

FAST DEPTH CODING IN 3D-HEVC USING DEEP LEARNING

A DISSERTATION
SUBMITTED TO THE
DEPARTMENT OF ELECTRONIC AND INFORMATION ENGINEERING
OF THE HONG KONG POLYTECHNIC UNIVERSITY
IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE DEGREE OF
MASTER OF SCIENCE

Zhen-xiang WANG

November 2017

CERTIFICATE OF ORIGINALITY

I hereby declare that this dissertation is my own work and that, to the best of my knowledge and belief, it reproduces no material previously published or written nor material which has been accepted for the award of any other degree or diploma, except where due acknowledgement has been made in the text.

_____(Signed)

_____(Date)

Abstract

The 3D Extension of the High Efficiency Video Coding standard (3D-HEVC), which has been finalized by the Joint Collaborative Team on Video Coding (JCT-VC) in February 2015, is the new industry standard for 3D applications. The 3D-HEVC provides plenty of advanced coding tools specifically for addressing the coding of auto-stereoscopic videos which have the format of multiple texture views along with the depth maps which are responsible for synthesising intermediate views with sufficient quality for auto-stereoscopic display. The provided tools take advantage of the statistical redundancies amongst texture views and depth maps in the video sequences, as well as the unique characteristics of depth maps to significantly shrink the bit-rate while preserving the objective visual quality of the 3D videos. However, those tools with high capability in terms of compression come with the high complexity of computation which has made the encoding time of the 3D video sequences much longer than ever by traversing a lot more candidates, calculating time-consuming RD Cost for each of them, especially in the wedgelet searching process for depth maps. While this full-search style method can promise to find the best candidate in depth intra mode decision, the time cost is expensive.

In this dissertation we address the time cost by presenting a new intra mode decision method for depth maps, leveraging the deep convolutional neural networks to predict the wedgelet angles for the depth blocks. The predictions from the learned models are capable of reducing the number of wedgelet candidates by half as well as the angular modes in depth map coding. The size of the neural network has been carefully designed to balance the trade-off between the time cost of model prediction and the model prediction accuracy. Confusion matrix is used to monitor the training process. Top-K criteria is employed for the prediction. We have integrated the learned models into the reference software of 3D-HEVC for the experiments. The compiled executable binaries are able to harness the power of the simultaneous computation of CPU, as well as the parallel computation of GPU to accelerate the predictions. The simulation results show that the proposed algorithm provides 64.6% time reduction in average while the BD performance has a tiny decrease comparing with the state-of-the-art 3D-HEVC standard.

Acknowledgments

First and foremost, I would like to give sincere thanks to my supervisor, Dr.Yui-Lam Chan, for his extremely generous support, most insightful advices and innumerable yet constructive feedback. I learned from him to first identify a problem, by reading a vast amount of articles to know what people have achieved and what bottlenecks they have encountered. I learned how to read papers, how to organize them to become the inner comprehension. He guided me to use the machine learning approach to solve the problem that has been found in the first stage. Without his guidance I will not have the idea to learn the deep learning technology and apply it to optimize the video coding. His encyclopedic knowledge and charming personalities made him my mentor in both research and life. I wish to thank Dr.Sik-Ho Tsang, for our in-depth discussions from which I can always find useful clues to proceed to next step. His great expertise in video coding significantly benefits me during my intensive period of learning. Also I would like to thank my friends Alex and Jacky, for our extensive discussions about artificial intelligence and their applications. Finally thank you my parents, for the great love and constant encouragement which give me confidence to face and handle all the challenges at every moment.

