage, the previous MB MV/DV information must already be available to avoid pipeline stalls. The MB time borders (indicated by the vertical dashed lines in Fig. 5.10b, c) interfaces are the same for both prediction modes allowing the mode exchange (Fast \leftrightarrow Ultra Fast) with no pipeline stalls.

5.3 Dynamic Search Window Formation

Before moving to the memory organization and the power-gating algorithm, we present the access pattern analysis that provides the basis for our memory access pattern prediction. After that, the search maps used to predict the memory access pattern are introduced followed by the Dynamic Search Window formation algorithm. The memory hierarchy and its application-aware power gating are also presented in this section with detailed results. Similar description of our dynamic search window formation can be found in (Zatt et al. 2011d).

5.3.1 *ME/DE Memory Access Pattern Analysis*

Real encoding systems do not implement the exhaustive search (Full Search, FS) but fast search algorithms. Fast ME/DE search algorithms usually are based on multiple search interactions following a given geometric shape and may employ start point selection and early stop criteria to reduce the computational effort. These algorithms can provide expressive speedup and reduced number of memory accesses at the cost of negligible quality loss. However, real systems may suffer due to the irregular memory access pattern of external memory.

As a case study we present the memory access pattern for two fast ME/DE algorithms: TZ Search and Log Search, considering low- and high-motion MBs (see Fig. 5.11). These search algorithms are implemented in the MVC reference software (JMVC) and their behavior represent a wide family of search algorithms. During ME the high-motion areas perform higher number of memory accesses compared to low-motion areas in which the search converges quickly. Analogous behavior happens for DE where objects with high disparity require more effort to find a good match. In Fig. 5.12 the memory access profile for one frame is presented. The flat regions represent the low-motion/disparity areas while the peaks are located at high-motion/disparity ones. Other important observation is that for a same image region or object the number of memory access and the search pattern



Fig. 5.11 ME/DE search pattern for TZ search and Log search

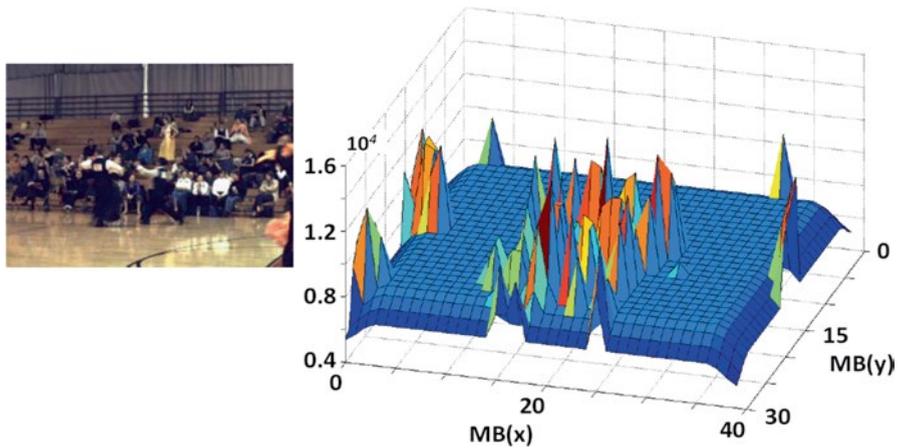


Fig. 5.12 Number of pixels accessed in external memory

behavior is similar, i.e., neighbor MBs that belong to the same object tend to have similar memory and search pattern behavior.

Even considering high-motion/disparity regions it is noticeable that big part of the search window is not accessed resulting in communication and storage energy wastage. Averagely, the ME/DE accesses 19.85 % of the total search window using TZ search and 1.01 % using Log Search. This represents that most of the search window is read and stored in vain, as detailed in Fig. 5.13. The search pattern also is of key importance in the accuracy vs. memory access trade-off. Compared to Log search, the TZ requires more memory accesses (see Fig. 5.13), reaches extended search area (see Fig. 5.11), and tends to provide more accurate search results.

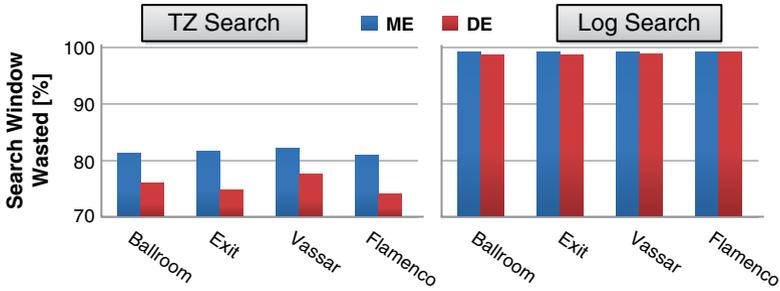


Fig. 5.13 ME/DE search window wastage

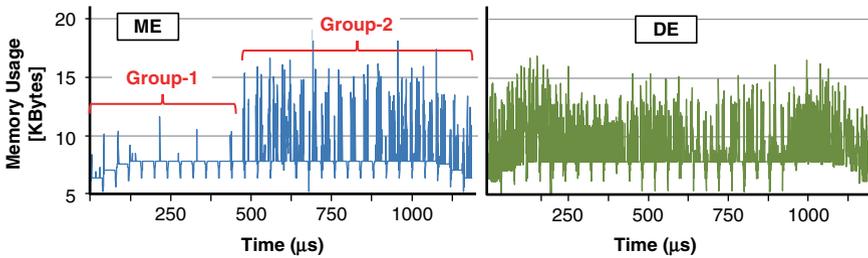


Fig. 5.14 Memory usage variation within one video frame

When further analyzing the memory requirements within a frame (see Fig. 5.14 for *Ballroom* sequence), two different variation zones are noticed in ME that correspond to two different groups of MB, where MBs in a group have similar spatial and temporal properties. MBs in the group-1 exhibit a low variation in their memory usage, while MBs in the group-2 exhibit high-variation. The distinction between two groups can be made by evaluating the average spatial and temporal properties of MBs. Depending upon the group-level variations, low-leakage or high-leakage sleep mode may be selected. The large variations for DE are primarily due to the bigger search performed by the TZ algorithm for capturing longer disparity vectors.

Summarizing, an application-aware power management scheme for an on-chip video memory needs to consider the knowledge of ME/DE algorithm, spatial and temporal video properties (at both frame and MB levels), and correlation in the 3D-Neighborhood to determine the number of idle sectors and an appropriate sleep mode for each sector.

5.3.2 Search Map Prediction

Figure 5.15 presents the *Search Map* for two neighboring MBs (denoted as MB_x and MB_{x+1}) using the Log Search algorithm. After the ME search is performed for the MB_x , a Search Map is built based on the search trajectory (i.e., the ID of the selected

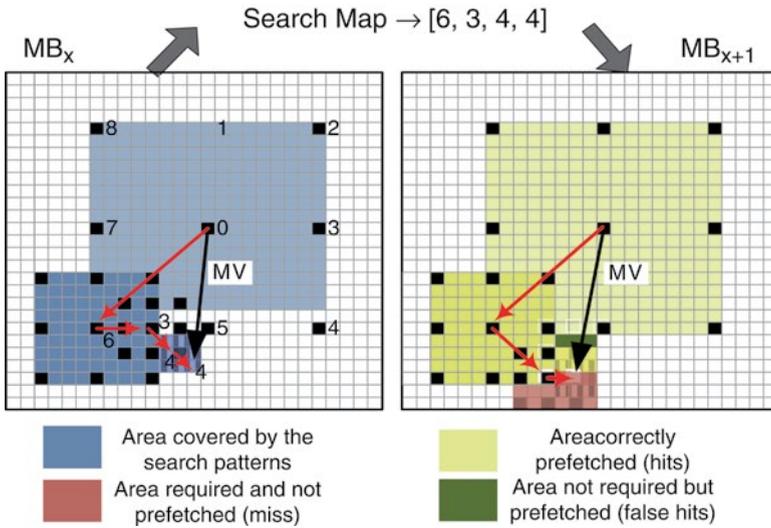


Fig. 5.15 Search Map prediction for the Log search

candidate search points at each search stage of the ME/DE scheme). As shown in Fig. 5.15a, the first search stage selects the candidate with ID 6 as the best candidate. Similarly, candidates with ID 3 (at stage 2), ID 4 (at stage 3), and ID 4 (at stage 4) are selected as the best candidates at their respective search stages. This provides a Search Map of [6,3,4,4] (the trajectory is shown by the red arrows). Note, for each search stage there is an entry in the Search Map with the ID of the candidate with minimum SAD at that particular search stage.

Considering the analysis of the MB neighborhood, a Search Map for the MB_{x+1} can be predicted from the Search Map of MB_x . In case there is a deviation in the search trajectory of these two MBs, there will be a miss in the on-chip memory due to the prefetch of the false region (see the green box in Fig. 5.15b). Typically, these misses are at the boundaries of the moving objects and occur in relatively few MBs along the whole video frame. In case of a miss, there will be a stall only for the prefetching of the first candidate data on the new trajectory (i.e., 16×16 pixel data). All other candidates on the new trajectory will be then prefetched correctly (before their respective SAD computations, thus not causing any stall) as the search pattern design of the fast ME/DE schemes is fixed at design time (see the brown box for the new prefetched data). Typically a miss in the trajectory depends upon the motion/disparity difference of two MBs, which is significantly smaller in most of the neighboring MBs due to high correlation between them.

5.3.3 Dynamic Search Window Formation

Figure 5.16 depicts the pseudo-code for the algorithm of the dynamic search window formation and expansion. The algorithm works in two major steps. First it

```

1. // Predict the Search Map from the Neighboring MBs
2. PredictorSet  $\leftarrow \emptyset$ ;
3. If (PredictorsAvailable) Then
4.   PredictorSet = {MVLeft, MVTop, MVTopRight, MVSpatialMedian};
5.   computeVariance (PredictorSet); //Compute Variance of all predictors
6.   getTemporalInfo (PredictorSet, currMB); //Get MV, DV, SADs
7.   MotioncurrMB = (SADcurrMB > THSAD)? 1: 0;
8.   For i = 0 to all Predictors //Compute the Similarity of predictors, i.e., check if predictors
   belong to the same object as of the current MB
9.     diffVarpredi = VarcurrMB - Varpredi;
10.    Motionpredi = (SADpredi > THSAD)? 1: 0;
11.    diffMotionpredi = MotioncurrMB - Motionpredi;
12.    predDiffpredi =  $\alpha$ *diffVarpredi +  $\beta$ *diffMotionpredi;
13. End For
14. PredictorSet = sortPredictors (predDiff, PredictorSet);
15. bestPred = determineBestPredictor (PredictorSet, currMB);
16. If (predDiffbestPred < THdiff) Then
17.   predSearchMap = SearchMapbestPred;
18. Else
19.   predSearchMap = findClosestSM (PredictorSet, THdiff);
20. End If
21. End If
22. // Perform Dynamic Formation and Expansion of the Search Window
23. For all SearchStages // Depending upon the fast ME/DE scheme
24. SM_Miss = checkSearchMap (searchStageID, predSearchMap);
25. If ((PredictorSet ==  $\emptyset$ ) or SM_Miss) Then
26.   SWBuffer = prefetchPartialWindow (refFrame, searchStageID,
   searchStagePattern);
27. Else
28.   SWBuffer = prefetchPartialWindow (refFrame, searchStageID,
   predSearchMap);
29. End If
30. bestCand = performMEDE (currMB, SWBuffer, SearchAlgorithm);
31. Build_CurrMB_SearchMap (bestCand, searchStageID);
32. If (earlyTermination) return; End If
33. End For
34. return;

```

Fig. 5.16 Algorithm for Search Map prediction and dynamic formation of the Search window

predicts the Search Map from the spatial predictors (lines 3–21). Afterwards, it checks if the search pattern matches the Search Map, prefetches the appropriate partial search window, and updates the *Search Map* (lines 23–33).

Four spatial predictor with presenting high correlation with the current MB are used to predict the Search Map (line 4). Afterwards, variance of these predictors is computed and motion and disparity information is obtained (lines 5–6). Based on the spatial, temporal, and view information, a matching function is computed that provides a hint that predictors may belong to the same object or may exhibit similar motion/disparity properties (lines 7, 9–12). Afterwards, the predictors are sorted with regard to their similarity to the current MB (line 14). The closest predictor is determined by computing the SAD with the current MB (line 15). In case the closest predictor also belongs to the same object or exhibit similar motion/disparity, its

Search Map is considered as the predicted Search Map (line 17). Otherwise, the closest map is found in the remaining predictor set (line 19). If none of the predictors exhibit similarity to the current MB, then the predicted Search Map is empty.

After the Search Map is predicted, it is used to form the search window. For each search stage, the partial search window is determined according to the Search Map and prefetched. In case the search candidates of the search pattern are present in the Search Map (i.e., the search trajectory falls in the predicted region), the partial search window is simply constructed according to the predicted Search Map and the prefetched data is used (i.e., a case of *hit*) (see line 28). Otherwise, if the Search Map is empty or does not contain the search candidate, the Search Map is ignored for this stage onwards (see line 26). In this case the prefetched data is wasted and it is considered as a *miss*. The partial search window is then constructed according to the search pattern for the miss parts (see line 31, it can also be seen in the example of Fig. 5.15b).

In the following section we discuss the architecture of our joint ME/DE scheme along with the design of the multibank on-chip memory and application-aware power gating.

5.4 On-Chip Video Memory

5.4.1 On-Chip Memory Design

On-chip storage of rectangular search windows incurs in increased area and leakage of on-chip memory, like those presented in (Chen et al. 2006, 2007a, b; Ding et al. 2010; Saponara and Fanucci 2004; Tsung et al. 2009; Tsai et al. 2007). The size of the dynamically formed search window is significantly lower compared to the rectangular search windows. This scenario becomes even more critical in MVC where ME and DE are performed for multiple views using larger search windows (for instance $[\pm 96, \pm 96]$ to capture high disparity regions in DE). Depending upon the MB properties, the sizes of dynamically expanding search windows may vary significantly. However, the size of on-chip memory that stores this window must be fixed at design time. Therefore, we firstly perform a design space exploration to obtain a reasonable size of the on-chip memory (that provides leakage power reduction and area savings). In case the MB exhibits low motion and the size of the prefetched window is still less than the on-chip memory, the remaining parts of the memory are power-gated to further reduce leakage.

Figure 5.17 demonstrates the design space exploration for the memory access distribution using *Ballroom* video sequence (a fast motion sequence). Figure 5.17a shows the number of MBs for which ME and DE require less than 96 MBs. Here, a MB-based memory fetching is considered. Please note that the reduced number is mainly due to the adaptive nature of fast ME/DE schemes and it does not mean that this is within a smaller search range. A rectangular search window of $[\pm 96, \pm 96]$

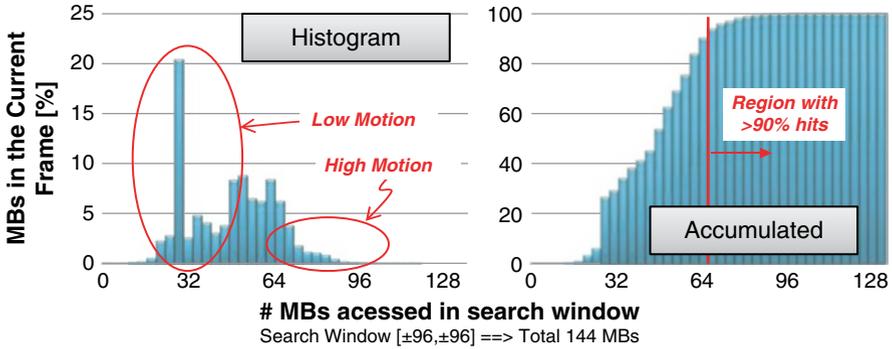


Fig. 5.17 Analyzing the memory requirements for ME/DE of different MBs in Ballroom sequence

size requires 37 KB of on-chip memory. Figure 5.17b shows that even for such a large search range, at most 96 candidates are evaluated per MB. This corresponds to an on-chip memory of 24 KB, i.e., a reduction of 35 % area. When scrutinizing Fig. 5.17b, it is noticed that in more than 95 % cases a storage of only 64 MBs is required (i.e., 16 KB \rightarrow 57 % savings). We have performed such an analysis for various video sequences with diverse motion (not shown here due to space reasons). Similar observations were made in all of the cases. Therefore, we have selected an on-chip memory of 16 KB, which provides significant leakage reduction in the on-chip memory. In rare cases, where the ME and DE may require more MBs, misses may happen (as we will show in Sect. 5.1). The on-chip memory is organized in 16 banks, where one 16 pixel row of an MB is stored in each of the banks, in order to guarantee high parallel throughput.

As discussed above, even 16 KB memory may not be completely used to store the dynamically expanding search window as the size of prefetched search window highly depends upon the MB properties and the fast ME/DE scheme (it can be seen in Fig. 5.17b that in more than 20 % of the cases storage for 32 MBs is used, i.e., only half of the memory). Therefore, each bank is partitioned into multiple sectors (eight sectors in this case) where each sector can be individually power-gated to further reduce the leakage (see Fig. 5.18). The main challenge here is to incorporate the application and video knowledge to determine the power-gate control, such that the power-gating signals may be determined depending upon the predicted memory requirements of the current MB.

5.4.2 Application-Aware Power Gating

A previously discussed, MB properties provides a relatively high potential for leakage reduction by accurately predicting the memory requirements of MBs before

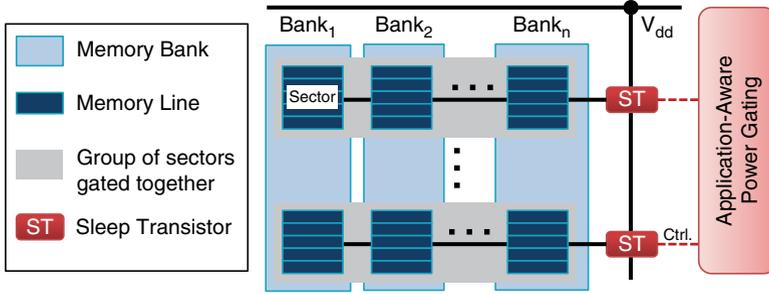


Fig. 5.18 Search window memory organization with power gating

their actual ME/DE. However, frequent on–off switching needs to be avoided to reduce power-gating wake-up energy overhead. Therefore, our scheme predicts the sleep time as function of n consecutive MBs whose sectors can be jointly power-gated. Considering the worst case of stationary MBs, to overcome the wake-up energy overhead, the condition defined in Eq. (5.3) must hold:

$$P_{leak_onChipMem} * T_{minMEDE} * n > E_{wake-up} \quad (5.3)$$

However, the minimum ME/DE time depends upon the deployed fast ME/DE scheme. For instance, in case of a stationary MB there will be a minimum of nine SAD computations for the Log Search and for the TZ Search it is 46 SADs. Therefore, considering the minimum number of SADs for a stationary MB, Eq. (5.3) can be rewritten as Eq. (5.4) where $minNumberSADs$ is 9 and 46 for Log and TZ Searches, respectively. For a given sleep transistor design and a given SRAM memory, the n can be determined. In reality, MBs exhibit diverse motion and spatial properties. Therefore, the number of consecutive MBs that require a certain amount of on-chip memory may be even less than n :

$$n > \left(E_{wakeUp} / P_{leak_onChipMem} * minNumberSADs * T_{SAD} \right). \quad (5.4)$$

Let us assume n consecutive MBs require at most R KB of on-chip memory for their search window prefetching. For a given on-chip memory of size S_{memory} KB with N_{Sec} number of S_{Sec} KB sectors, the amount of power-gateable memory sectors is computed by Eq. (5.5):

$$N_{gateableSectors} = (S_{memory} - R) / S_{Sec}. \quad (5.5)$$

The control signal is generated by the Power-Gating Control unit by simply reading the motion and disparity vectors from 3D-Neighborhood and counting the number of consecutive low-motion/disparity MBs.

5.5 Hardware Architecture Evaluation

5.5.1 Dynamic Window Formation Accuracy

For the detailed experimental results presented in this section a set video sequences with four views each was used. The search algorithm used were *TZ Search* (Yang 2009) and *Log Search* (JVT 2009a) considering three QP values (22,27,32) and search in the four possible directions with a search window of $[\pm 96, \pm 96]$. The thresholds set used were $N=6$, $\alpha=1$, $\beta=500$, and $TH_{SAD}=400$.

Figure 5.19 presents details for the Search Map and on-chip memory evaluation. Figure 5.19 shows that the accuracy of Search Map prediction is higher for low-motion sequences (e.g., *Vassar*) compared to high-motion sequences (e.g., *Flamenco2*) as the search trajectory is shorter and easier to be predicted (due to a higher number of stationary/slow-moving MBs). However, even for the worst case the hits are around 80 % (see Fig. 5.19a). In case of off-chip memory accesses, the misses are higher for low-motion sequences because the search trajectory tends to converge to the center (only the central region of search window is accessed) reducing the overlapping accessed area with the neighboring MBs. The higher number of memory misses for low-motion sequences, however, does not limit off-chip energy savings achieved for the same sequences. The reason is that the percentage of misses is calculated over a much smaller number of total memory accesses for low-motion sequences.

5.5.2 Hardware Architecture Evaluation

Figure 5.20 shows the hardware architecture with our proposed dynamic search window formation scheme. It employs the above-discussed dynamically expanding

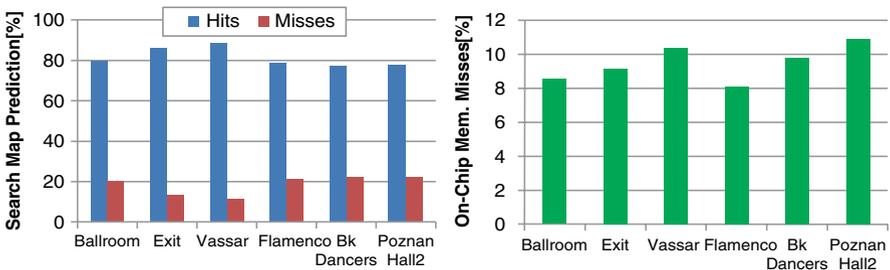


Fig. 5.19 Search Map prediction accuracy and on-chip memory misses

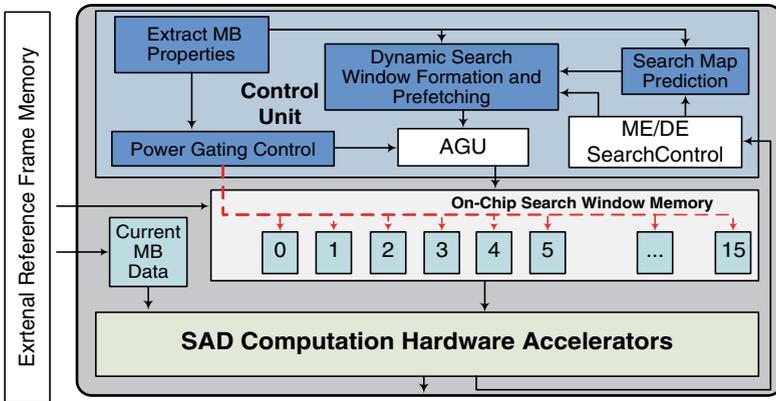


Fig. 5.20 ME/DE hardware architecture block diagram

Table 5.2 Comparison of our fast ME/DE algorithm

	Fast ME/DE architecture w/Dynamic search window formation
Technology	ST 65 nm Low-Power 7 metal layer
Gate count	102 k
SRAM	512 Kbits
Max. frequency	300 MHz
Power	74 mW, 1.0 V
Proc. capability	4-Views HD1080p

search window prefetching and a multibanked on-chip memory with application-aware power-gating control. In order to obtain high throughput, a set of 64 (4×4-pixel) SAD operators and SAD trees is provided as the main computation block. An ME/DE search control unit is integrated which can be programmed to realize various fast ME/DE schemes. This unit controls the search stages and patterns, and it provides the required algorithmic information to various other modules. The search window formation unit predicts the Search Map and dynamically constructs the search window structure. This data corresponding to the window is prefetched in the multibank search window memory which consists of various sectors that can be individually power-gated (Zhang et al. 2005) depending upon the ME/DE requirements of the current MB.

Table 5.2 presents the ASIC implementation results of our architecture. The hardware implementation executes at 300 MHz and provides real-time ME/DE for up to 4-views HD1080p consuming 74 mW. Reduced power is reached mainly due to the employment of dynamic search window formation, power gating, smaller logic, and fast ME/DE scheme. Please refer to Sect. 6.2 for comparison to state-of-the-art ME/DE architectures.

5.6 Summary of Energy-Efficient Algorithms for Multiview Video Coding

Our architectural solution to enable real-time ME/DE is presented along this chapter. Initially, the architectural template and the basic hardware building blocks are described in Sect. 5.1. Based on this structure a pipelined architecture able to implement both regular search patterns and the proposed Fast ME/DE algorithm is described in detail.

Targeting the reduction of the energy consumption related to the external memory accesses and on-chip video memory leakage, the Dynamic Search Window Formation strategy is proposed in Sect. 5.3. This solution observes the search patterns of the neighboring MBs in order to anticipate the data required for the current MB. It allows accurate external memory data prefetching while reducing the on-chip memory size by avoiding the entire search window storage.

Targeting further energy reduction an application-aware power-gating scheme is integrated to the ME/DE architecture. Assuming an on-chip memory with sector-level power-gating capabilities, the application-aware power gating considers the memory usage characteristics of the future MBs along with the wake-up cost to define the sectors for power gating. By doing so, this architecture is able to significantly reduce the overall energy consumption through minimizing on-chip memory leakage.

Chapter 6

Results and Comparison

In this chapter the overall results of this work and the comparison with the latest state-of-the-art approaches are presented. Before moving to the actual comparison, a description of the experimental setup is presented discussing the fairness of comparison in relation to the related works. The benchmark video properties, common test conditions, simulation environment, and synthesis tool chain are also introduced in this chapter. The results for energy-efficient algorithms are discussed in terms of complexity reduction while considering the coding efficiency and video quality in relation to state-of-the-art and optimal solutions. The video quality control algorithm based on rate control is compared to other rate control techniques described in the current literature. Energy-efficient architectures are evaluated against the latest hardware solutions for ME/DE on MVC with emphasis on the overall energy consumption for both memory access and processing datapath. Additionally, throughput and IC footprint area are discussed.

6.1 Experimental Setup

In this section are described the simulation, design, and synthesis environment employed during the development of this work. Afterwards is presented a discussion on the test conditions and benchmark video sequences followed by the fairness of comparison with the state-of-the-art approaches. The hardware design method and synthesis tool chain are also presented in this section.

6.1.1 *Software Simulation Environment*

Each algorithm proposed along this monograph was implemented and evaluated using the reference software platform provided by the Joint Video Team

Table 6.1 Video encoder settings

Parameter	Setting
Entropy encoder	CABAC
FRExt	Yes
QP (experiments w/o rate control)	22, 27, 32, 37, 42
Bitrate (experiments w/rate control)	256, 392, 512, 768, 1024, 2048, 4096
GOP Size	8
Anchor period	8
Temporal coding structure	IBP (Hierarchical B Prediction)
#Views	4/8
View coding structure	IBP (0-2-1-3 or 0-2-1-4-3-6-5-7)
Number of reference frames	Up to 4 (one per temporal/view direction)
Inter-frame/Inter-view prediction pictures first	Inter-frame
B Pictures reference	Yes
Search mode	TZ Search
Search range	Up to $[\pm 96, \pm 96]$
Distortion metric	SAD
Weighted prediction	No
Deblocking filter	Yes

(JVT 2009a), the Joint Model for MVC, also known as JMVC. The JMVC is provided in order to prove the concepts behind the MVC standard and facilitate the experimentation and integration of new tools to the MVC.

Initially, implementations were described on top of the JMVC 6.0, the latest version available by the time this work was started. In face of limitations related to the simulation of HD1080p sequences (note the use of these sequences were normalized in March 2011 (ISO/IEC 2011) after this work was started), our algorithms were migrated to a more recent version, the JMVC 8.5, in order to extend our experimental results. Details on the JMVC software structure and the implemented modifications are detailed in Appendix A.

In Table 6.1 is presented a summary of the video encoder settings and parameters most commonly used for experimentation along this monograph. Note that some settings may vary depending on the experiments nature. These changes, however, are mentioned along the results discussion. Table 6.2 describes the computational processing resources used for simulation.

6.1.2 Benchmark Video Sequences

To allow other researchers to easily compare their results against ours and, consequently, make our results more meaningful to the current literature, the benchmark video sequences used in our experimental section were derived from the common

Table 6.2 Simulation infrastructure

Desktop for simulation	
Processor	Intel Core 2 Duo-6600@2.4 GHz
Main memory	3.25 GB DDR2
Operational system	Windows XP SP2
Mobile device for battery-aware experiments	
Device	HP Pavillion DV6000 Series
Processor	Intel Core-2 Duo T5500 @1.66 GHz
Main memory	2 GB DDR2
Operational system	Windows XP SP2
Battery	6-Cell lithium ion 4400 mAh 10.8 V

Table 6.3 Benchmark video sequences

Sequence	Resolution	# Views	Cameras organization
Ballroom	640×480	8	20 cm spacing; 1D/parallel
Exit	640×480	8	20 cm spacing; 1D/parallel
Vassar	640×480	8	20 cm spacing; 1D/parallel
Race1	640×480	8	20 cm spacing; 1D/parallel
Rena	640×480	100	5 cm spacing; 1D/parallel
Akko & Kayo	640×480	100	5 cm horizontal and 20 cm vertical spacing; 2D array
Flamenco2	640×480	5	20 cm spacing; 2D/parallel (Cross)
Ballet	1024×768	8	20 cm spacing; 1D/arc
Breakdancers	1024×768	8	20 cm spacing; 1D/arc
Uli	1024×768	8	20 cm spacing; 1D/parallel convergent
GT Fly	1920×1088	9	Computer generated
Poznan Hall2	1920×1088	9	13.75 cm spacing; 1D/parallel
Poznan Street	1920×1088	9	13.75 cm spacing; 1D/parallel
Undo Dancer	1920×1088	9	Computer generated

test conditions recommendations provided by JVT (Su et al. 2006) and ISO/IEC (2011). In Table 6.3 are presented the video sequence names along with the number of views, cameras organization, and resolution. The considered video resolutions are VGA (640×480), XGA (1024×768), and HD1080p (1920×1088—typically cropped to 1920×1080) featuring distinct number of cameras, camera spacing, and organization. Although some sequences have up to 100 cameras, our experiments are constrained to four or eight depending on the algorithm under evaluation. Please consider that the main goal of this monograph is on MVC encoding for mobile devices that are not expected to feature more than eight cameras. Nevertheless, the concepts behind the energy reduction algorithms proposed in this monograph are scalable to increased number of views for applications such 3DTV and FTV.

To support the reader that is not familiar with these video sequences, is provided, in Fig. 6.1, the spatial, temporal, and disparity indexes (SI, TI, and DI) for each video sequence referred in Table 6.3. The higher the index the more complex the

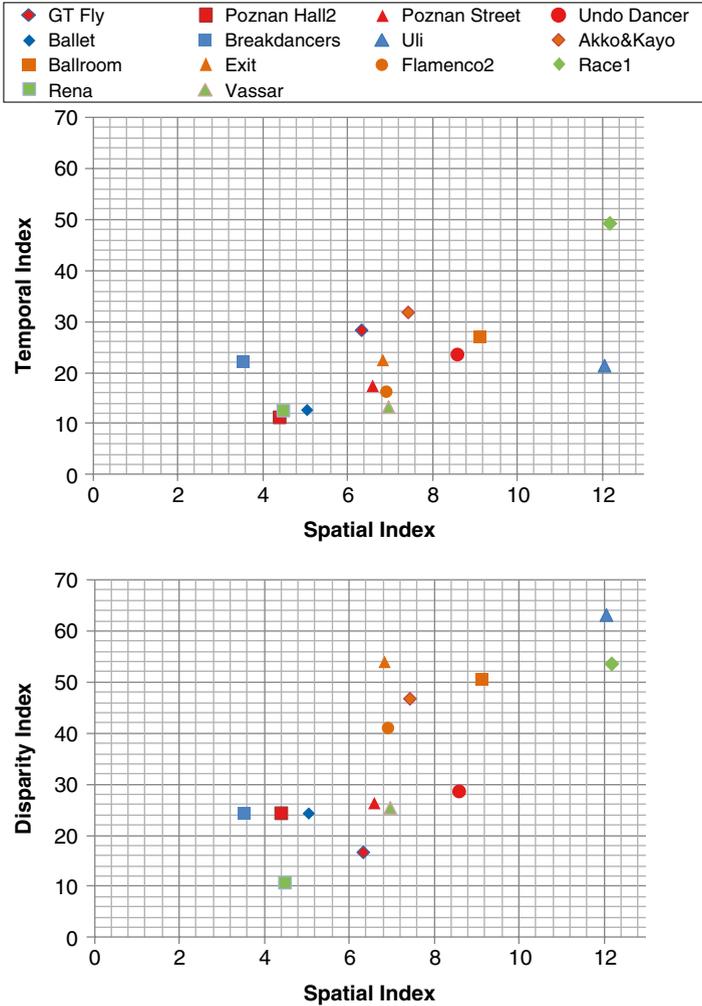


Fig. 6.1 Spatial–temporal–disparity indexes for the benchmark multiview video sequences

sequence is in that specific dimension. The goal is to better understand why some sequences perform better than others under certain coding conditions and/or algorithms. The spatial and temporal indexes were proposed in (ITU-T 1999) and have been used to classify the benchmark video sequences (Naccari et al. 2011) used to test the next video coding standard, the High Efficiency Video Coding (HEVC/H.265). Equations (6.1) and (6.2) give the equations that define SI and TI extended for multiview videos where $\rho(i,j)$ represents the pixel luminance value in coordinates i and j , $Sobel$ denotes the Sobel filter operator, and n is the frame temporal index. Additionally, in order to further adapt to multiview special needs we define the

disparity index (DI) based on the same metric used for TI according to Eq. (6.3) where v is the view index:

$$SI = \max_{\text{view}} \{ \max_{\text{time}} \{ \text{std}_{\text{space}} [\text{Sobel}(\rho(i, j))] \} \}, \quad (6.1)$$

$$TI = \max_{\text{view}} \{ \max_{\text{time}} \{ \text{std}_{\text{space}} [\rho_n(i, j) - \rho_{n-1}(i, j)] \} \}, \quad (6.2)$$

$$DI = \max_{\text{view}} \{ \max_{\text{time}} \{ \text{std}_{\text{space}} [\rho_v(i, j) - \rho_{v-1}(i, j)] \} \}. \quad (6.3)$$

6.1.3 Fairness of Comparison

Although the experimental results were generated using standard benchmark video sequences and standard coding settings, it is frequently not possible to directly compare our algorithms with the results provided by the published related works. For this reason, all state-of-the-art competitors were implemented using our infrastructure based on the information available in the referred literature. This approach requires significant implementation effort overhead; however, it ensures that all algorithms are tested under the same conditions and guarantees the fairness of comparison between all proposed solutions. The simulation infrastructure and modifications applied to the JMVC are presented in Appendix A.

6.1.4 Hardware Description and ASIC Synthesis

The architectural contribution proposed along this monograph includes complete RTL (Register Transfer Level) description, functional verification, and logical and physical synthesis. The hardware was described using VHDL hardware description language followed by functional verification with Mentor Graphics ModelSim (Mentor Graphics 2012) using real video test vectors. The standard-cell ASIC synthesis for 65-nm technologies was performed using the Cadence ASIC Tool chain (Cadence Design Inc. 2012). Two distinct processes and standard-cell libraries were considered in our hardware results, the *IBM 65 nm LPe LowK* (Synopsys Inc. 2012) and *ST 65 nm Low Power* (Circuits Multi-Projects 2012). For preliminary results, FPGA synthesis targeting Xilinx FPGAs was performed using the Xilinx ISE tool (Xilinx Inc. 2012).

As mentioned above, the presented architecture was completely designed, integrated, and tested. The only exception is the on-chip SRAM memories featuring power gating. As far as our memory libraries and memory compiler were not able to generate such memories, regular SRAM memories were instantiated instead for connectivity and area approximation. The SRAM energy numbers were extracted from the related works that describe, implement, and characterize the multiple power states SRAM memories for 65 nm (Fukano et al. 2008; Zhang et al. 2005; Singh et al. 2007). With the energy numbers and ME/DE memory traces, a memory simulator was designed to provide the energy saving results.

6.2 Comparison with the State of the Art

6.2.1 Energy-Efficient Algorithms

In this section is presented the comparison between the energy-efficient mode decision algorithms proposed in this monograph and the state of the art for fast mode decision. The efficiency of the algorithms is measured in terms of time savings compared to the JMVC using RDO-MD. Also, the video quality (PSNR in dB) and bitrate (BR in # of bits) variations are presented using RD curves and the Bjøntegaard rate-distortion metric (Tan et al. 2005).

6.2.1.1 Comparing Our Mode Decision Algorithms to the State of the Art

Figure 6.2 presents the percentage time savings compared to RDO-MD for the early SKIP mode decision (section “Early SKIP Prediction” in Chap. 4), the two strengths of our multilevel fast mode decision (*Relax* and *Aggressive*, Sect. 4.3.1), and two related works (Han and Lee 2008; Shen et al. 2009b). Each bar represents the average for all QPs for that specific video sequence and mode decision algorithm. Even our simplest solution, the early SKIP algorithm, is able to outperform (Shen et al. 2009b) for most of the cases. The work proposed in (Han and Lee 2008) provides time savings superior to the early SKIP but pays a price in terms of video quality, as will be discussed soon. The multilevel fast mode decision shows a superior performance compared to all competitors and provides up to 79 % time reduction. Additionally, it provides two complexity reduction strengths that allow handling the energy saving vs. quality trade-off according to the system state and video content.

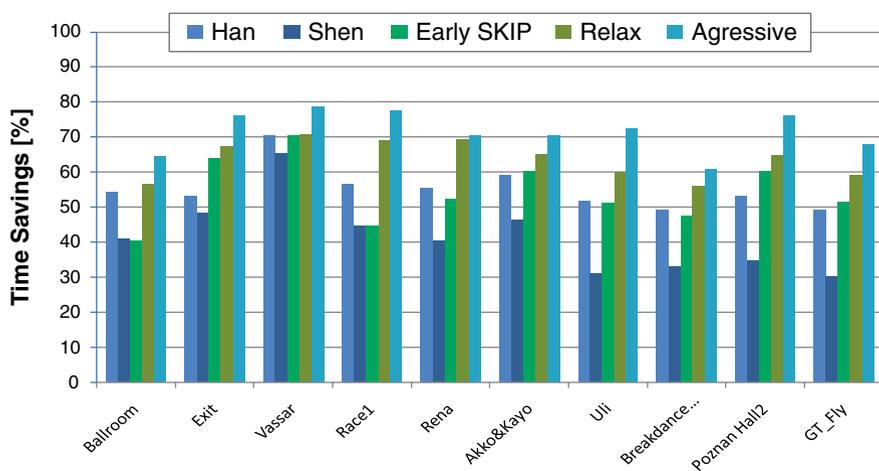


Fig. 6.2 Time savings comparison with the state of the art

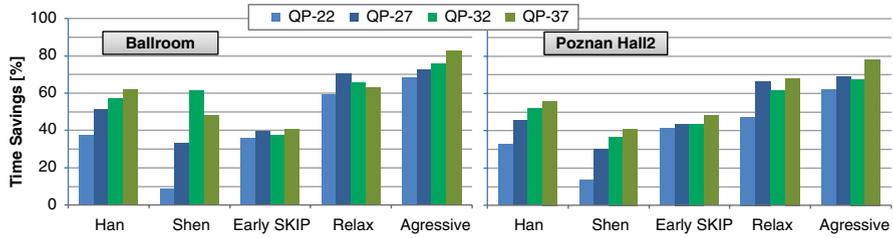


Fig. 6.3 Time savings considering the multiple QPs

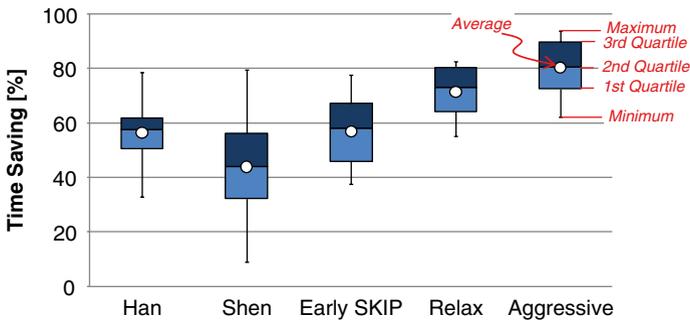


Fig. 6.4 Time savings distribution summary

The multilevel mode decision outperforms the state of the art for all cases while keeping the video quality losses within an acceptable range, as discussed below.

The graph in Fig. 6.3 brings the time savings information detailing its behavior for multiple QPs considering two video sequences, one VGA and one HD1080p. This plot shows that our fast MD algorithms are able to sustain the time savings for the whole QP range due to the QP-based thresholding employed. For instance, the work of Shen et al. (2009b) employs fixed threshold and suffers from reduced time savings specially for low QPs. At high QPs, the fixed thresholds tend to incur increased quality drop. To summarize the complexity reduction results, Fig. 6.4 depicts the distribution of time savings provided by each competitor algorithm considering all video sequences and QPs tested. In summary, the algorithms proposed along this monograph provide averagely higher complexity reduction while sustaining significant complexity reduction for any encoding scenario. While the work of Shen et al. (2009b) shows scenarios with 10 % reduction, the early SKIP prediction provides at least 38 %. The multilevel fast mode decision ensures time savings between 55 % and 90 %.