Contents

Abstract	iii
Acknowledgments	iv
1 Introduction	1
1.1 Motivation	5
1.2 Contribution and Dissertation Outline	8
2 Background	11
2.1 Video Coding	11
2.2 Deep Learning	13
3 Prepare the Data for Deep Learning	16
3.1 Vicodeo Coding	16
4 Train the Deep Model for Prediction	17
4.1 Video Cnbvnbvnboding	17
5 Evaluate the Learned Deep Model	18
5.1 Video Codhgvingg	18
6 Employ the Learned Deep Model	19
6.1 Video Cobnbvnbvnbvnding	19
7 Conclusion	20
7.1 Video Codvbnvnbvnbvnbng	20
Bibliography	21

List of Tables

1.1	Characteristics Comparison of Stereoscopic Display and Autostereoscopic Display . .	2
1.2	Impact of Available View Amount for Autostereoscopic Display	3
1.3	The summary of the time percentages occupied by DMM1 searching in the process for compressing CUs	7
1.4	The summary of the time percentages occupied by VSO in DMM1 searching	7

List of Figures

1.1	System Structure for transmitting videos targeting stereo display	2
1.2	System Structure for transmitting videos of Multi-view Plus Depth format	3
1.3	Wedgelet partition illustration	4
1.4	Contour partition illustration	5
1.5	An example showing a piece of the command line outputs during the encoding process for Shark sequence	6
1.6	A screen capture of the time profiling information for Newspaper sequence	6
1.7	Flowchart for proposed fast depth coding algorithm	9
2.1	The brief history of the video coding standards	12

Chapter 1

Introduction

Video is the medium to record, copy, playback, broadcast and display the motion images in an electronic style [1]. Watching videos is becoming an important way for our entertainment as well as education. The high definition (HD) and ultra high definition (UHD) video are increasingly demanding nowadays. People prefer videos with higher definitions than those with lower resolutions because the former one provides much better viewing experience. However, challenges emerged for delivering videos with high definition. HD videos typically contain much more information in every picture frame than the standard definition videos. More data needs to be squeezed into the same capacity for transmission. For example, the uncompressed video with the dimension 720 x 480 at 30 frames per second requires 0.03 gigabytes per second, while the uncompressed video with the dimension 2880 x 2048 at 120 frames per second requires 2.12 gigabytes per second. Since bit rate is proportional to system bandwidth for transmission [2], and expanding the bandwidth in a large scale is too expensive, the significantly increased bit rate for transmitting the video data is becoming one of the major obstacles for HD video services.

To cope with the growing need for higher compression of moving pictures [3], Joint Collaborative Team on Video Coding (JCT-VC) [4] has developed the High Efficiency Video Coding standard which is the newest international video coding standard for substantially ameliorate the compression performance against the previous standards. Comparing with the H.264 Advanced Video Compression Standard [5], the H.265 High Efficiency Video Coding Standard provides fifty percent bit rate reduction while maintaining the objective video quality at the same level.

While Two-dimensional video is the most common video type, Three-dimensional (3D) video has been brought to market via lots of ways, including Blu-Ray disc, cable and satellite transmission, terrestrial broadcast, and streaming or downloading from the Internet [6]. 3D video provides the perception of depth information which augments the vividness of the video contents. Currently most 3D videos in the market are using stereo display technology. Two similar views, one for left eye, the other for right eye, are presented at the same time with the multiplexing techniques