Besides providing complexity reduction, fast mode decision algorithms must avoid significant video quality losses. In Fig. 6.5 the rate-distortion curves show that for most of the tested video sequences there is a small displacement compared to the RDO-MD solution. The *Relax* level of our multilevel mode decision scheme

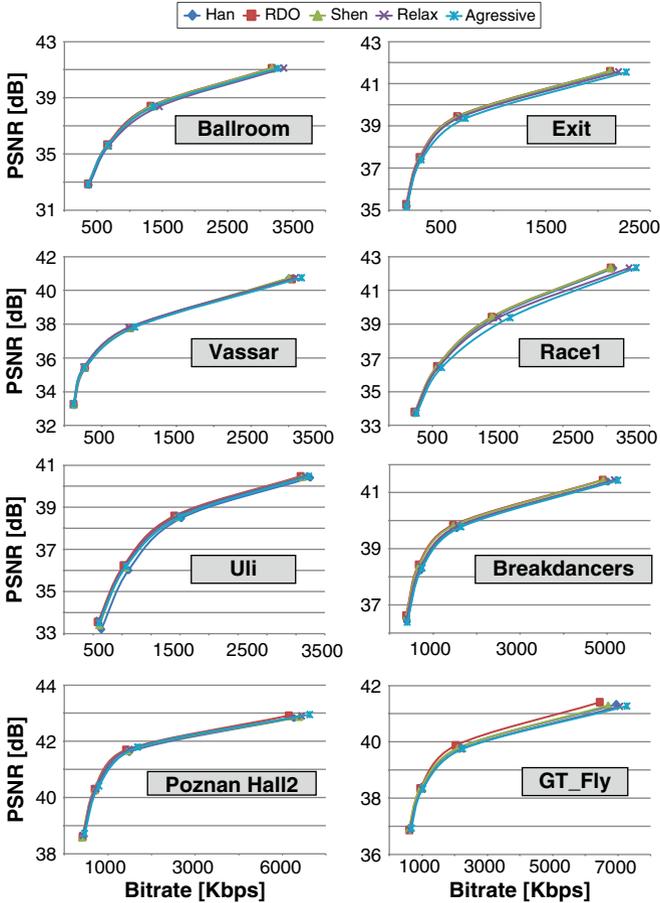
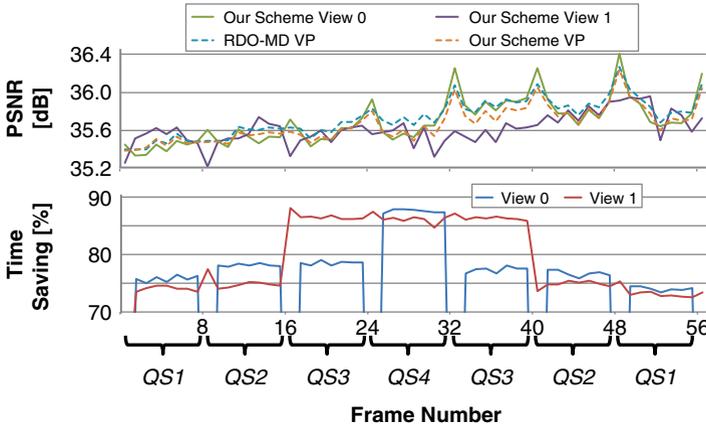


Fig. 6.5 Rate-distortion results for fast mode decision algorithms

provides RD results very close to the exhaustive RDO-MD for most of the cases. The *Aggressive* level incurs slightly worse RD results, especially for *Race1* and *Rena* sequences. Please note that our scheme with both *Relax* and *Aggressive* levels provides much higher complexity reduction compared to all schemes, as discussed earlier. The usage of the *Aggressive* level is recommended if high-complexity reduction is desired (e.g., when the battery level of a mobile device is low). Under normal execution conditions, the *Relax* level is recommended as it provides superior complexity reduction compared to the state of the art while keeping the RD performance close to the RDO-MD. In Table 6.4 is summarized the rate-distortion performance for the discussed mode decision algorithms. Averagely, the early SKIP and relax solutions present the best RD results. The *Aggressive* variant of the multilevel fast mode decision sacrifices RD performance, compared to other competitors, in order to provide the higher complexity reduction.

Table 6.4 Bjøntegaard PSNR and BR for fast mode decision algorithm

Video sequences	Han		Shen		Proposed relax		Proposed aggressive	
	BD-PSNR	BD-BR	BD-PSNR	BD-BR	BD-PSNR	BD-BR	BD-PSNR	BD-BR
Ballroom	-0.163	4.412	-0.054	1.458	-0.106	2.749	-0.272	7.221
Exit	-0.1234	5.278	-0.041	1.749	-0.097	3.960	-0.281	12.047
Vassar	-0.182	8.311	-0.122	5.582	-0.037	1.709	-0.172	8.189
Race1	-0.112	2.868	-0.024	0.600	-0.222	5.890	-0.514	14.019
Rena	-0.156	3.672	-0.022	0.514	-0.467	10.917	-1.031	25.585
Akko & Kayo	-0.298	6.444	-0.091	1.944	-0.278	5.852	-0.735	16.260
Breakdancers	-0.229	13.688	-0.039	2.301	-0.154	9.044	-0.268	15.314
Uli	-0.424	12.400	-0.149	4.234	-0.084	2.242	-0.202	5.521
Poznan_Hall2	-0.112	7.781	-0.027	3.137	-0.042	3.242	-0.140	8.780
GT_Fly	-0.134	7.273	-0.107	5.614	-0.232	12.886	-0.276	14.697
Average	-0.193	7.212	-0.067	2.713	-0.171	5.849	-0.389	12.763

**Fig. 6.6** Complexity adaptation for MVC for changing battery levels

6.2.1.2 Comparing the Energy-Aware Complexity Adaptation to the State of the Art

The evaluation of the energy-aware complexity adaptation algorithm was done by experimentation on a battery-powered HP laptop (DV6000, Core-2 Duo). For accessing the battery level, we have used the *CallNiPowerInformation* windows API. In this experiment, the *Quality States* were forced to switch from *QS1* to *QS4* (simulating a battery discharge) and from *QS4* back to *QS1* (simulating battery charging). Figure 6.6 shows the frame-wise quality and time savings of our scheme encoding the *Ballroom* sequence. Two views are presented in Fig. 6.6. Compared to the RDO-MD, the *Quality States* *QS1* and *QS2* incur a negligible quality loss while providing a TS of up to 75 %. For *QS3* and *QS4* the TS go up to 79 % and 88 %, respectively.

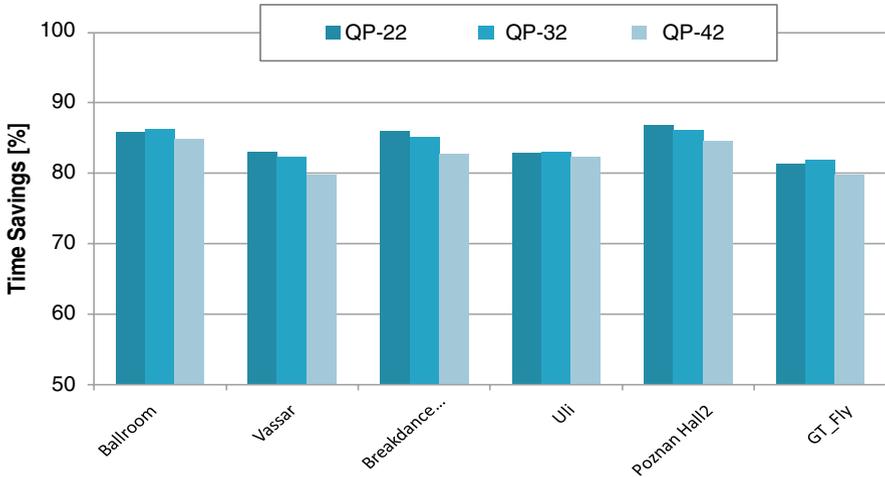


Fig. 6.7 Complexity reduction for the fast ME/DE

respectively. Due to the binocular suppression $QS3$ maintains a negligible PSNR loss. The resulting quality for the resulting viewpoint (VP) is measured according to Eq. (6.4) (Ozbek et al. 2007):

$$PSNR^{VP} = (1 - \alpha) \times PSNR^{HighQuality} + \alpha \times PSNR^{LowQuality}; \alpha = 1/3. \quad (6.4)$$

The energy-aware complexity adaptation for MVC enables run-time trade-off between complexity and video quality using different *Quality–Complexity Classes* (QCCs). Our scheme facilitates encoding of even and odd views using different QCCs (i.e., asymmetric view encoding) such that the overall perceived video quality is close to that of the high-quality view. Our scheme is especially beneficial for next-generation battery-operated mobile devices with a support of 3D-multimedia.

6.2.1.3 Comparing the Fast Motion and Disparity Estimation to the State of the Art

The comparison of our Fast ME/DE with the *TZ Search* algorithm and state-of-the-art complexity reduction schemes for ME/DE is presented in this section. Figure 6.7 shows the time savings of our Fast ME/DE algorithms for multiple video sequences and QPs for 4-view sequences. The *TZ Search* is used for comparison as it is 23× faster compared to the *Full Search*, while providing the similar rate-distortion (RD) results. Compared to the *TZ Search*, our fast ME/DE provides 83 % execution time saving. In the best case, the execution time savings go up to 86 %, which represents a significant computation reduction. These results are possible through drastic reduction in the number of SAD operations required, as shown in Fig. 6.8.

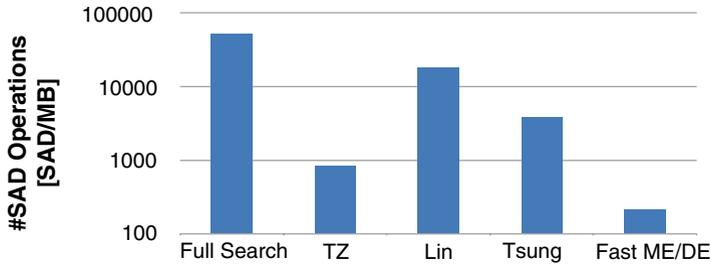


Fig. 6.8 Average number of SAD operations

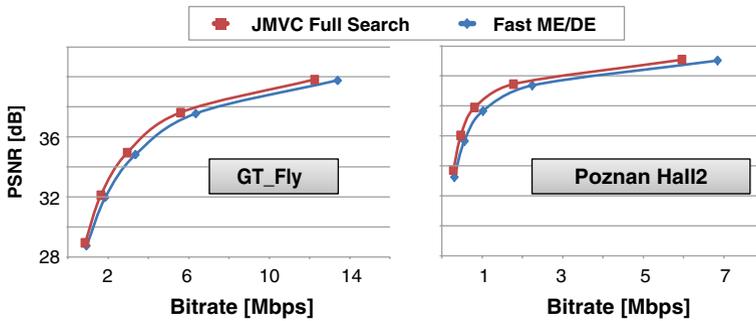


Fig. 6.9 Fast ME/DE RD curves

Compared to (Lin et al. 2008) and (Tsung et al. 2009), the number of SAD operations is reduced in 99 % and 94 %, respectively. It also represents 86 % complexity reduction compared to the original *TZ Search*.

The Fast ME/DE algorithm was designed to avoid high-quality drops and bitrate increases that surpass 10 %; for this reason, it does not result in high rate-distortion losses. The RD curves presented in Fig. 6.9 summarize the average 0.116 dB quality reduction and 10.6 % bitrate increase (see detailed table in Sect. 4.4.2) resulting from the aggressive complexity reduction provided by the proposed algorithm. Also, the Fast ME/DE demonstrate its robustness in terms of complexity reduction for all the tested video resolutions and QPs. This characteristic is desirable for real-time hardware architectures design.

6.2.2 Video Quality Control Algorithms

To deal with the quality losses posed by our fast algorithms, we propose, in Sect. 4.5, a complete rate control (RC) solution in order to efficiently manage the video quality vs. energy trade-off. An efficient RC is supposed to sustain the bitrate as close as possible to the target bitrate (optimizing the bandwidth usage) while avoiding sudden bitrate variations. To measure the RC accuracy, that is, how close the

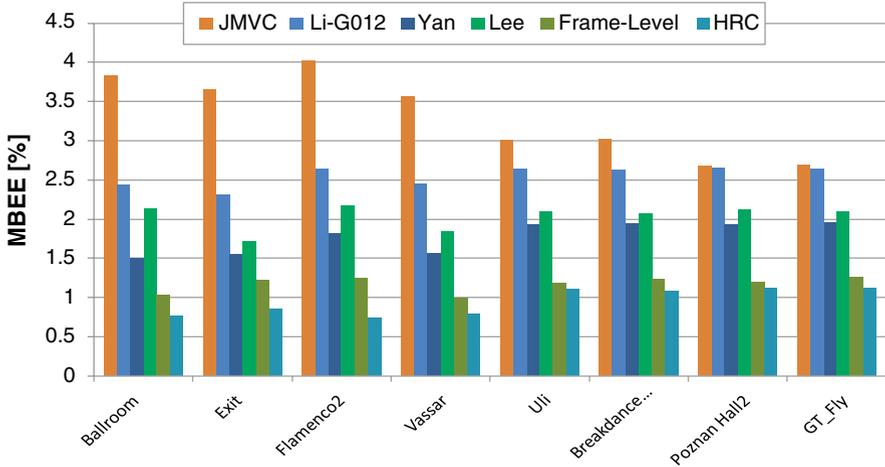


Fig. 6.10 Bitrate prediction accuracy

actual generated bitrate (R_a) is in relation to the target bitrate (R_t), we use the Mean Bit Estimation Error (MBEE) [see (4.2)] metric. The average is calculated over all Basic Units (N_{BU}) along 8 views and 13 GOPs for each video sequence.

Figure 6.10 presents the accuracy in terms of MBEE (less is better) for our HRC compared to the state-of-the-art solutions (Li et al. 2003; Yan et al. 2009a; Lee and Lai 2011), and our frame-level RC. On average, our Hierarchical Rate Control provides 0.95 % MBEE, while ranging from 0.7 % to 1.37 %. The competitors (Li et al. 2003; Yan et al. 2009a; Lee and Lai 2011) and the frame-level RC present, on average, 2.55 %, 1.78 %, 2.03 %, and 1.18 %, respectively. The HRC reduces the state-of-the-art error on 0.83 %, on average. The superior accuracy is a result of the ability to adapt the QP jointly at frame and BU levels while considering the 3D-Neighborhood correlation and the video content properties.

In Fig. 6.11 the long-term behavior of distinct Rate Control schemes is presented in terms of accumulated bitrate. A more accurate RC maximizes the use of available bandwidth and, consequently, the accumulated bitrate tends to stay closer to the target bitrate line. After a few initial GOPs required for control stabilization, our HRC curve better fits to the target bitrate followed by our frame-level RC, as shown in Fig. 6.11. JMVC without RC presents the worst bandwidth usage, as expected.

Once the accuracy of our HRC is proven we present the rate-distortion (RD) results to show that overall video quality and quality smoothness are not compromised. Table 6.5 summarizes the objective rate-distortion in terms of BD-PSNR (Bjontegaard Delta PSNR) and BD-BR (Bjontegaard Delta Bitrate) in relation to JMVC without RC. The HRC provides 1.86 dB BD-PSNR increase along with BD-BR reduction of 40.05 %, on average. If compared to the work of Lee and Lai (2011), which presents the best RD performance among the related works, the HRC delivers 0.06 dB increased BD-PSNR and 3.18 % reduced BD-BR. Remember, besides superior RD performance, the HRC also outperforms (Lee and Lai 2011) in terms of accuracy (1.08 % MBEE).

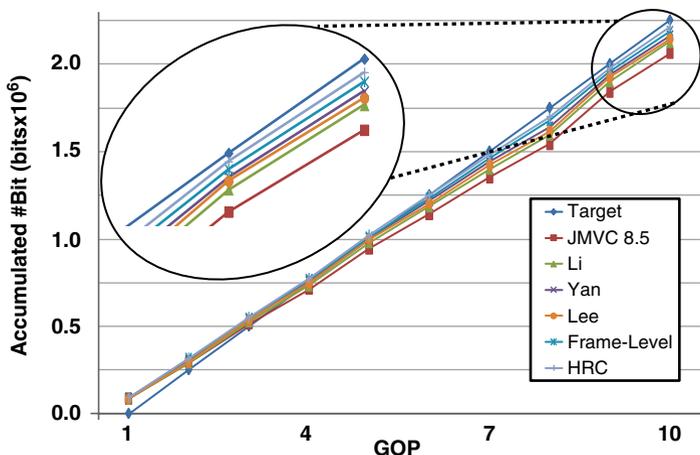


Fig. 6.11 Accumulated bitrate along the time

Table 6.5 Bjøntegaard PSNR and BR for the HRC

JMVC vs.	VGA				XGA			HD1080p		
	Ballroom	Exit	Flamenco2	Vassar	Bdancer	Uli	Poznan	GT_Fly	AVG	
Li	BD-PSNR	0.328	0.368	0.217	0.183	0.215	0.208	0.254	0.012	0.223
	BD-BR	-9.831	-10.348	-8.784	-6.116	-8.963	-9.805	-12.186	-6.711	-9.093
Yan	BD-PSNR	-0.090	0.073	0.114	0.051	-0.086	0.155	0.169	-0.118	0.034
	BD-BR	-4.156	-5.463	-3.346	-1.958	22.819	-5.293	-0.671	16.953	2.361
Lee	BD-PSNR	2.056	2.058	1.292	1.509	2.019	1.879	1.928	1.721	1.808
	BD-BR	-35.446	-43.167	-26.643	-33.474	-43.445	-39.110	-40.931	-38.134	-37.544
Frame-Level	BD-PSNR	0.939	1.089	0.880	0.596	0.881	0.670	0.750	0.614	0.802
	BD-BR	-22.241	-26.965	-22.989	-16.897	-22.818	-20.964	-17.184	-20.872	-21.366
HRC	BD-PSNR	1.585	2.375	2.103	1.176	2.060	1.870	2.086	2.056	1.914
	BD-BR	-31.588	-47.458	-38.199	-27.335	-46.112	-49.660	-48.766	-47.258	-42.047

Figure 6.12 shows the RD curves for different video sequences considering videos from distinct spatial, temporal, and disparity indexes. The HRC shows its superiority in relation to the state of the art for most of the RD curves. It is also important to highlight that HRC does not insert visual artifacts such as blurring and blocking noise. Moreover, our RC does not compromise the borders sharpness typically lost in case of bad QP selection.

6.2.3 Energy-Efficient Hardware Architectures

Results and comparison to the state of the art for the proposed energy-efficient hardware architecture are presented in this section. Table 6.6 summarizes the hardware implementation results with details on gate count, size of on-chip memory,

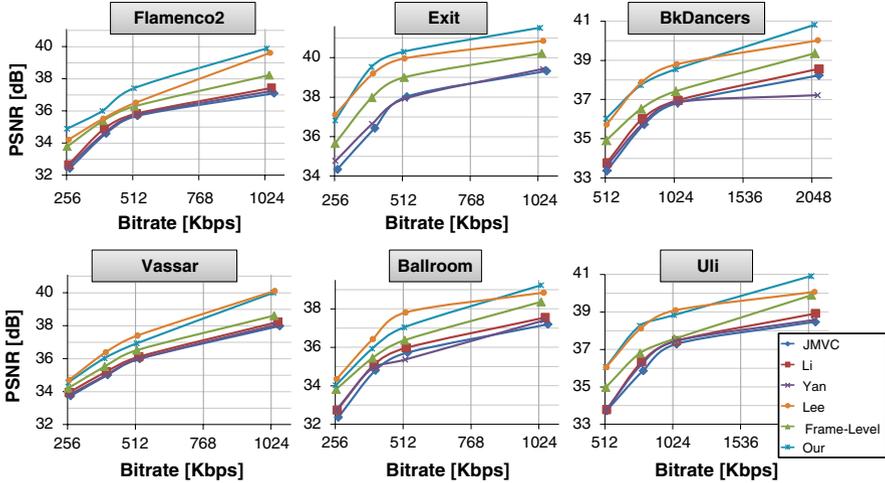


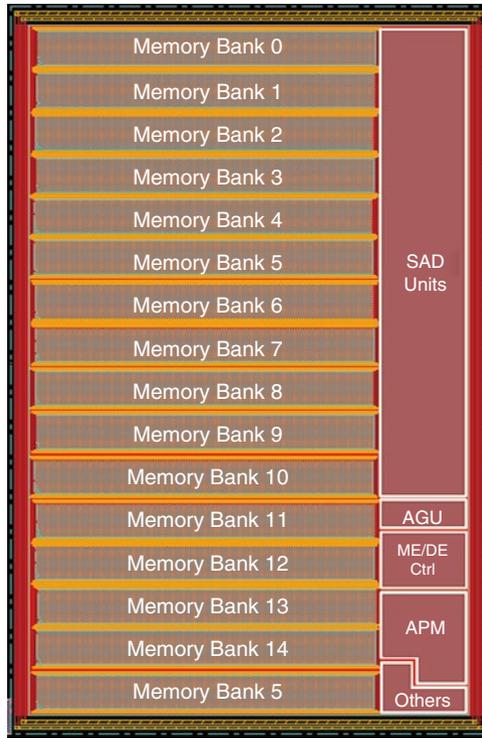
Fig. 6.12 Rate-distortion results for the HRC

Table 6.6 Motion and disparity estimation hardware architectures comparison

	Chang et al. (2010)	Tsung et al. (2009)	Zatt et al. (2011c)	ME/DE Architecture w/dynamic search window (B)
Technology	UMC 90 nm	TSMC 90 nm Low power LowK Cu	IBM 65 nm LPe LowK	ST 65 nm LP 7 metal layer
Gate count	562 k	230 k	211 k	102 k
SRAM	170 Kbits	64 Kbits	737 Kbits	512 Kbits
Max. frequency	95 MHz	300 MHz	300 MHz	300 MHz
Power	n/a	265 mW, 1.2 V	81 mW, 0.8 V	74 mW, 1.0 V
Search range		[±16, ±16]	[±64, ±64]	[±96, ±96]
Proc. capability	CIF @ 42fps	4-views 720p	4-views HD1080p	4-views HD1080p

performance, and power consumption (on-chip). Figure 6.13 shows the physical layout of the ASIC implementing architecture. To the current stage, the IC was completely synthesized down to physical level but not fabricated.

Compared to our architecture, the one presented in Chang et al. (2010) requires more hardware resources while providing significantly low performance, attending only CIF resolution (352×258) requirements. Even assuming a frequency extrapolation, the performance provided by Chang et al. (2010) is not comparable to the other discussed solutions. Comparing to (Tsung et al. 2009), our design is able to provide real-time ME/DE for up to 4-views HD1080p videos compared to HD720p provided by Tsung et al. (2009) at the same operation frequency. This represents a throughput increase (in terms of the processed MBs) of 2.26× obtained through Fast ME/DE and careful pipelining and scheduling in architecture. Additionally, our architecture reduces the gate count and power consumption compared to Tsung et al. (2009).



SAD Units:	Sum of Absolute Differences Operators
ME/DE Ctrl:	Motion/Disparity Estimation Control
AGU:	Address Generation Unit
DPM:	Dynamic Power Management

Fig. 6.13 ME/DE architecture with application-aware power gating physical layout

The proposed architecture brings improvements in terms of control flow and memory design. Our solution reduces the gate count, number of memory bits, and power consumption by 52 %, 9 %, and 30 %, respectively, compared to Zatt et al. (2011c). It also increases the maximum search range from $[\pm 64, \pm 64]$ to $[\pm 96, \pm 96]$. Compared to Tsung et al. (2009) the area and power reductions are 66 % and 72 %, respectively, while providing higher throughput. This significant power reduction is mainly due to the employment of dynamic search window formation, on-chip memory power gating, smaller logic, and fast ME/DE scheme. Note that the standard-cell library and fabrication technologies are different compared to (Tsung et al. 2009; Zatt et al. 2011c). Our architecture is implemented in 65 nm at 0.8 V while Tsung et al. (2009) use a 90 nm low-power technology at 1.8 V. To the best of our knowledge, the Motion and Disparity Estimation HW Architecture with Application-Aware Power Gating represents the most efficient architectural solution available in the current literature and guarantees the processing of 4-view HD1080p running at 300 MHz and dissipation 74 mW.

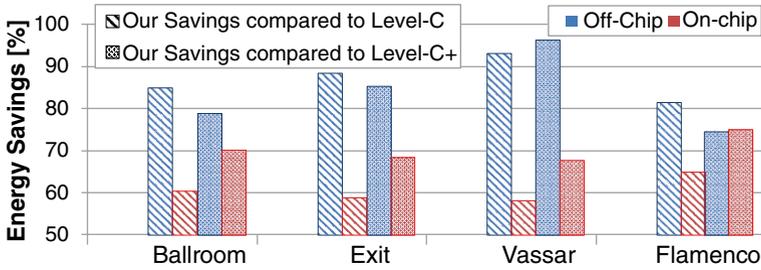


Fig. 6.14 Memory-related energy savings employing dynamic search window technique

At first analysis, the main drawback of the proposed ME/DE architecture lies in the increase on-chip video memory in comparison to the state of the art. The on-chip memory in our hardware is relatively larger as it supports a much bigger search window of up to $[\pm 96, \pm 96]$ compared to $[\pm 16, \pm 16]$ in Tsung et al. (2009) (which is insufficient to capture larger disparity vectors). However, the larger on-chip memory does not imply an increased power dissipation because of the dynamic power management and power-gating techniques employed in our solutions. Compared to Zatt et al. (2011c) an on-chip memory reduction of about 30% is achieved.

The authors of Tsung et al. (2009) use a rectangular data reuse technique such as Level-C (Chen et al. 2006), which compared to our proposed solutions (search as dynamic window formation) perform inefficiently. Note, Level-C (Chen et al. 2006) with a search window of $[\pm 96, \pm 96]$ would require 4 memories of 288 Kb (i.e., a total of 1,115 Mb) to exploit the reusability in four possible prediction directions available in MVC. Our approach implements it with 512 Kbits. To perform a fair comparison, we have deployed the Level-C and Level-C+ (Chen et al. 2006) techniques in our hardware architecture.

Figure 6.14 shows the energy benefit of employing our dynamically expanding search window and multibank on-chip memory with power gating. Compared to Level-C and Level-C+ (Chen et al. 2006) prefetching techniques (based on rectangular search windows), our approach presents energy reduction in on-chip and off-chip memories as shown in Fig. 6.14. For a search window of $[\pm 96, \pm 96]$, our approach provides an energy reduction of up to 82–96 % and 57–75 % for off-chip and on-chip memory access, respectively. These significant energy savings are due to the fact that Level-C and Level-C+ (Chen et al. 2006) suffer from a high data retransmission for every first MB in the row. Additionally, our approach provides higher data reuse and incurs reduced leakage due to a smaller on-chip memory and power gating of the unused sectors.

6.3 Summary of Results and Comparison

To cope with comparison fairness issues, along this chapter were detailed all settings and videos used for comparison along this monograph. The video benchmark sequences were classified using the spatial, temporal, and disparity indexes

to quantify their complexity along these domains. Additionally, the simulation infrastructure and tools employed along this work were presented in Sect. 6.1.

A complete comparison to the state of the art was presented in Sect. 6.2 and showed the superiority of our solutions. The mode decision complexity reduction algorithms are able to provide average 71 % complexity reduction with 0.17 dB quality loss or 5.8 % bitrate increase (Bjøntegaard). In turn, the Fast ME/DE contributes with additional 83 % complexity reduction with a drawback of 0.116 dB quality loss and 10.6 % bitrate increase. The energy/complexity versus quality trade-off can be managed using the presented complexity adaptation algorithm whose stability and fast reaction to changing scenarios were demonstrated in section “Comparing the Energy-Aware Complexity Adaptation to the State-of-the-Art”. The quality drawback posed by our algorithms may be recovered by employing our HRC that provides an average video quality increase of 1.9 dB (Bjøntegaard).

Compared to the state-of-the-art architectures for ME/DE the solution presented in this monograph is able to reduce the gate count in 56 % while increasing the performance 2.26x. Most importantly, our architecture is able to provide real-time 4-views HD1080p at 300 MHz with 74 mW. It represents a 73 % power reduction compared to Tsung et al. (2009), the most prominent related work. This reduction is mainly due to intelligent on-chip video memory power gating. The chip physical layout is showed in Sect. 6.2.3.

The results presented in this chapter for both energy-efficient algorithms and energy-efficient architectures demonstrate the superior performance of our solutions in face of the related works. Moreover, the results demonstrate that is possible to provide solutions able to encode MVC at real time while respecting embedded devices energy constraints.

Chapter 7

Conclusion and Future Works

The presented monograph focuses on the energy reduction of the Multiview Video Coding (MVC) encoder to enable the realization of real-time high-definition 3D-video encoding running on mobile embedded devices with battery-constrained energy. For that, novel energy-efficient techniques are proposed at both algorithmic and architectural abstraction levels. The joint consideration of algorithms and underlying hardware architecture is the key enabler to provide improved energy efficiency, as demonstrated along this monograph.

The strong correlation within the 3D-Neighborhood domain, concept defined in this work, has been the base for designing most of the algorithms and hardware architecture adaptation schemes proposed. An extensive study based on statistical analysis correlating MVC coding side information (such as coding modes, motion/disparity vectors, and RDCost) to the video content properties is provided to justify the importance of the 3D-Neighborhood understanding and to demonstrate its potential to support energy reduction in the MVC encoder.

A set of *energy-efficient algorithms for MVC* compose one of the major contributions to the state of the art proposed in this work. A multilevel fast mode decision algorithm with 6 levels is described targeting energy efficiency through complexity reduction. The Early SKIP prediction, one stage of our scheme, exploits the high occurrence of SKIP coded MBs to accelerate the encoding process by employing statistical methods that define if each MB is in the high SKIP probability region in order to avoid other coding modes evaluation. Our algorithm eliminates coding modes evaluation even in the case where an early SKIP is not detected by analyzing the coding modes available within the 3D-Neighborhood while considering a video/RDCost-based mode ranking. The video properties are also used to define block sizes and prediction modes orientation. To protect the multilevel fast MD algorithm from inserting excessive quality losses an early termination test is inserted between each prediction step. This algorithm defines QP-based thresholds for two distinct energy reduction strengths: the *relax* and *aggressive* strengths. By employing two operation modes it is possible to select the best energy vs. quality trade-off for a given system state and video content. Moreover, multiple fast MD modes enable the

integration of an energy-aware complexity-adaptation control scheme. The multi-level fast MD algorithm evaluation, results, and comparison demonstrate a complexity reduction of up to 79% at the cost of 0.32 dB quality loss and 10 % bitrate increase, for *aggressive* mode, and 0.1 dB quality loss and 3 % bitrate increase, for *relax* mode when compared to RDO-MD.

This work demonstrated that the coding properties and coding effort highly depend on the video content. Moreover, when considering embedded applications, the processing power is constrained by energy resources available in the embedded battery. From this observation, it is proposed an energy-aware complexity-adaptation algorithm. The goal is to jointly consider the video input characteristics and the battery state to sustain the highest possible video quality by selecting the appropriated MD algorithm and quality states. In case of battery discharging, further energy reduction is necessary leading to quality reduction. Thus, the complexity-adaptation algorithm delivers graceful quality degradation by employing the binocular suppression theory knowledge. For binocular displaying, the Human Visual System (HVS) tends to perceive the highest quality view, so the proposed algorithms firstly drop the quality of odd views guaranteeing a high perceived quality while reducing energy consumption for encoding these odd views. Experimental results show the beneficial effect of the complexity adaptation for energy consumption and smooth quality variation along the time under battery charging and discharging scenarios.

The motion and disparity estimation consumes more than 90 % of the total MVC encoding energy and represents the main target for energy reduction. In this work, a novel Fast ME/DE was detailed. It uses the motion and disparity vectors available in the 3D-Neighborhood to avoid, for multiple frames in the prediction structure, the complete motion/disparity search pattern. There are defined two classes for frames, key and non/key frames, where the key frames are encoded using off-the-shelf fast search patterns and the non/key frames employ our Fast ME/DE. According to the confidence, defined using image properties, on the vectors inferred from the neighborhood, each MB in the non/key frames selects between fast mode or the ultra-fast mode. These modes test only 3 or 13 candidate blocks, respectively. The proposed Fast ME/DE algorithms are able to reduce 83 % of the total encoding time at the cost of 0.116 dB and 10 % bitrate increase.

To compensate eventual losses posed by the energy-efficient algorithms, a video quality management based on our hierarchical rate control (HRC) algorithm was proposed. The HRC operates in two actuation levels, the frame level and the basic unit (BU) level, and features a coupled closed feedback loop. The frame-level RC employs a Model Predictive Controller (MPC) to predict the bitrate for future frames based on the bit allocation in the frames belonging to the 3D-Neighborhood. The multiple stimuli coming from temporal, disparity, and phase neighboring frames compose the MPC input. The bitrate prediction is then used to define the optimal QP for that frame. The QP is further refined inside the frame by a Markov Decision Process (MDP)-based BU-level rate control. It considers Regions of Interest to prioritize hard-to-encode image regions. Reinforcement learning is used to update the MDP parameters. The HRC provides smooth bitrate and video quality variations along time and view axes, while respecting bandwidth constraints and providing improved video quality. Compared to the fixed QP solution, the video quality was improved in 1.9 dB (Bjontegaard). In comparison to the state of the art,

the bitrate prediction error is reduced in 0.83 % in addition to 0.106 dB PSNR increase or 4.5 % Bjontegaard bitrate reduction.

In addition to the energy-efficient algorithms, the severe energy restrictions and performance requirements demanded by the MVC encoder require hardware dedicated acceleration able to employ sophisticated application-aware power-gating techniques. One *energy-efficient hardware architecture* for motion and disparity estimation was proposed in order to provide throughput to encode, at real time, 4-view HD1080p video sequences.

It was proposed a novel ME/DE architecture that incorporates a multibank video on-chip memory and the dynamic search window-based data prefetching technique for jointly reducing the on/off-chip memory energy consumption. A dynamically expanding search window is constructed at run time based on the neighborhood-extracted search map to reduce the off-chip memory accesses. Considering the multistage processing nature of advanced fast ME/DE schemes, the reduced-size multibank on-chip memory is partitioned in multiple sectors which can be power-gated depending upon the video properties while enabling fine-grained tuning for leakage current reduction. A novel processing scheduling was designed exploiting the multiple parallelism levels available in the MVC coding structure, view, frame, reference frame, and MB levels, to deal with data dependencies.

The architectural contribution presented in this monograph involves the architecture design, management schemes, complete RTL (VHDL) coding, and ASIC synthesis down to physical layer using 65-nm fabrication technology. From experimental results for multiple video sequences, the proposed architectures provide a dynamic energy reduction of 82–96 % for the off-chip memory and up to 80 % on-chip leakage energy reduction compared to state of the art. From this contribution, it is possible to demonstrate the feasibility of performing motion and disparity estimation for up to 4-view HD1080p at 30 fps with a power dissipation of 74 mW running at 300 MHz on an IC footprint with 102 k gates.

The overall results and benchmarks demonstrate the energetic efficiency of the proposed algorithms and architectures in front of the state-of-the-art solutions. This supports one of our claims: for attending the 3D-video coding requirements for embedded systems, it is required to jointly consider and optimize the coding algorithms and the underlying dedicated hardware architectures. Additionally, run-time adaptation is required to better predict the system behavior and react to changing video input, coding parameters, and battery-level scenarios. For that, deep MVC application knowledge coming from extensive analysis, such as the correlation available within the 3D-Neighborhood, must be employed.

7.1 Future Works

Beyond the contribution brought in this monograph work, there are multiple research topics related to 3D-video coding and video processing that were not addressed in this volume. The algorithms and architectures here presented were centered in mode decision and motion and disparity estimation once these are the

most energy-hungry coding units in the MVC encoder. Additionally, focusing on video quality issues the rate control was discussed. The MVC, however, brings a big set of other research challenges if embedded applications are considered. 3D-video pre- and post-processing also play key roles in the 3D-video system and present plenty of novel challenges. Finally, next-generation 3D-video coding algorithms are under study for future standardization. The next 3D-video generation is expected to bring innovative tools and provide good perspective for future research opportunities in the 3D-multimedia field.

7.1.1 Remaining MVC Challenges

Although the main challenges in terms of complexity and energy consumption are related to the MD and ME/DE blocks, attending to the MVC demands while respecting energy constraints presents challenges related to other MVC processing blocks. The entropy encoder, for instance, may become the bottleneck of the encoder system if no proper parallelization is employed. The block-level data dependencies in intra prediction also require research attention. Finding efficient solutions to deal with data dependencies and parallelization issues provide interesting research opportunities for future works.

7.1.2 3D-Video Pre- and Post-processing

Video encoding is one single stage in the 3D-video system. Between video capturing and video coding phases, there is a need for preprocessing such as geometrical calibration (for correcting the alignment of the multiple videos) and color correction (responsible for equalizing the brightness level and color gamut). After the transmission and decoding, the video is processed for displaying depending on the application and display technology. This post-processing phase includes color space mapping (in a system using color polarization), resolution scaling, and viewpoint synthesis (generation of intermediate viewpoints for displaying). The pre- and post-processing implement complex and data-intensive algorithms (especially for viewpoint synthesis) that run concurrently with the video encoder/decoder and require real-time performance. Therefore, the embedded energy and hardware resources must be shared to attend both video coding and pre/post-processing demands.

7.1.3 Next-Generation 3D-Video Coding

The next generation for 3D-video coding is currently referred as 3DV (3D-Video) (ISO/IEC 2009) and is based on the Video+Depth concept that defines distinct

channels to transmit video and the depth maps. The 3DV is expected to be defined as an extension to the HEVC/H.265 (Sullivan and Ohm 2010). The 3DV tools will bring a completely new set of challenges boosting the research topics related to 3D-multimedia. Moreover, the video coding standards lifetime is expected to reduce for future standard generations resulting in the simultaneous coexistence of multiple coding standards. Thus, there is a need to support multiple complex coding standards in the same device by employing flexible and adaptive solutions.

Appendix A

JMVC Simulation Environment

The JVT (Joint Video Team), formed from the cooperation between the ITU-T Study Group 16 (VCEG) and ISO/IEC Motion Picture Experts group (MPEG), responsible for the standardization of the H.264, SVC (scalable video coding), and MVC provides software models used for algorithms experimentation and for standards proof of concept. The JMVC (Joint Model for MVC), currently on version 8.5, is the reference software available for experimentation on the MVC standard. Along this work the JMVC software, coded using C++, was used and modified to implement the proposed algorithms. Initially, the version 6.0 was used followed by an upgraded to version 8.5. Considering the length and complexity of the software, a high-level overview of the interaction between the main encoder classes is presented here. Afterwards are shown the classes modified to enable our algorithms experimentation. For in-depth details of the class structure refers to JMVC documentation .

A.1 JMVC Encoder Overview

The JMVC classes are hierarchically structured as shown in Fig. A.1. The JMVC encodes each view at a time requiring as many calls as number of view to be encoded. The reference views are stored in temporary files. The class *H264AVCEncoderTest* represents the top encoder entity; it initializes the encoder, and calls the *CreateH264AVCEncoder* class to initialize the other coding classes. At this level, the *PicEncoder* is initialized and the frame-level loop is controlled. The *PicEncoder* loops over the slices inside each frame and resets the RDCost. The slice encoder controls the MB-level loop and sets the reference frames for each slices. For MB encoding there are two main classes: the *MbEncoder* and the *MBCoder*. *MbEncoder* encapsulates all the prediction, transforms, and entropy steps. It implements the mode decision by looping over and encoding all possible coding modes (in case of RDO-MD) and determining the minimum RDCost. At this point no MB

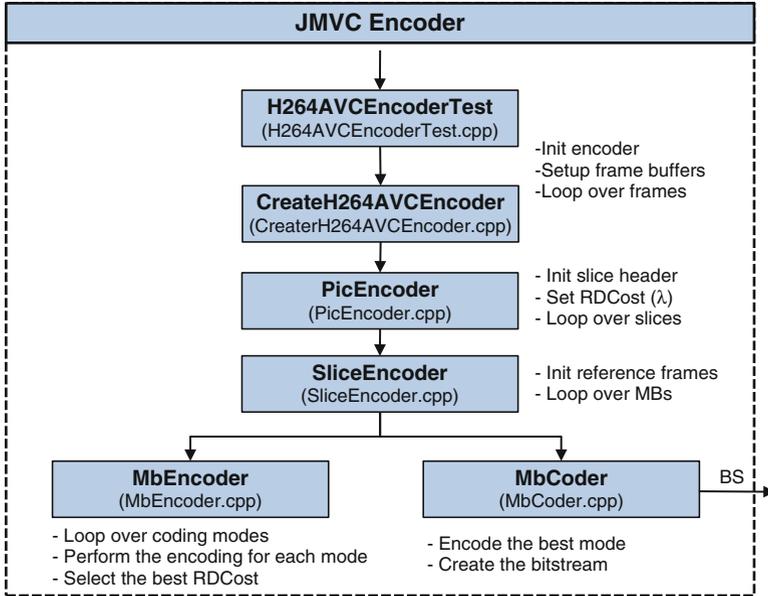


Fig. A.1 JMVC encoder high-level diagram

coding data is written to the bitstream. Once the best mode is selected, the *SliceEncoder* calls the *MbCoder* to write the MB-level side information and residues to the bitstream.

Figure A.2 (Tech et al. 2010) depicts the hierarchical call graph of methods inside the mode decision process implemented in *MbEncoder* class. Firstly, the SKIP and Direct modes are evaluated,;along this monograph these modes are jointly referred as SKIP MBs. In the following, all inter-prediction block sizes are evaluated. For each partition size a call to the method *MotionEstimation::estimateBlockWithStart* (see Fig. A.3 discussion) is performed. The same happens for the sub-partitions in case of 8×8 partitioning. *EstimateMb8x8Frect* is only called in case the FRect flag is set. Finally, the intra-frame coding modes including PCM, intra 4×4 , intra 8×8 (FRect only), and intra 16×16 are called. The *Estimate<mode>* methods call the complete coding loop for that specific mode including prediction, transforms, quantization, entropy encoding, and reconstruction. It allows a precise definition of the minimum RDCost (λ) and an optimal best mode selection at the cost of elevated coding complexity. The *MbCoder* is called to entropy encode the best mode and write the data into the bitstream output buffer.

The motion and disparity estimation search itself is defined in the method *estimateBlockWithStart* and is composed of three basic steps. The ME/DE dataflow is represented by the arrows in Fig. A.3. Once the *estimateBlockWithStart* is called, for instance in *EstimateMb16x16*, the search runs for each reference frame list

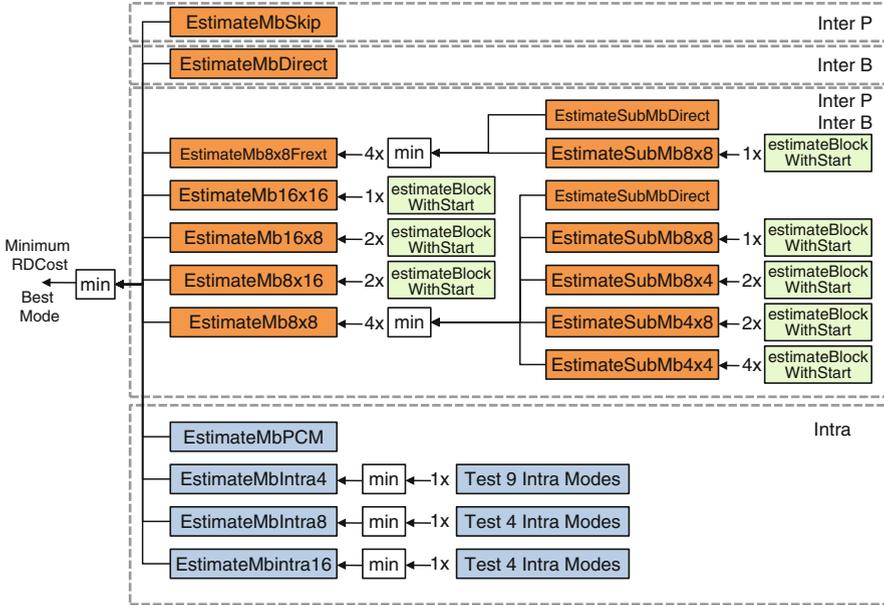


Fig. A.2 Mode decision hierarchy in JMVC

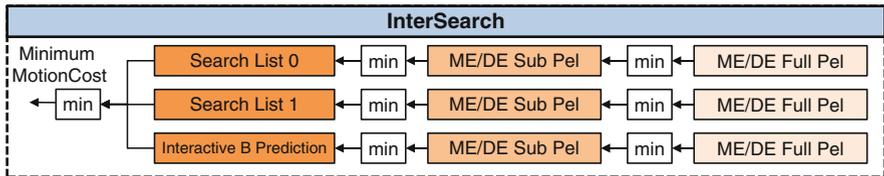


Fig. A.3 Inter-frame search in JMVC

(List 0 and List 1) and for an interactive B search mode that exploits both lists in an interactive fashion (in case the interactive B is active). At the software perspective, there is no distinction between ME and DE. List 0 and List 1 store both temporal and disparity reference frames. The search for a given reference frame firstly finds the best candidate block among the integer pixels (ME/DE Full Pel) and then refines the result considering half and quarter pixels (ME/DE Sub Pel). The search pattern depends on the search algorithm, and JMVC implements TZ Search, Full Search, Spiral Search, and Log Search. The goal is to find the candidate block that minimizes the Motion Cost (λ_{Motion}) in terms of SAD, SAD-YUV (considering chroma channels), SATD, or SSE according to user-defined coding parameters. The position of the best matching candidate block position defines the motion or disparity vector.

A.2 Modifications to the JMVC Encoder

A.2.1 JMVC Encoder Tracing

In order to generate the statistics used for coding modes and motion/disparity vector, some modifications were done in the original JMVC code. The point selected for this tracing is inside the entropy encoder to guarantee that the extracted data is the same actually encoded and transmitted. The entropy encoder is declared as the virtual class *MbSymbolWriterf*, but the actual implementation is in *CabacWriter* and *UvlcWriter*, depending on the entropy encoder selected in the configuration file. The methods monitored are *skipFlag* that encodes the SKIP (and Direct)-coded MBs and *mbMode* that encodes all other modes. Note, MB coding mode codes (*uiMbMode*) vary with the slice type as defined in Tables 7-11, 7-12, 7-12, 7-13, and 7-14 of the MVC standard (JVT 2008).

A.2.2 Communication Channels in JMVC

Multiple algorithms proposed in this monograph employ the information from the 3D-neighborhood. For that, there is a need to build communication channels between neighboring MBs in the special, temporal, and disparity domains. In other words, an infrastructure to send and receive data at MB level, at frame level (in same view), and at view level (frames in different views). Therefore, a hierarchical communication infrastructure was designed and implemented. Figure A.4 presents graphically the modified classes along with the new member data structures and communication methods.

The *MbDataAccess* already provides direct access to the left and upper neighbors (A, B, C, and D in Fig. A.4). This access was extended to the right and bottom neighbors (A*, B*, C*, and D*) enabling access to data from all spatial neighboring

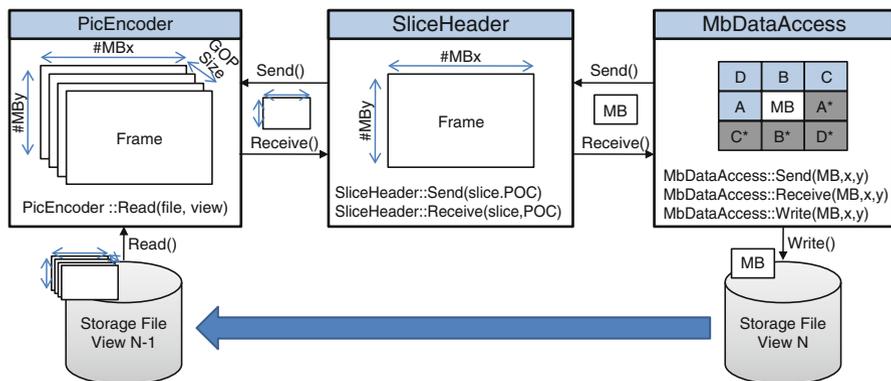


Fig. A.4 Communication in JMVC

MBs. For temporal neighboring MB access the current MB data is sent to *SliceHeader* (using *Send()* methods) where a 2D array stores the information from the MBs belonging to the current slice. Once the slice is completely processed, the 2D array is sent to the *PicEncoder* class. *PicEncoder* maintains the data for the whole current GOP. For reading the data communication channel writes the requested data from *PicEncoder* to *SliceHeader* and finally to *MbDataAccess* (using *Receive()* methods). As far as the views are processed in distinct encoder calls there is a need to use external temporary files to transmit the disparity neighboring information. The current MB data is written in these files from *MbDataAccess*, while the data from previous views is read in *PicEncoder*, as shown in Fig. A.4.

A.2.3 Mode Decision Modification in JMVC

The mode decision is programmed in a very simple becoming easy to find and modify. MD is handled in *MbEncoder::encodeMacroblock*. To find the exact point, search for the *xEstimateMb* methods responsible for calling the modes evaluation. Before this point are implemented the 3D-neighborhood communication calls and the calculation required to take the fast decisions.

A.2.4 ME/DE Modification in JMVC

The modifications for fast ME/DE are inserted in two distinct classes. For modifications at higher level such as avoiding interactive B search, search direction, and reference frames the modifications are done in the *MbEncoder* by modifying the *xEstimateMb* methods. If the modifications are in the search step itself, *MotionEstimation* class is the right point for modification. *estimateBlockWithStart* method is responsible for fetching the image data, prediction SKIP vectors, and calling the search methods (*xPelBlockSearch*, *xPelSpiralSearch*, *xPelLogSearch*, and *xTZSearch*). By modifying these methods it is possible to reach low-level modifications on the ME/DE search.

A.2.5 Rate Control Modification in JMVC

The JMVC does not implement any rate control algorithms. Therefore, to implement the hierarchical rate control (HRC) scheme one new class is created, the *RateControl*. Three files are used to better partition the RC hierarchy. File *RateCtlCore.cpp* describes the behavior of the whole HRC while *RateCtlMPC.cpp* and *RateCtlUB.cpp* are responsible for the calculations relative to the MPC and MDP controllers. *RateCtl.h* file is used to define the MPC and MDP actuation parameters. The QP

history is read from *CodingParameter* class and the generated bitrate is accessed via *BitWriteBuffer* and *BitCounter*. The QPs defined for the next frames or BU are sent back to *CodingParameter*. Additional modifications were required in files *MbCoder.cpp*, *CodingParameter.cpp*, *RateDistortion.cpp*, *ControlMngH264AVCEncoder.cpp*, *Multiview.cpp*, and *ControlMngH264AVCEncoder.h*.

Appendix B

Memory Access Analyzer Tool

The MVC Viewer software is used as part of this work to plot and analyze the memory accesses that are required by the motion and disparity estimation (ME/DE). The goal of this tool is to help the researchers in their projects in the visual and statistical analysis of the communication between the multiview video encoder and the reference samples memory. It provides a set of final statistics and several plots using the original input video.

The MVC Viewer was designed to be adapted to different encoder parameters. In a configuration file, the user should specify (a) the number of views, (b) the GOP size, (c) the video resolution, (d) the original YUV video files path, and, finally, (e) the memory tracing input files path. The tracing file is an intermediated way to communicate the video encoder output, like JM or $\times 264$, with the MVC Viewer tool. In this file, all memory accesses performed by ME/DE are listed.

This tool runs over the JVM (Java Virtual Machine) and provides a simple interface to the analysis. Figure B.1 presents the overview of the MVC Viewer main screen. The main parts are as follows:

1. Encoding parameters: GOP Size, number of coded frames, number of coded views, and video resolution (directly defined in the configuration files)
2. Tracing files path where all accessed regions of reference frames are listed
3. Original YUV videos
4. Program mode selection: the MVC Viewer has mainly two possible analysis tools:
 - (a) Current macroblock-based analysis
 - (b) Reference frame-based analysis
5. Listbox with all memory access that will be plotted in the output

The two analyses that are allowed by the MVC Viewer tool will be explained in the next sections.

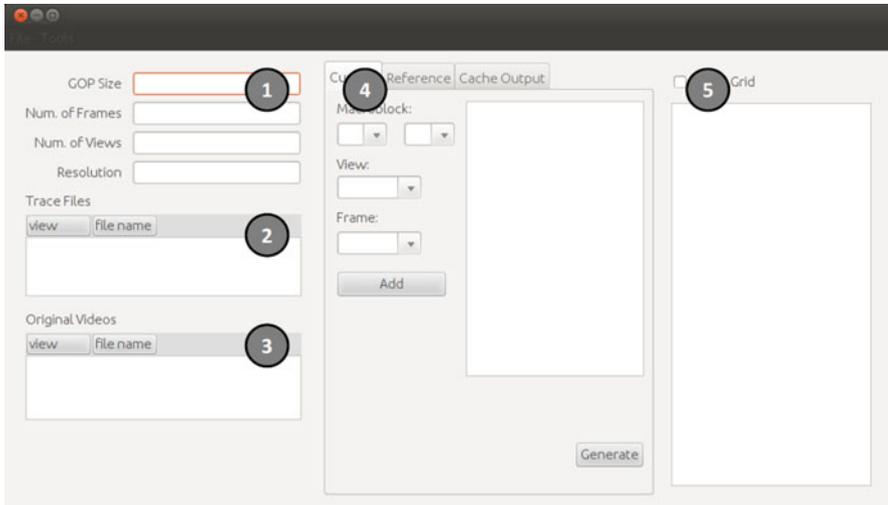


Fig. B.1 MVC viewer main screen

B.1 Current Macroblock-Based Analysis

In this analysis, the goal is to trace all accessed reference frame samples when the ME/DE is performed for one or more current macroblocks. The MVC massively uses multiple reference frames, and then the MVC Viewer will generate several plots that will determine the accessed regions for each reference frame (temporal and disparity neighbors). The Fig. B.2 shows a MVC Viewer screenshot when it is running this analysis. The main parts are as follows:

1. Selection of the target macroblocks that will be traced
2. List of all selected macroblocks
3. List of all memory access caused by the ME/DE for the selected macroblocks

Fig. B.3 presents an output example for one macroblock that reflects in samples accesses in the four directions: past and future temporal reference frames, and right and left disparity reference frames.

B.2 Search Window-Based Analysis

This analysis selects one specific frame and traces all accesses performed by the ME/DE when the selected frame is used as reference. This way, it is possible to determine the most accessed regions of the frame. The knowledge about this behavior is important to define strategies to save memory bandwidth. Figure B.4 presents the MVC Viewer during this analysis, where the main parts are as follows:

1. Reference frame selection: the user must define the frame identification (the view and frame positions) to be traced.

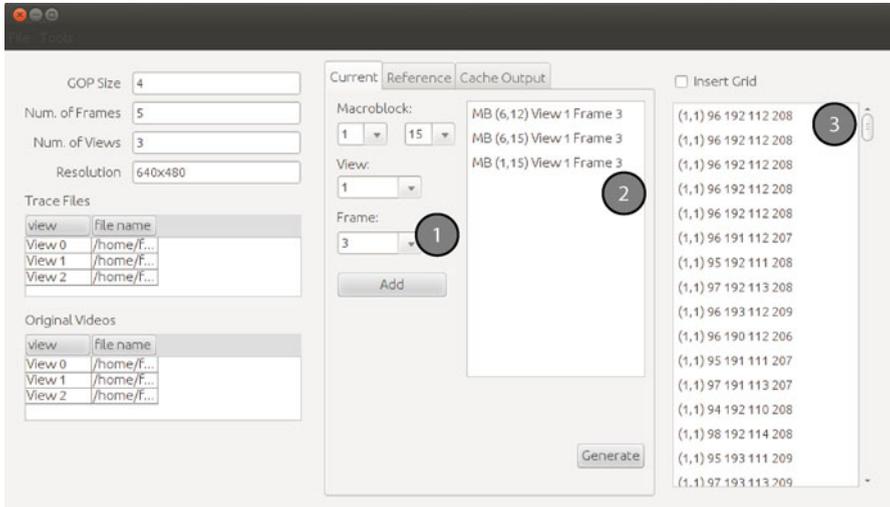


Fig. B.2 Current macroblock-based analysis screenshot

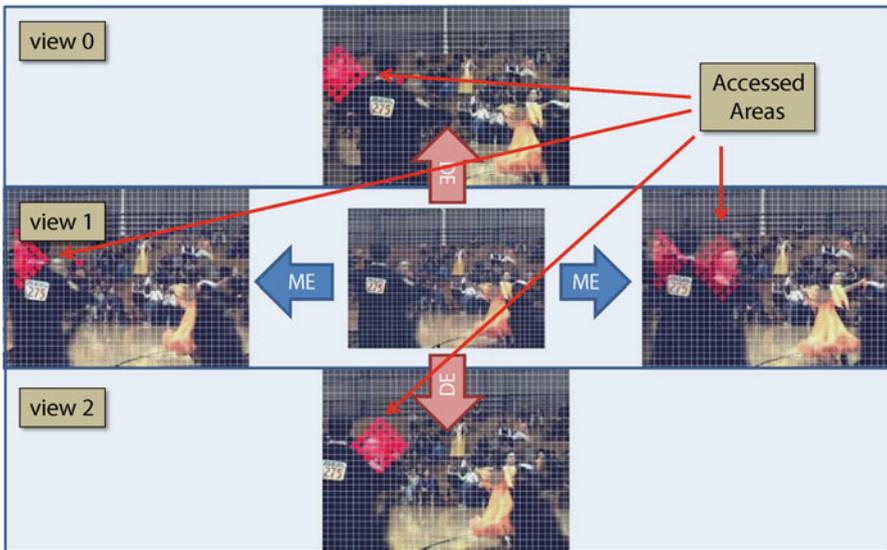


Fig. B.3 Output example: four prediction directions and their respective accessed areas

2. Current MBs tracing option: the user has the possibility to delimit an area inside the reference frame to discover which are the current blocks processed by the ME/DE that cause the accesses.
3. List of all memory access caused by the ME/DE in the selected reference frame.

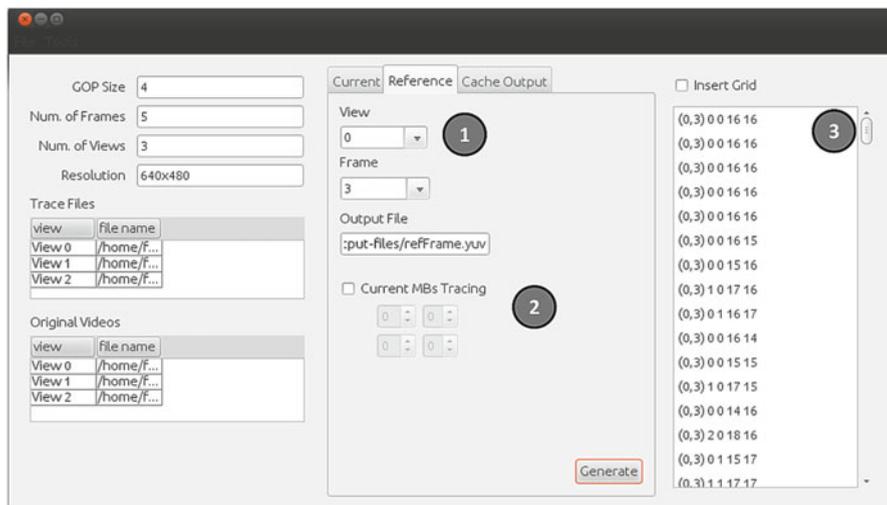


Fig. B.4 Current macroblock-based analysis screenshot



Fig. B.5 Output example: reference frame access index considering two block matching algorithms: full search and TZ search

The Fig. B.5 presents two different examples of the reference frame-based analysis considering two search algorithms: (a) full search and (b) TZ search.

The full search has a regular access pattern where all samples inside the search window are fetched. On the other hand, the TZ Search has a heuristic behavior and the access index varies in accordance with the video properties (low/high motion/disparity). These two different cases are represented in the plots of the Fig. B.5.

Appendix C

CES Video Analyzer Tool

The CES Video Analyzer tool was developed in house targeting the displaying and analyzing of video properties. It was described in C# programming language and features the graphic user interface presented in Fig. C.1. The goal of the original tool is to support the decision making during novel coding algorithms design. The tool support diverse displaying modes including luminance-only mode and applying MB grids. Also, the CES Video Analyzer implements image filters such as Sobel, Laplace, Kirsch, and Prewitt filters besides luminance, gradient, and variance maps. An additional information window summarizes all image properties. Figure C.2 exemplifies the tool features presenting the original frame with the MB grid, the Sobel filtered image, and the variance map.

To facilitate the development of the algorithms proposed in this volume, the CES Video Analyzer was extended to support and provide better visualization for MVC videos. Figure C.3 shows the visualization of a frame differentiating SKIP, inter, and intra MBs. In Fig. C.4 all MBs, including SKIPs, are classified in disparity estimation or motion estimation for different time instants.

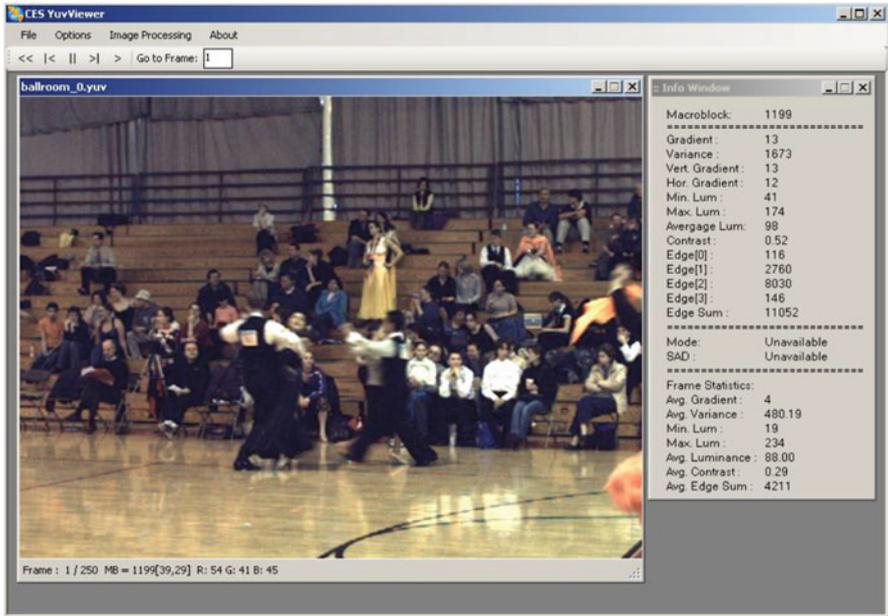


Fig. C.1 CES video analyzer user interface



Fig. C.2 CES video analyzer features

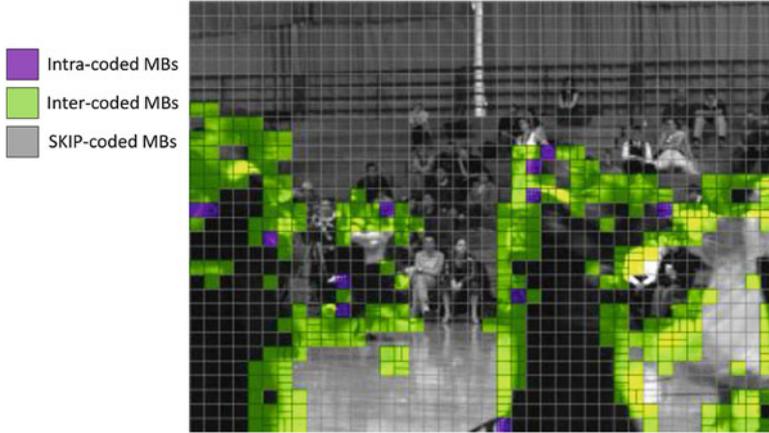


Fig. C.3 Coding mode analysis using CES video analyzer

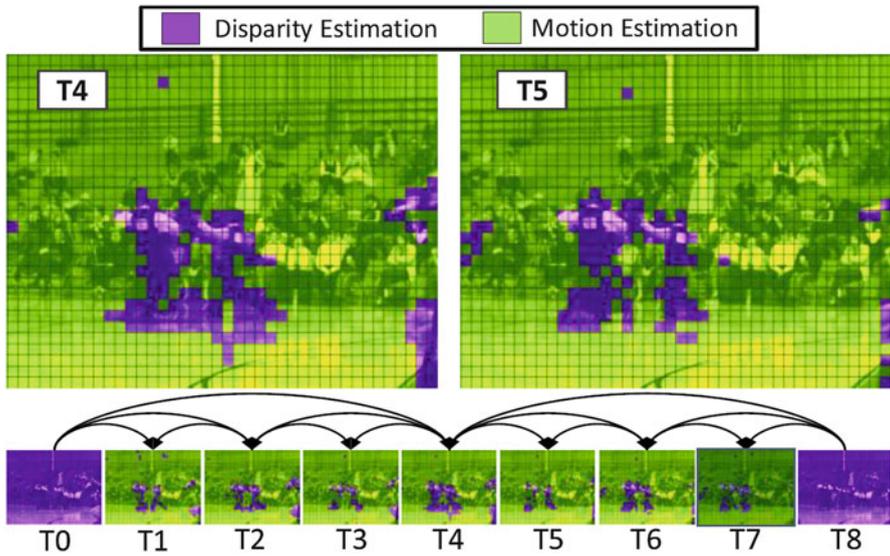


Fig. C.4 ME/DE analysis using CES video analyzer

References

- Abbo AA et al (2008) Xetal-II: a 107 GOPS, 600 mW massively parallel processor for video scene analysis. *IEEE J Solid-State Circuits* 43:192–201
- Agarwal K et al (2006) Power gating with multiple sleep modes. In: International symposium on quality electronic design, 27-29 March 2006, San Jose, CA, pp 633–637
- Agrafiotis D et al (2006) Multiple priority region of interest coding with H.264. In: IEEE international conference on image processing, ICIP, October 2006, Atlanta pp 53–56
- Akin A, Sayilar G, Hamzaoglu I (2010) A reconfigurable hardware for one bit transform based multiple reference frame motion estimation. In: Design, automation and test in europe, EDAA, Milão, Milan, Italy, pp 393–398
- Arapostathis A, Kumar R, Tangirala S (2003) Controlled Markov chains with safety upper bound. *IEEE Trans Automat Contr* 48:1230–1234
- ARM Ltd. (2012) ARM—the architecture for the digital world. <http://www.arm.com/>
- Arsura E et al (2005) Fast macroblock intra and inter modes selection for H.264/AVC. In: International conference on multimedia and expo (ICME), July 6-8, 2005, Amsterdam, The Netherlands, pp 378–381
- Barto Andrew G (1994) Reinforcement learning control. *Curr Opin Neurobiol* 4:888–893
- Bauer L et al (2007) RISPP: rotating instruction set processing platform. In: 44th Design automation conference, San Diego, CA, pp 791–796 [s.n.]
- Bauer L et al (2008a) Run-time system for an extensible embedded processor with dynamic instruction set. In: Design, automation and test in Europe, Munich, Germany, pp 752–757 [s.n.]
- Bauer L, Shafique M, Henkel J (2008b) Run-time instruction set selection in a transmutable embedded processor. In: 45th Design automation conference, pp 56–61
- Beck Antonio Carlos S et al (2008) Transparent reconfigurable acceleration for heterogeneous embedded applications. In: Design automation conference, 8-13 June, Anaheim, CA, pp 1208–1213
- Bennett Kyle (2011) Intel Core i7-3960X—Sandy Bridge E Processor Review, HardOCP, November 2011. http://hardocp.com/article/2011/11/14/intel_core_i73960x_sandy_bridge_e_processor_review/4
- Berekovic M et al (2008) Mapping of nomadic multimedia applications on the ADRES reconfigurable array processor. *Microprocess Microsyst* 33:290–294
- Bhaskaran V, Konstantinides K (1999) Image and video compression standards: algorithms and architectures. Kluwer Academic, Boston, MA
- Blanche P-A et al (2010) Holographic three-dimensional telepresence using large-area photorefractive polymer. *Nature* 468:80–83

- Blu-ray Disc Association (2010) White paper blu-ray disc read-only format. http://www.blu-raydisc.com/assets/Downloadablefile/BD-ROM_Audio_Visual_Application_Format_Specifications-18780.pdf
- Cadence Design Systems, Inc. (2012) Digital implementation. <http://www.cadence.com/products/di/Pages/default.aspx>
- Cao Z et al (2010) Optimality and improvement of dynamic voltage scaling algorithms for multimedia applications. *IEEE Trans Circuits Syst I Regul Pap* 57:681–690
- Chan C-C, Tang C-W (2012) Coding statistics based fast mode decision for multi-view video coding. *J Vis Commun Image Represent*. doi: [10.1016/j.jvcir.2012.01.004](https://doi.org/10.1016/j.jvcir.2012.01.004)
- Chang H-C et al (2009) A dynamic quality-adjustable H.264 video encoder for power-aware video applications. *IEEE Trans Circuits Syst Video Technol* 12:1739–1754
- Chang NY-C et al (2010) Algorithm and architecture of disparity estimation with mini-census adaptive support weight. *IEEE Trans Circuits Syst Video Technol* 20:792–805
- Chen JC, Chien S-Y (2008) CRISP: coarse-grained reconfigurable image stream processor for digital still cameras and camcorders. *IEEE Trans Circuits Syst Video Technol* 18:1223–1236, 1051–8215
- Chen Z, Zhou P, He Y (2002) Fast integer pel and fractional pel motion estimation for JVT, JVT-F017, Joint Video Team (JVT) of ISO/IECMPEG & ITU-T VCEG 6th Meeting, Awaji, JP
- Chen C-Y et al (2006) Level C+ data reuse scheme for motion estimation with corresponding coding orders. *IEEE Trans Circuits Syst Video Technol* 16:553–558
- Chen T-C et al (2007) Fast algorithm and architecture design of low-power integer motion estimation for H.264/AVC. *IEEE Trans Circuits Syst Video Technol* 17:568–577
- Chen Y et al (2009a) Coding techniques in multiview video coding and joint multiview video model. In: *Picture coding symposium, IEEE, Piscataway, NJ*, pp 313–316
- Chen Y et al (2009b) The emerging MVC standard for 3D video services. In: *3DTV conference, May 04-06 2009, Potsdam, Germany, vol 2009*, pp 1–13
- Chen Y-H et al (2009c) Algorithm and architecture design of power-oriented H.264/AVC baseline profile encoder for portable devices. *IEEE Trans Circuits Syst Video Technol* 19:1118–1128, 1051–8215
- Chien S-Y et al (2008) An 8.6 mW 25 Mvertices/s 400-MFLOPS 800-MOPS 8.91 mm multimedia stream processor core for mobile applications. *IEEE J Solid-State Circuits* 43:2025–2035
- Chiu J-C, Chou Y-L (2010) Multi-streaming SIMD multimedia computing engine. *Microprocess Microsyst* 34:247–258
- Chuang T-D et al (2010) A 59.5mW scalable/multi-view video decoder chip for Quad/3D Full HDTV and video streaming applications. In: *IEEE international conference on solid-state circuits (ISSCC), 7-11 Feb, San Francisco, CA*, pp 330–331
- Circuits Multi-Projects (2012) STMicroelectronics deep sub-micron processes. http://cmp.imag.fr/aboutus/slides/slides2007/04_KT_ST.pdf
- CISCO (2012) Cisco visual networking index: global mobile data traffic forecast update, 2011–2016. http://www.cisco.com/en/US/solutions/collateral/ns341/ns525/ns537/ns705/ns827/white_paper_c11-520862.pdf
- Cong J et al (2009) Automatic memory partitioning and scheduling for throughput and power optimization. In: *International conference on computer aided design (ICCAD), November 2-5, 2009, San Jose, CA*, pp 697–704
- de-Frutos-López M et al (2010) An improved fast mode decision algorithm for intraprediction in H.264/AVC video coding. *Signal Process Image Commun* 25:709–716
- Deng Z-P et al (2009) A fast view-temporal prediction algorithm for stereoscopic video coding. In: *International congress on image and signal processing (CISP), 17-19 October 2009, Tianjin, China*. pp 1–5
- Díaz-Honrubia Antonio J, Martínez José Luis, Cuenca Pedro (2012) HEVC: a review, trends and challenges. In: *Workshop on multimedia data coding and transmission, Sep. 2012, WMDCT, Alicante, Spain*
- Ding L-F et al (2008a) Content-aware prediction algorithm with inter-view mode decision for multiview video coding. *IEEE Trans Multimedia* 10:1553–1564

- Ding L-F et al (2008b) Fast motion estimation with inter-view motion vector prediction for stereo and multiview video coding. In: International conference on acoustics speech and signal processing (ICASSP), Las Vegas, NV, March 30 – April 4, 2008, pp 1373–1376
- Ding L-F et al (2010a) A 212 MPixels/s 4096 2160p multiview video encoder chip for 3D/Quad full HDTV applications. *IEEE J Solid-State Circuits* 45:46–58
- Dodgson NA (2005) Autostereoscopic 3D displays. *IEEE Comput* 38:31–36
- Dolby (2012) Dolby 3D. <http://www.dolby.com/us/en/consumer/technology/movie/dolby-3d.html>
- Erdayandi K (2009) JMVC documentation. In: JMVC—JVT-AD207. http://students.sabanciuniv.edu/~kerdayandi/jmvc/index_jmvc.html
- Finchelstein DF, Sze V, Chandrakasan AP (2009) Multicore processing and efficient on-chip caching for H.264 and future video decoders. *IEEE Trans Circuits Syst Video Technol* 19:1704–1713
- Fujifilm (2011) FinePix REAL 3D W3 | FujiFilm Global. http://www.fujifilm.com/products/3d/camera/finepix_real3dw3/
- Fujii T (2010) Panel Discussion 1 (D1)—3DTV/FTV. In: Picture coding symposium (PCS), December 7–10, 2010, Nagoya, Japan
- Fukano G et al (2008) A 65nm 1Mb SRAM macro with dynamic voltage scaling in dual power supply scheme for low power SoCs, In: Joint Non-Volatile Semiconductor Memory Workshop and International Conference on Memory Technology and Design (NVSMW/ICMTD), Opio, France, pp 97–98
- García Carlos E, Prett David M, Morari M (1989) Model predictive control: theory and practice—a survey. *Automatica* 25:335–348
- Gassée J-L (2010) Intel's bold bet against ARM: visionary or myopic? Monday Note. <http://www.mondaynote.com/2010/06/27/intel%E2%80%99s-bold-bet-against-arm-visionary-or-myopic/>
- Ghanbari M (1990) The cross-search algorithm for motion estimation. *IEEE Trans Commun* 38:950–953
- Grecos C, Yang MY (2005) Fast inter mode prediction for P slices in the H264 video coding standard. *IEEE Transactions On Broadcasting*, Vol. 51, No. 2, June 2005. 256–263
- Grun P, Balasa F, Dutt N (1998) Memory size estimation for multimedia applications, In: International Workshop on Hardware/Software Codesign (CODES/CASHE), California University Press, Irvine, CA, pp 145–149
- Han D-H, Lee Y-L (2008) Fast mode decision using global disparity vector for multiview video coding. In: Future generation communication and networking symposia (FGCNS), December 13–December 15 2008. Proceedings of the 2008 Second International Conference on Future Generation Communication and Networking Symposia - Volume 01. IEEE Computer Society, Washington, DC, USA. pp 209–213
- He Z, Cheng W, Chen X (2008) Energy minimization of portable video communication devices based on power-rate-distortion optimization. *IEEE Trans Circuits Syst Video Technol* 18:596–608
- Yu-Wen Huang; Bing-Yu Hsieh; Shao-Yi Chien; Shyh-Yih Ma; Liang-Gee Chen, Analysis and complexity reduction of multiple reference frames motion estimation in H.264/AVC, *IEEE Transactions on Circuits and Systems for Video Technology*, vol.16, no.4, pp.507,522, April 2006. doi: 10.1109/TCSVT.2006.872783
- Huang Y-H, Ou T-S, Shen H (2009) Fast H.264 selective intra mode decision for inter-frame coding. In: Picture coding symposium (PCS), Chicago, Illinois, USA, on May 6–8, 2009, pp 377–380
- Huff HR, Gilmer DC (2004) High dielectric constant materials: VLSI MOSFET applications. Springer, New York, NY
- IC Insights (2012) IC Insights raises forecast for tablets, notebooks, total PC shipments. Electronic specifier. <http://www.electronicsspecifier.com/Tech-News/IC-Insights-Raises-Forecast-Tablets-Notebooks-Total-PC-Shipments.asp>
- IMAX (2012) IMAX3D. <http://www.imax.com/about/imax-3d/>
- ISO/IEC (2009) Vision on 3D video. <http://mpeg.chiariglione.org/visions/3dv/index.htm>

- ISO/IEC (2011) Common test conditions for MVC W12036. ISO/IEC JTC1/SC29/WG11, Meeting MPEG 96, Geneva, Switzerland, March de 2011
- ITU-T (1999) Subjective video quality assessment methods for multimedia applications—P.910 // Series P: telephone transmission quality, telephone installations, local line networks
- Javed H et al (2011) Low-power adaptive pipelined MPSoCs for multimedia: an H.264 video encoder case study. In: Design automation conference, JUNE 2-6, San Diego, CA, pp 1032–1037
- Jeon BW, Lee JY (2003) Fast mode decision for H.264—Document JVT-J033. Joint Video Team (JVT) of ISO/IECMPEG & ITU-T VCEG 8th Meeting, Waikoloa, HI
- Ji W, Li P, Chen M, and Chen Y (2009) Power Scalable Video Encoding Strategy Based on Game Theory. In Proceedings of the 10th Pacific Rim Conference on Multimedia: Advances in Multimedia Information Processing (PCM '09), Paisarn Muneesawang, Feng Wu, Itsuo Kumazawa, Athikom Roeksabutr, Mark Liao, and Xiaou Tang (Eds.). Springer-Verlag, Berlin, Heidelberg, 1237–1243.
- Jiang M, Yi X, Ling N (2004) Improved frame-layer rate control for H.264 using MAD ratio. In: International symposium on circuits and systems, ISCAS, 23-26 May, Vancouver, Canada
- Jing X, Chau L-P (2001) An efficient three-step search algorithm for block motion estimation,” *Multimedia, IEEE Transactions on*, vol.6, no.3, pp. 435,438, June 2004. doi: 10.1109/TMM.2004.827517
- Jing X, Chau L-P (2004) Fast approach for H.264 inter mode decision, *Electronics Letters* , vol.40, no.17, pp.1050,1052, 19 Aug. 2004. doi: 10.1049/el:20045243
- JVT (2003) Draft ITU-T Rec. and final draft international standard of joint video specification
- JVT (2008) Joint draft 8.0 on multiview video coding—JVT-AB204
- JVT (2009a) JMVC 6.0 [garcon.ient.rwthachen.de]
- JVT (2009b) Joint multiview video coding
- Kamat SP (2009) Energy management architecture for multimedia applications in battery powered devices. *IEEE Trans Consum Electron* 55:763–767
- Kauff P et al (2007) Depth map creation and image-based rendering for advanced 3DTV services providing interoperability and scalability. *Signal Process Image Commun* 22:217–234
- Kay Roger (2011) Forbes. Is the PC dead? <http://www.forbes.com/sites/rogerkay/2011/02/28/is-the-pc-dead/>
- Khailany BK et al (2008) A programmable 512 GOPS stream processor for signal, image, and video processing. *IEEE J Solid-State Circuits* 43:202–213
- Kim B-G, Cho C-S (2007) A fast inter-mode decision algorithm based on macro-block tracking for P slices in the H.264/AVC video standard. In: International conference on image processing (ICIP), September 16-19, 2007, San Antonio, Texas, USA, pp V-301–V-304
- Kim C, Kuo C-CJ (2007) Feature-based intra-/intercoding mode selection for H.264/AVC. *IEEE Trans Circuits Syst Video Technol* 17:441–453
- Kim D-Y, Lee Y-L (2011) A fast intra prediction mode decision using DCT and quantization for H.264/AVC. *Signal Process Image Commun* 26:455–465
- Kim Y, Kim J, Sohn K (2007a) Fast disparity and motion estimation for multi-view video coding. *IEEE Trans Circuits Syst Video Technol* 53:712–719
- Kim Y, Kim J, Sohn K (2007b) Fast disparity and motion estimation for multi-view video coding. [s.l.]. *IEEE Trans Circuits Syst Video Technol* 53: 712–719
- Ko H, Yoo K, Sohn K (2009) Fast mode-decision for H.264/AVC based on inter-frame correlations. *Signal Process Image Commun* 24:803–813
- Kollig P, Osborne C, Henriksson T (2009) Heterogeneous multi-core platform for consumer multimedia applications. In: Design, automation test in Europe conference, April 20-24, Nice, France, pp 1254–1259
- Kondo H et al (2009) Heterogeneous multicore SoC with SiP for secure multimedia applications. *IEEE J Solid-State Circuits* 44:2251–2259, 0018-9200
- Koo H-S, Jeon Y-J, Jeon B-M (2007) MVC Motion skip mode - Doc. JVT-W081
- Krolikoski Stan (2004) Chipvision design systems. Orinoco saves power. <http://www.eda.org/edps/edp04/submissions/presentationKrolikoski.pdf>

- Krügner J et al (2005) Image based 3DSurveillance for flexible Man-Robot-cooperation. *CIRP Ann Manuf Technol* 54:19–22
- Kuhn P (1999) Algorithms, complexity analysis and VLSI architectures for MPEG-4 motion estimation. Kluwer Academic, Boston, MA
- Kume Hideyoshi (2010) Panasonic's new Li-Ion batteries use Si Anode for 30% higher capacity. *TechOn*. <http://techon.nikkeibp.co.jp/article/HONSHI/20100223/180545/>
- Kwon D-K, Shen M-Y, Kuo C-CJ (2007) Rate control for H.264 video with enhanced rate and distortion models. *IEEE Trans Circuits Syst Video Technol* 17:517–529, 1051–8215
- Lee P-J, Lai Y-C (2011) Vision perceptual based rate control algorithm for multi-view video coding. In: International conference on system science and engineering (ICSSE), 8-10 June, 2011, Macao, China pp 342–345
- Lee S-Y, Shin K-M, Chung K-D (2008) An object-based mode decision algorithm for multi-view video coding. In: International symposium on multimedia (ISM), December 15-17, 2008, Berkeley, California, USA pp 74–81
- Li Z. G. et al (2003) Adaptive basic unit layer rate control for JVT - JVT-G012. Joint Video Team (JVT) of ISO/IEC MPEG & ITU-T VCEG 7th Meeting, Pattaya, Thailand
- Liang Y, Ahmad I (2009) Power and distortion optimization for pervasive video coding. *IEEE Trans Circuits Syst Video Technol* 19:1436–1447
- Lim KP (2003) Fast inter mode selection—Document JVT-I020. In: 9th JVT Meeting
- Lin J-P, Tang Angela C-W (2009) A fast direction predictor of inter frame prediction for multi-view video coding. In: IEEE international symposium on circuits and system, IEEE, Piscataway, pp 2598–2593
- Lin Y-K et al (2008) A 242mW 10mm² 1080p H.264/AVC high profile encoder chip. In: Design automation conference, DAC Anaheim, CA, USA, June 8-13, pp 78–83
- Ling N (2010) Expectations and challenges for next generation. In: Conference on industrial electronics and applications, ICIEA, 15-17 June 2010, Taichung, Taiwan pp 2339–2344
- Liu X, Shenoy P, Corner MD (2008) Chameleon: application level power management. *IEEE Trans Mobile Comput* 7:995–1010, 1536–1233
- Liu A et al (2010) Just noticeable difference for images with decomposition model for separating edge and textured regions. *IEEE Trans Circuits Syst Video Technol* 20:1648–1652
- Lu X et al (2005) Fast mode decision and motion estimation for H.264 with a focus on MPEG-2/H.264 transcoding. In: International conference on circuits and systems (ISCAS), 23-26 May 2005, Kobe, Japan, pp 1246–1249
- Ma S, Gao W, Lu Y (2005) Rate-distortion analysis for H.264/AVC video coding and its application to rate control. *IEEE Trans Circuits Syst Video Technol* 15:1533–1544
- Marlow S, Ng J, McArdle C (1997) Efficient motion estimation using multiple log searching and adaptive search windows, IPA, July 1997. Dublin, Ireland In: International conference on image processing and its applications, pp 214–218
- McCann K et al (2012) Technical Evolution of the DTT Platform—an independent report by ZetaCast, commissioned by Ofcom. <http://stakeholders.ofcom.org.uk/binaries/consultations/uhf-strategy/zetacast.pdf>. Accessed Jan 2012
- Meng B et al (2003) Efficient intra-prediction mode selection for 4x4 blocks in H.264. In: International conference on multimedia and expo (ICME), 6-9 July 2003, Baltimore, MD, USA. pp III-521–III-524
- Mentor Graphics (2012) ModelSim—advanced simulation and debugging. <http://model.com/>
- Merkle P et al (2007) Efficient prediction structures for multiview video coding. *IEEE Trans Circuits Syst Video Technol* 17:1461–1473
- Merkle P et al (2009) Stereo video compression for mobile 3D services. In: 3DTV, May 04-06 2009, Potsdam, Germany Conference, pp 1–4
- Merritt L, Vanam R (2007) Improved rate control and motion estimation for H.264 Encoder. In: IEEE international conference on image processing ICIP, September 16-19, 2007, San Antonio, Texas, USA, pp V-309–V-312
- Miano J (1999) Compressed image file formats: Jpeg, Png, Gif, Xbm, Bmp. ACM Press, Boston, MA

- Mondal S, Ogrenci Memik S (2005) Fine-grain leakage optimization in SRAM based FPGAs. In: ACM great lakes symposium on VLSI, Chicago, Illinois, USA, April 17-19, 2005, pp 238–243
- Morari M, Lee JH (1999) Model predictive control: past, present and future. *Comput Chem Eng* 23:667–682
- Muller K et al (2005) 3-D reconstruction of a dynamic environment with a fully calibrated background for traffic scenes. [s.l.]. *IEEE Trans Circuits Syst Video Technol* 15: 538–549
- Naccari M et al (2011) Low Complexity Deblocking Filter Perceptual Optimization For The HEVC codec. In: International conference on image processing ICIP, Brussels, Belgium, September 11-14, 2011, pp 737–740
- Nintendo (2011) Nintendo 3DS. <http://www.nintendo.com/3ds>
- Nvidia (2012a) Nvidia GeForce GX690. <http://www.geforce.com/hardware/desktop-gpus/geforce-gtx-690>
- Nvidia (2012b) Tegra 3 super processors. <http://www.nvidia.com/object/tegra-3-processor.html>
- Nvidia Corp. (2012) Tegra 2 and Tegra 3 super processors. <http://www.nvidia.com/object/tegra-3-processor.html>
- Oh K-J, Lee J, Park D-S (2011) Multi-view video coding based on high efficiency video coding. In: Pacific Rim conference on advances in image and video technology, PSIVT 2011, November 2011, Gwangju, South Korea, pp 371–380
- Ostermann J et al (2004) Video coding with H.264/AVC: tools, performance, and complexity. *IEEE Circuits Syst Mag* 4(1st Quarter):7–28
- Otero A et al (2010) Run-time scalable systolic coprocessors for flexible. In: International conference on field programmable logic and applications (FPL), August 31 2010-September 2, 2010, Milano, Italy pp 70–76, 1946-1488
- Ou T-S, Huang Y-H, Chen HH (2009) Efficient MB and prediction mode decisions for intra prediction of H.264 high profile. In: Picture coding symposium (PCS), 6-8 May, Chicago, IL, USA, pp 1–4
- Ozbek N, Tekalp AM, Tunali ET (2007) Rate allocation between views in scalable stereo video coding using an objective stereo video quality measure. In: International conference on acoustics speech and signal processing (ICASSP), April 15-20, 2007, Honolulu, Hawaii, USA, pp 1045–1048
- Pan F et al (2005) Fast mode decision algorithm for intraprediction in H.264/AVC video coding. *IEEE Trans Circuits Syst Video Technol* 15:813–822
- Panasonic Panasonic HDC-SDT750K (2011). <http://www2.panasonic.com/consumer-electronics/support/Cameras-Camcorders/Camcorders/3D-CAMCORDERS/model.HDC-SDT750K>
- Panda PR, Dutt ND, Nicolau A (1997) Architectural exploration and optimization of local memory in embedded systems. In: International symposium on system synthesis ISSS, September 17-19, 1997, Antwerp, Belgium, vol 10, pp 90–97
- Park I, Capson DW (2008) Improved inter mode decision based on residue in H.264/AVC. In: International conference on multimedia and expo (ICME), June 23-26 2008, Hannover, Germany, pp 709–712
- Park S, Sim D (2009) An efficient rate-control algorithm for multi-view video coding. In: IEEE international symposium on consumer electronics (ISCE), May 25-28, 2009, Mielparque-Kyoto, Kyoto, Japan, pp 115–118
- Park JS, Song HJ (2006) Selective intra prediction mode decision for H.264/AVC encoders. *Int J Appl Sci Eng Technol* 13:214–218
- Payá-Vayá G et al (2010) VLIW architecture optimization for an efficient computation of stereoscopic video applications. In: International conference on green circuits and systems (ICGCS), 21-23 June, Shanghai, china, pp 457–462
- Pei G et al (2002) FinFET design considerations based on 3-D simulation and analytical modeling. *IEEE Trans Electron Devices* 49:1411–1419, 0018-9383
- Peng Z et al (2008a) Fast macroblock mode selection algorithm for multiview video coding. In: *EURASIP J Image Video Process*, Volume 2008:393727 doi:10.1155/2008/393727

- Pourazad M, Nasiopoulos P, Ward R (2009) An efficient low random-access delay panorama-based multiview video coding scheme. In: IEEE conference on image processing ICIP, 7-10 November 2009, Cairo, Egypt, IEEE, Cairo, pp 2945–2948
- Qualcomm Inc. (2011) Snapdragon S4 processors: system on chip solutions for a new mobile age—white paper. <https://developer.qualcomm.com/download/snapdragon-s4-processors-system-on-chip-solutions-for-a-new-mobile-age.pdf>
- Rajamani K et al (2006) Application-aware power management. In: IEEE international symposium on workload characterization, IISWC-2006, October 25-27, San Jose, California, USA, pp 39–48
- RealD (2012) RealD 3D. <http://reald.com/>
- Research and Markets (2010) 3D TV market and future forecast worldwide (2010–2014). http://www.researchandmarkets.com/reports/1525112/3d_tv_market_and_future_forecast_worldwide_2010
- Richardson I (2010) The H. 264 advanced video compression standard. Wiley, New York [s.l.]
- Roy S, Ranganathan N, Katkooi S (2011) State-retentive power gating of register files in multi-core processors featuring multithreaded in-order cores. IEEE Trans Comput 60:1547–1560, 0018-9340
- Salgado L, Nieto M (2006) Sequence independent very fast mode decision algorithm on H.264/AVC baseline profile. In: International conference on image processing ICIP, October 8-11, Atlanta, Georgia, USA, pp 41–44
- Samsung (2012) Samsung Galaxy SIII. <http://www.samsung.com/global/galaxy3/>
- Samsung Electronics Co. Ltd. (2012) Samsung Exynos 4 Quad. <http://www.samsung.com/global/business/semiconductor/minisite/Exynos/products4quad.html>
- Saponara S, Fanucci L (2004) Data-adaptive motion estimation algorithm and VLSI architecture design for low-power video systems. IEE Comput Digit Tech 151:51–59
- Shafique M, Bauer L, Henkel J (2008) 3-tier dynamically adaptive power-aware motion estimator for h.264/AVC video encoding. In: International symposium on low power electronics and design (ISLPED), August 11-13, Bangalore, India, pp 147–152
- Shafique M, Bauer L, Henkel J (2010) enBudget: a run-time adaptive predictive energy-budgeting scheme for energy-aware motion estimation in H.264/MPEG-4 AVC video encoder. In: Design, automation and test in Europe (DATE), Dresden, Germany, March 8-12
- Shafique M, Molkenthin B, Henkel J (2010a) An HVS-based adaptive computational complexity reduction scheme for H.264/AVC video encoder using prognostic early mode exclusion. In: IEEE design, automation and test in Europe (DATE), Dresden, Germany, March 8-12, pp 1713–1718
- Shafique M et al (2010b) Power-aware complexity-scalable multiview video coding for mobile devices. In: 28th picture coding symposium (PCS'10), 8-10 December, Nagoya, Japan, pp 350–353
- Sharp (2011) Lynx 3D SH-03C. <http://www.sharp.co.jp/products/sh03c/index.html>
- Shen L et al (2009a) Selective disparity estimation and variable size motion estimation based on motion homogeneity for multi-view coding. IEEE Trans Broadcast 55:761–766
- Shen L et al (2009b) Fast mode decision for multiview video coding. In: International conference on image processing (ICIP), 7-10 November 2009, Cairo, Egypt, pp 2953–2956
- Shen L et al (2010a) Early SKIP mode decision for MVC using inter-view correlation. Signal Process Image Commun 25:88–93
- Shen L et al (2010b) View-adaptive motion estimation and disparity estimation for low complexity multiview video coding. IEEE Trans Circuits Syst Video Technol 20:925–930
- Shim H, Kyung C-M (2009) Selective search area reuse algorithm for low external memory access motion estimation. IEEE Trans Circuits Syst Video Technol 19:1044–1050
- Shimpi AL (2011) ARM's Mali-T658 GPU in 2013, Up to 10x faster than Mali-400. AnadTech. <http://www.anandtech.com/show/5077/arms-malit658-gpu-in-2013-up-to-10x-faster-than-mali400>
- Singh H et al (2007) Enhanced leakage reduction techniques using intermediate strength power gating. IEEE Trans VLSI Syst 15:1215–1224

- Smolic A et al (2007) Coding algorithms for 3DTV—a survey. *IEEE Trans Circuits Syst Video Technol* 17:1606–1621
- Social Times (2011) Social times. Cisco predicts that 90% of all internet traffic will be video in the next three years. http://socialtimes.com/cisco-predicts-that-90-of-all-internet-traffic-will-be-video-in-the-next-three-years_b82819
- Softpedia (2010) ARM wants a share out of the server and desktop pc market by 2015. Softpedia. <http://news.softpedia.com/newsImage/ARM-Wants-a-Share-of-the-Server-and-Desktop-PC-Market-by-2015-5.png/>
- Sony (2011) HDR-TD10—full HD 3D Camcorder. <http://www.sonymstyle.com/webapp/wcs/stores/servlet/ProductDisplay?catalogId=10551&storeId=10151&langId=-1&productId=8198552921666294297>
- Stelmach LB, Tam JW (1998) Stereoscopic image coding: effect of disparate image-quality in left- and right-eye views. *Signal Process Image Commun* 14:111–117, 0923–5965
- Stelmach LB, Tam WJ (1999) Stereoscopic image coding: effect of disparate image-quality in left- and right-eye views. [s.l.]. *Signal Process Image Commun* 14: 111–117
- Stoykova E et al (2007) 3-D time-varying scene capture technologies: a survey. *IEEE Trans Circuits Syst Video Technol* 17:1568–1586
- Su Y, Vetro A, Smolic A (2006) Common test conditions for multiview video coding—Doc. JVT-T207
- Sullivan GJ, Ohm J-R (2010) Recent developments in standardization of high efficiency video coding (HEVC). *Proc. SPIE 7798, Applications of Digital Image Processing XXXIII, 77980V* (September 07, 2010); doi:10.1117/12.863486
- Sullivan G, Wiegand T (1998) Rate-distortion optimization for video compression. *IEEE Signal Process Mag* 15:74–90
- Sullivan GJ, Wiegand T (2005) Video compression—from concepts to the H.264/AVC standard. *Proc IEEE* 93:18–31
- Synopsys, Inc. (2012) IBM—65NM. <http://www.synopsys.com/dw/emlselector.php?f=IBM&g=65>
- Tan TK, Sullivan G, Wedi T (2005) Recommended simulation conditions for coding efficiency experiments—VCEG-AA10. Nice, [s.n.]
- Tang Xiu-li, Dai Sheng-kui and Cai Can-hui (2010) An analysis of TZSearch algorithm in JMVC. In: International conference on green circuits and systems (ICGCS), 21-23 June 2010, Shanghai, china, pp 516–520
- Tanimoto M (2005) FTV (free viewpoint television) creating ray-based image engineering. In: International conference on image processing (ICIP), Genoa, Italy, September 11-14, pp 25–28
- Tatjewski P (2010) Supervisory predictive control and on-line set-point optimization. *Int J Appl Math Comput Sci* 20:483–495, 1641–876X
- Tech G et al (2010) Final report on coding algorithms for mobile 3DTV // MOBILE3DTV—Project No. 216503. http://sp.cs.tut.fi/mobile3dvtv/results/tech/D2.6_Mobile3DTV_v1.0.pdf
- Texas Instruments Inc. (2012) OMAP™ mobile processors: OMAP™ 5 platform. <http://www.ti.com/general/docs/wtbu/wtbuproductcontent.tsp?templateId=6123&navigationId=12862&contentId=101230>
- The Digital Entertainment Group (2009) 3D White Paper
- Tian L et al (2010) Analysis of quadratic R-D model in H.264/AVC video coding. In: 17th IEEE international conference on image processing (ICIP), September 26-29, Hong Kong, China, pp 2853–2856
- Tourapis AM (2002) Enhanced predictive zonal search for single and multiple frame motion estimation. In: Visual communication and image processing conference (VCIP), Huang Shan, An Hui, China, 11-14 July, 2010
- Tsai C-Y et al (2007) Low power cache algorithm and architecture design for fast motion estimation in H.264/AVC encoder system. In: International conference on acoustics speech and signal processing (ICASSP), April 15-20, 2007, Honolulu, Hawaii, USA, pp II-97–II-100
- Tsung P-K et al (2009) Cache-based integer motion/disparity estimation for quad-HD H.264/AVC and HD multiview video coding. In: International conference on acoustics, speech and signal processing, IEEE (ICASSP), 19-24 April 2009, Taipei, Taiwan Taipei, pp 2013–2016

- Tuan T, Kao S, Trimberger S (2006) A 90nm low-power FPGA for battery-powered applications. In: International symposium on field programmable gate arrays (FPGA), pp 3–11
- Vimeo (2012) Vimeo 3D. <http://vimeo.com/channels/stereoscopy>
- Vizzotto BB et al (2012) A model predictive controller for frame-level rate control in multiview video coding. IEEE international conference on multimedia & expo (ICME'12), Melbourne, Australia, July 9-13, 2012, pp 485–490
- Wang W, Mishra P (2010) Leakage-aware energy minimization using dynamic voltage scaling and cache reconfiguration in real-time systems. In: 23rd international conference on VLSI design, 3-7 Jan. 2010, Bangalore, India, vol 23, pp 357–372
- Wang X et al (2007) Fast mode decision for H.264 video encoder based on MB motion characteristic. In: International conference on multimedia and expo (ICME), July 2-5, 2007, Beijing, China, pp 372–375
- Wang S-H, Tai S-H, Chiang T (2009) A low-power and bandwidth-efficient motion estimation IP core design using binary search. IEEE Trans Circuits Syst Video Technol 19:760–765
- Wei Z, Ngan KN, Li H (2008) An efficient intra-mode selection algorithm for H.264 based on edge classification and rate-distortion estimation. Signal Process Image Commun 23:699–710
- Welch G et al (2005) Remote 3D medical consultation, vol 2, pp 1026–1033
- Wen J et al (2010) ESVD: an integrated energy scalable framework for low-power video decoding systems. EURASIP J Wirel Commun Netw 5:5:1–5:13
- Wiegand T et al (2003) Overview of the H.264/AVC video coding standard. IEEE Trans Circuits Syst Video Technol 13:560–576
- Willner K et al (2008) Mobile 3D video using MVC and N800 internet tablet. In: 3DTV conference 28-30 MAY 2008, ISTANBUL, TURKEY
- Woo J-H et al (2008) A 195 mW/152 mW mobile multimedia SoC with fully programmable 3-D graphics and MPEG4/H.264/JPEG. IEEE J Solid-State Circuits 43:2047–2056
- Wu C-Y, Su P-C (2009) A region of interest rate-control scheme for encoding traffic surveillance videos. In: International conference on intelligent information hiding and multimedia signal processing (IIH-MSP), September 12 - 14, 2009 Kyoto, Japan, pp 194–197
- Wu D et al (2004) Block inter mode decision for fast encoding of H.264. In: International conference on acoustics speech and signal processing (ICASSP), May 17-21, 2004, Montreal, Canada, vol iii, pp 181–184
- Xilinx, Inc. (2012) ISE design suite. <http://www.xilinx.com/products/design-tools/ise-design-suite/index.htm>
- Xu X, He Y (2008) Fast disparity motion estimation in MVC based on range prediction. In: IEEE international conference on image processing, October 12-15, 2008, San Diego, California, USA, 2008, San Diego. IEEE, Piscataway, pp 2000–2003
- Xu L et al (2011) Priority pyramid based bit allocation for multiview video coding. In: IEEE visual communications and image processing (VCIP), Tainan, Taiwan, November 6-9, 2011, pp 1–4
- Yamaoka M, Shinozaki Y, Maeda N, Shimazaki Y et al (2004) A 300MHz 25 μ A/Mb leakage on-chip SRAM module featuring process-variation immunity and low-leakage-active mode for mobile-phone application processor. In: IEEE international solid-state circuits conference (ISSCC), February 9-11, 2004, San Francisco, CA, pp 494–542
- Yamaoka M et al (2005) A 300-MHz 25 μ A/Mb-leakage on-chip SRAM module featuring process-variation immunity and low-leakage-active mode for mobile-phone application processor. [s.l.]. IEEE 40: 186–194
- Yan T et al (2009a) Frame-layer rate control algorithm for multi-view video coding. In: ACM/SIGEVO, Shanghai, China — June 12 - 14, 2009, summit on genetic and evolutionary computation, pp 1025–1028
- Yan T et al (2009b) Rate control algorithm for multi-view video coding based on correlation analysis. In: Symposium on photonics and optoelectronics (SOPO), Aug 23, 2009 - Aug 25, 2009 Wuhan, China pp 1–4
- Yang J (2009) Multiview video coding based on rectified epipolar lines. In: International conference on information, communication and signal processing (ICICS), 7 – 10 December 2009, Macau, China, pp 1–5

- Yang S, Wolf W, Vijaykrishnan N (2005) Power and performance analysis of motion estimation based on hardware and software realizations. *IEEE Trans Comput* 54:714–726
- YouTube (2011) YouTube—Broadcast yourself. <http://www.youtube.com/>
- YouTube 3D (2011) YouTube—3D Channel. <http://www.youtube.com/user/3D>
- Yu AC (2004) Efficient block-size selection algorithm for inter-frame coding in H.264/MPEG-4 AVC. In: International conference on acoustic, speech and signal processing (ICASSP), May 17–21, 2004, Montreal, Canada, pp III169–III172
- Zatt B et al (2007) Memory hierarchy targeting bi-predictive motion compensation for H.264/AVC decoder. In: IEEE computer society annual symposium on VLSI (ISVLSI), May 9–11, 2007, Porto Alegre, Brazil, pp 445–446
- Zatt B et al (2010) An adaptive early skip mode decision scheme for multiview video coding. In: Picture coding symposium (PCS), Nagoya, Japan, 8–10 December, pp 42–45
- Zatt B et al (2011a) A low-power memory architecture with application-aware power management for motion & disparity estimation in multiview video coding. In: IEEE/ACM 29th international conference on computer-aided design (ICCAD'11), November 7–11, 2010, San Jose, CA, USA, vol 29, pp 40–47
- Zatt B et al (2011b) A multi-level dynamic complexity reduction scheme for multiview video coding. *IEEE 18th international conference on image processing (ICIP'11)*, Brussels, Belgium, September 11–14, 2011, vol 18, pp 761–764
- Zatt B et al (2011c) Multi-level pipelined parallel hardware architecture for high throughput motion and disparity estimation in multiview video coding. In: IEEE/ACM 14th design automation and test in europe conference (DATE'11), Grenoble, France, March 14–18, 2011, vol 14, pp 1448–1453
- Zatt B et al (2011d) Run-time adaptive energy-aware motion and disparity estimation in multiview video coding. In: ACM/IEEE/EDA 48th design automation conference (DAC'11), San Diego, California, USA, June 5–10, pp 1026–1031
- Zeng H, Ma K-K, Cai C (2011) Fast mode decision for multiview video coding using mode correlation. *IEEE Trans Circuits Syst Video Technol* 21:1659–1666
- Zhang K, Bhattacharya U, Zhanping Chen, Hamzaoglu F, Murray D, Vallepalli N, Yih Wang, Zheng B, Bohr M, “SRAM design on 65-nm CMOS technology with dynamic sleep transistor for leakage reduction,” *Solid-State Circuits, IEEE Journal of*, vol.40, no.4, pp.895,901, April 2005. doi: 10.1109/JSSC.2004.842846
- Zhang Y et al (2009) ASIP approach for multimedia applications based on a scalable VLIW DSP architecture. *Tsinghua Sci Technol* 14:126–132
- Zhou Y et al (2011) PID-based bit allocation strategy for H.264/AVC rate control. *IEEE Trans Circuits Syst II Express Briefs* 58:184–188
- Zhu H, Luican II, Balasa F (2006) Memory size computation for multimedia processing applications. In: Asia and South Pacific conference on design automation, ASP-DAC 2006, Yokohama, Japan, January 24–27, pp 802–807
- Zone R (2007) *Stereoscopic cinema and the origins of 3-D film, 1838–1952*. ISBN 0813124611

Index

A

- Address generation unit (AGU), 129–130, 134–135
- Application-aware power gating, 9, 70, 128, 146–147
- Application-specific integrated circuits (ASIC), 30, 32, 36, 155
- Availability, 84

B

- Ballroom sequence, 80, 142, 145–146, 159
- Basic unit-level bitrate distribution, 86–87
- Basic unit-level rate control
 - MDP-based, 113
 - MVC, 65
 - for video-quality management
 - coupled reinforcement learning, 121
 - diagram, 119
 - Markov decision process, 119–121
 - regions of interest, 120
- Benchmark video sequence, 152–155
- Bitrate correlation analysis
 - basic unit-level bitrate distribution, 86–87
 - frame-level bitrate distribution, 86
 - view-level bitrate distribution, 85

C

- CES video analyzer tool, 185–187
- Coding mode correlation analysis
 - analysis, 75–76
 - coding mode analysis summary, 80–81
 - coding mode distribution analysis, 74–75
 - RDCost analysis, 79–80
 - video property analysis, 76, 78–79
- Coupled reinforcement learning, 113, 121

D

- 3D Blu-Ray, 29
- 3D-cinema systems, 29
- 2D/3D digital video
 - 1D-array/2D-array, 12
 - FTV system, 14
 - HVS, 11
 - macroblocks, 12
 - multiview video sequence, 12–13
 - RGB space, 11
 - YUV space, 11
- Decoded picture buffer (DPB), 59, 129
- Disparity domain correlation, 15, 16, 19
- 3D multimedia processing
 - 3D-video pre-and post-processing, 172
 - dynamic search window, 171
 - energy reduction, 169
 - fast ME/DE, 170–171
 - frame and basic unit (BU) level, 170
 - HRC algorithm, 170
 - HVS, 170
 - issues and challenges, 6–7
 - MDP, 170
 - mode decision algorithm, 169–170
 - MPC, 170
 - multibank video on-chip memory, 171
 - MVC challenges, 172
 - next-generation 3D-video coding, 172–173
 - relax and aggressive strengths, 169–170
 - requirements and trends, 3–5
 - state of the art, 170–171
 - video sequence, 171
- 3D-neighborhood correlation analysis
 - bitrate correlation analysis
 - basic unit-level bitrate distribution, 86–87

- 3D-neighborhood correlation analysis (*cont.*)
 - frame-level bitrate distribution, 86
 - view-level bitrate distribution, 85
 - coding mode correlation analysis
 - analysis, 75–76
 - coding mode analysis summary, 80–81
 - coding mode distribution analysis, 74–75
 - RDCost analysis, 79–80
 - video property analysis, 76–79
 - concept, 7–8
 - energy-efficient algorithms and architectures, 67–68
 - motion correlation analysis
 - basic prediction structure, 82
 - MV/DV error distribution, 82–83
 - predictors hit rate and availability, 83–84
 - 3D-video applications, 2–3
 - 3D-video coding
 - 3D-video pre-and post-processing, 172
 - dynamic search window, 171
 - energy reduction, 169
 - fast ME/DE, 170–171
 - frame and basic unit (BU) level, 170
 - HRC algorithm, 170
 - HVS, 170
 - MDP, 170
 - mode decision algorithm, 169–170
 - MPC, 170
 - multibank video on-chip memory, 171
 - MVC challenges, 172
 - next-generation 3D-video coding, 172–173
 - relax and aggressive strengths, 169–170
 - state of the art, 170–171
 - video sequence, 171
 - 3D-video system, 29–30
 - Dynamic complexity adaptation, 62
 - Dynamic power management, 63
 - Dynamic search window formation-based date reuse, 9
 - Dynamic voltage scaling (DVS), SRAM, 36
- E**
- Energy-efficient algorithms
 - 3D-neighborhood correlation analysis
 - bitrate correlation analysis, 84–87
 - coding mode correlation analysis, 74–81
 - motion correlation analysis, 82–84
 - fast motion and disparity estimation algorithm, 107–109
 - fast ME/DE algorithm results, 109–111
 - multilevel mode decision-based complexity adaptation, 8
 - energy-aware complexity adaptation, 95–100
 - energy-aware complexity adaptation results, 105–107
 - multilevel fast mode decision, 90–95
 - multilevel fast mode results, 100–105
 - multiview video coding
 - mode decision, 40–42
 - motion and disparity estimation, 42–44
 - vs. state-of-the-art
 - energy-aware complexity adaptation, 159–160
 - fast ME/DE, 160–161
 - mode decision, 156–159
 - thresholds
 - probability density function, 87
 - quantization parameter, 88
 - RDCost property of SKIP MB, 88–89
 - video-quality management
 - basic unit-level rate control, 119–121
 - frame-level rate control, 113–118
 - hierarchical rate control, 111–113
 - hierarchical rate control results, 121–126
- Energy-efficient architecture
- hardware, 9
 - multimedia processing
 - DVS, 36
 - dynamic power management, memories, 35–36
 - power-rate–distortion model, 36
 - SRAM dynamic voltage-scaling infrastructure, 34–35
 - video architecture, 37–38
 - video memories, 34
- MVC. (*see* Multiview video coding (MVC))
- vs. state-of-the-art
 - application-aware power gating, 164–165
 - CIF resolution, 164
 - dynamic search window technique, 166
 - motion and disparity estimation, 163–164
- Energy-efficient mode decision scheme, 62
- Energy-efficient motion and disparity estimation, 62
- Energy-efficient on-chip video memory hierarchy, 63
- Energy/power consumption, 38–40

F

- Fast motion and disparity estimation
 - algorithm, 68–69, 107–109
 - algorithm results, 109–111
 - flow diagram, 108
 - vs. TZ search, 109
- Finite state machine (FSM), 131–133
- Frame-level bitrate distribution, 86, 87, 124, 125
- Frame-level rate control
 - bitrate prediction, 114–115
 - diagram, 115
 - evaluation, 117–118
 - model predictive controller, 113–114
 - MPC-based, 113
 - for quality related changes, 65
 - quantization parameter definition, 117
 - rate model, 116
- Free-viewpoint television (FTV), 3, 14

H

- Hierarchical rate control (HRC), 170
 - MBEE, 162
 - model predictive control-based rate control, 69
 - for motion and disparity estimation, 170–171
 - for MVC, 8, 111–113
 - results
 - BD-PSNR and BD-BR comparison, 123
 - bitrate accuracy, 122
 - bitrate and PSNR distribution, 124, 125
 - bitrate distribution at BU level, 124–126
 - view-level bitrate distribution, 124
- Human visual system (HVS), 11, 170

J

- Joint model for multiview video coding (JMVC)
 - communication channels, 178–179
 - CreateH264AVCEncoder class, 175–176
 - encoder tracing, 178
 - H264AVCEncoderTest, 175–176
 - inter-frame search, 176–177
 - ME/DE modification, 179
 - mode decision modification, 179
 - motion estimation, 176
 - PicEncoder, 175–176
 - rate control modification, 179–180
 - SliceEncoder, 176

K

- Key frames (KF)
 - fast ME/DE algorithm scheduling, 138–139
 - MVC encoder, parallelism in, 136

M

- Markov decision process (MDP), 48–49, 113, 119–121, 170
- Mean absolute differences (MAD), 45
- Mean bit estimation error (MBEE), 162
- Memory design methodology, 9
- Mode decision (MD), 55
- Model predictive control (MPC)
 - based frame-level rate control, 46–47, 113, 170
 - goal, 113
- Monograph
 - 3D-neighborhood correlation analysis, 7–8
 - energy-efficient hardware architecture, 9
 - energy-efficient MVC algorithm, 8
 - goal, 7
- Motion and disparity estimation (ME/DE)
 - Bi-prediction, 25–26
 - design
 - AGU, 134–135
 - application-aware power-gating scheme, 128–129
 - architectural template, 129–130
 - DPB, 129
 - dynamic search window formation algorithm, 128
 - energy/complexity-aware control unit, 129
 - multibank on-chip memory, 128
 - on-chip video memory, 133–134
 - programmable search control unit, 131–133
 - SAD calculator, 130–131
 - diamond search, 25
 - enhanced predictive zonal search, 25
 - full search, 24–25
 - hardware architecture, 9
 - log search, 25
 - motion/disparity vector prediction, 26–27
 - multiple block sizes, 26
 - multiple reference frames and reference view, 26
 - pipeline scheduling
 - common vector predictor, 139–140
 - fast operation modes, 139
 - frame-level evaluation, 140
 - generic search pattern, 137–138

- Motion and disparity estimation (ME/DE) (*cont.*)
 - GOP-level, 138
 - SKIP predictor, 140
 - TZ module, 138–139
 - ultra fast, 139
 - zero vector, 139
- quarter-sample motion vector accuracy, 26
- reference frame, 22–24
- search range (SR), 24
- search window (SW), 24
- temporal frame, 22–23
- three step search, 25
- weighted prediction, 26
- Motion correlation analysis
 - basic prediction structure, 82
 - MV/DV error distribution, 82–83
 - predictors hit rate and availability, 83–84
- Multibank on-chip video memory, 9, 70
- Multilevel mode decision-based complexity adaptation, 8
 - energy-aware complexity adaptation algorithm, 99–100
 - employing asymmetric view quality, 95–96
 - mode decision algorithm, 97–99
 - quality-complexity classes, 96–97
 - results, 105–107
 - multilevel fast mode decision
 - early mode decision terminator, 92–94
 - early SKIP prediction, 92
 - high-confidence modes and low-confidence modes, 94
 - RDCost confidence-level ranking, 90–91
 - texture direction-based mode prediction, 94–95
 - video properties-based mode prediction, 94
 - multilevel fast mode results
 - frame-level time saving evaluation, 103–104
 - multilevel mode decision algorithm overhead, 105
 - tested modes evaluation, 100, 102–103
 - view-level time saving evaluation, 100, 102
- Multimedia embedded systems, 5–6
- Multimedia processing architecture
 - ASIC, 32
 - energy management
 - DVS, 36
 - dynamic power management, memories, 35–36
 - power-rate-distortion model, 36
 - SRAM dynamic voltage-scaling infrastructure, 34–35
 - video architecture, 37–38
 - video memories, 34
 - heterogeneous multicore SoC, 33
 - multimedia processors/DSP, 30–31
 - reconfigurable processors, 31–32
- Multiview video coding (MVC)
 - adaptivity in MVC video encoder, 60–61
 - application analysis for energy and quality, 71
 - definition, 16
 - disparity domain correlation, 15, 16, 19
 - 3D-neighborhood, 67–68
 - dynamic search window formation algorithm
 - search map prediction, 142–143
 - search window, 143–145
 - TZ/Log search, 140–142
 - encoding process
 - deblocking filter (DF), 22
 - double-Z processing order, transform module, 21
 - encoder block diagram, 18
 - inter-frame prediction/motion estimation (ME), 19
 - inter-view prediction/disparity estimation (DE), 19–20
 - intra-frame prediction, 19–20
 - macroblock borders filtering, 22
 - mode decision (MD), 20–21
 - PSNR, 20
 - quality vs. efficiency trade-off, 20–21
 - rate-distortion (RD) trade-off, 20
 - Zigzag scan order, 21–22
 - energy-efficient algorithm
 - mode decision, 40–42
 - motion and disparity estimation, 42–44
 - energy-efficient algorithms
 - fast motion and disparity estimation, 68–69
 - hierarchical rate control, 69
 - multilevel mode decision-based complexity adaptation, 68
 - thresholds definition methodology, 69
 - energy-efficient algorithms for
 - 3D-neighborhood correlation analysis, 74–87
 - fast motion and disparity estimation, 107–110
 - multilevel mode decision-based complexity adaptation, 90–107
 - thresholds, 87–89
 - video-quality management, 111–126

- energy-efficient architecture
 - application-aware power gating, 70
 - dynamic search window formation-based date reuse, 70
 - memory design methodology, 70
 - motion and disparity estimation
 - hardware architecture, 69
 - multibank on-chip video memory, 70
 - energy-related challenges in, 62–63
 - energy requirements for
 - component blocks energy breakdown, 54, 55
 - distinct mode decision schemes, 55–56
 - energy breakdown, 56–57
 - energy consumption and battery life, 53–54
 - multiple search window sizes, 54–55
 - evaluation, 148–149
 - H.264 standard, 17
 - ME/DE design
 - AGU, 134–135
 - application-aware power-gating scheme, 128–129
 - architectural template, 129–130
 - DPB, 129
 - dynamic search window formation algorithm, 128
 - energy/complexity-aware control unit, 129
 - multibank on-chip memory, 128
 - on-chip video memory, 133–134
 - programmable search control unit, 131–133
 - SAD calculator, 130–131
 - ME/DE pipeline scheduling
 - common and predictor vector, 139–140
 - fast operation modes, 139
 - frame-level evaluation, 140
 - generic search pattern, 137–138
 - GOP-level, 138
 - SKIP predictor, 140
 - TZ module, 138–139
 - ultra fast, 139
 - zero vector, 139
 - memory access, 59–60
 - mode decision (MD), 27
 - motion and disparity estimation
 - Bi-prediction, 25–26
 - diamond search, 25
 - enhanced predictive zonal search, 25
 - full search, 24–25
 - log search, 25
 - motion/disparity vector prediction, 26–27
 - multiple block sizes, 26
 - multiple reference frames and reference view, 26
 - quarter-sample motion vector accuracy, 26
 - reference frame, 22–24
 - search range (SR), 24
 - search window (SW), 24
 - temporal frame, 22–23
 - three step search, 25
 - weighted prediction, 26
 - objective quality analysis for, 63–65
 - on-chip video memory
 - application-aware power gating, 146–147
 - design, 145–147
 - parallelism, 136–137
 - quality-related challenges in, 65
 - rate control (RC), 28–29
 - vs. simulcast approach, 16, 17
 - vs. simulcast computational complexity, 57–58
 - spatial domain correlation, 14–15
 - temporal correlation, 15
 - video quality
 - frame-and BU-level RC, 45–46
 - MAD, 45
 - MDP, 48–49
 - MPC, 46–47
 - pyramid-based priority structure, 45–46
 - quantization parameters (QP), 45–46
 - region of interest, 49–50
 - reinforcement learning model, 49
 - viewer tool
 - Java virtual machine, 181
 - macroblocks, 182–183
 - main screen, 181–182
 - search window-based analysis, 182–184
 - YUV, RGB and YCgCo, 16–17
 - MVC. *See* Multiview video coding (MVC)
- N**
- Non/key frames (NKF)
 - fast ME/DE algorithm scheduling, 138–139
 - MVC encoder, parallelism in, 136
- O**
- On-chip video memory
 - application-aware power gating, 146–147
 - as cache memory, 133–134
 - design, 145–147

On-chip video memory (*cont.*)
 energy-efficient hierarchy, 63
 multibank, 9, 70
 on-chip memory design, 145–146
 organization, 134

P
 Picture order counter (POC), 134
 Prefetching technique, 63

Q
 Quality-complexity classes (QCC)
 mode decision algorithm for
 pseudo-code of, 98
 RDCost characterization, 97–100
 types, 96–97
 Quality states (QS), 99–100, 105
 Quantization parameter (QP)
 based thresholding, 65
 definition, 117

R
 Region of interest (RoI), 45, 49–50, 87,
 119, 120
 Rotating instruction set processing platform
 (RISPP), 31–32

S
 Search pattern memory (SPM), 131–133
 Sobel filter, 154–155
 Spatial domain correlation, 14–15
 Spatial–temporal–disparity index, 153–155
 SRAM. *See* Static random access
 memory (SRAM)
 State-of-the-art
 benchmark video sequence, 152–155
 vs. energy-efficient algorithm
 energy-aware complexity adaptation,
 159–160
 fast ME/DE, 160–161
 mode decision, 156–159
 vs. energy-efficient hardware architecture,
 163–166
 fairness of comparison, 155
 fast ME/DE algorithms, 43
 hardware description and ASIC synthesis,
 155
 software simulation, 151–153
 vs. video quality control algorithm, 161–164

Static random access memory (SRAM)
 application-aware power gating, 147
 ASIC synthesis, 155
 energy-efficient architecture, 164
 fast ME/DE algorithm, 149
 on-chip memory, 34
 on-chip video memory, 130, 133–134
 power management, 35–36
 voltage-scaling infrastructure, 34–35

T
 Temporal domain correlation, 15
 Three-dimensional surveillance, 3
 Three-dimensional telemedicine, 3
 Three-dimensional telepresence, 3
 Three-dimensional television (3DTV), 2
 Three-dimensional video personal
 recording, 2
 Thresholds
 probability density function, 87
 quantization parameter, 88
 RDCost property of SKIP MB, 88–89

V
 Video-quality management
 basic unit-level rate control
 coupled reinforcement learning, 121
 diagram, 119
 Markov decision process, 119–121
 regions of interest, 120
 frame-and BU-level RC, 45–46
 frame-level rate control
 bitrate prediction, 114–115
 diagram, 115
 evaluation, 117–118
 model predictive controller, 113–114
 quantization parameter
 definition, 117
 rate model, 116
 hierarchical rate control, 111–113
 hierarchical rate control results, 121–126
 MAD, 45, 46
 MDP, 48–49
 MPC, 46–47
 pyramid-based priority structure, 45–46
 quantization parameters (QP), 45–46
 region of interest, 49–50
 reinforcement learning model, 49
 vs. state-of-the-art, 161–164
 Video scaling, 4
 View-level bitrate distribution, 85