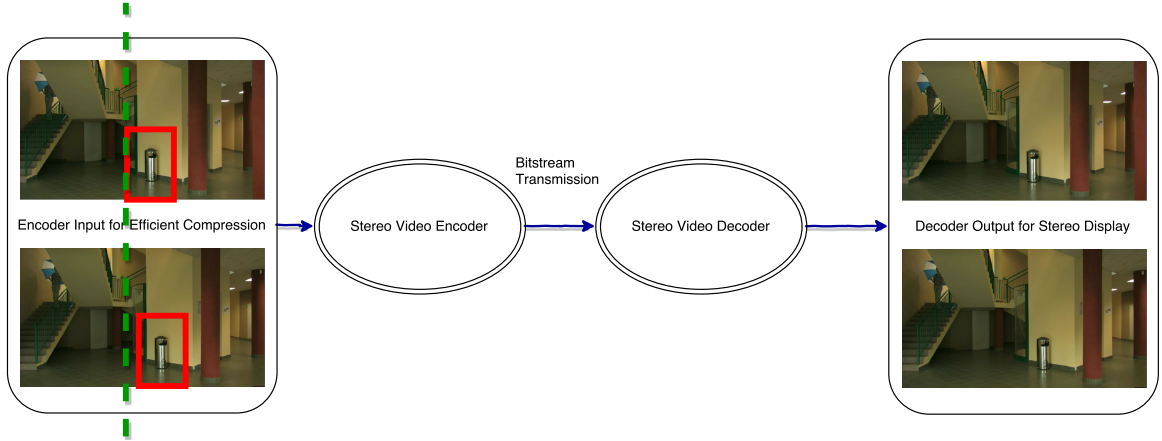


Figure 1.1: System Structure for transmitting videos targeting stereo display.

enabling the adjustments of video geometry information [7] to provide the 3D effect. Figure 1.1 illustrates the typical system structure for transmitting videos targeting stereo display. It can be observed that there exists a displacement between the two views. The green vertical left margins of the red rectangles in the two views at encoder side are different. Such a displacement is the visual disparity for 3D perception. Stereoscopic videos [8] have achieved great profitability for movie theatres in recent years. For example, IMAX 3D has become the most popular one that offering the immersing multimedia experiences around the world. Special 3D glasses are needed for watching the IMAX 3D movies. The current 3D film industry is very successful in terms of attracting customers, however, it is not the end of the story. Myopic people do not like to wear one more pair of glasses when watching 3D movies. Some people will experience discomfort after wearing the 3D glasses for a period of two hours. To get rid of the undesired 3D glasses, autostereoscopic multi-view technology [8] is coming to our rescue. The two major different characteristics between stereo display and autostereoscopic display are listed in Table 1.1 [9]. The impact of different view numbers for autostereoscopic display is shown in Table 1.2 [9]. Comparative ease can be brought to the 3D video audience since they do not need to wear 3D glasses for watching autostereoscopic videos. At each different view position, scenes with minor differences are available from multiple stereo pairs which are provided by autostereoscopic display [9]. As a result, when audience make a move for various view positions, scenes not viewable from the previous locations are revealed during the movement.

Table 1.1: Characteristics Comparison of Stereoscopic Display and Autostereoscopic Display

Characteristic	Stereo Display	Autostereoscopic Display
Glass-Free	No	Yes
Multiple Stereo Pairs	No	Yes

Table 1.2: Impact of Available View Amount for Autostereoscopic Display

Characteristic	Small Number of Views	Large Number of Views
Seamless View Transition	No	Yes
High Quality of Scene Depth	No	Yes

The autostereoscopic multi-view display demands more than two views. With a sufficient amount of views present in autostereoscopic display, the disparities between every two adjacent views can be small enough to offer seamless transitions from scene to scene, such that when multiple views meet eyes sequentially, the scenes as a whole can be gorgeous. The visual quality of the autostereoscopic display is highly proportional to the number of available views. Due to limited available bandwidth, transmitting arbitrary number of views is not practical. Researchers have proposed a new format which only requires limited number of view and their associated depth maps for the capability of generating arbitrary amount of views theoretically. The typical system structure for using this new format to compress and supply 3D video resources is shown in Figure 1.2. An enormous amount of views in the medium positions which are able to guarantee the high quality of the 3D video can be synthesized from the decoded texture frames in combination with decoded depth maps.

To employ multi-view plus depth format for 3D video, efficient compressing methods are desired, which has led to the 3D Video Coding Extension of the High Efficiency Video Coding Standard (3D-HEVC) by the Joint Collaborative Team on 3D Video Coding Extension Development (JCT-3V) [10]. The 3D Extension of the HEVC standard gives extra coding efficiency for encoding a few texture views along with the corresponding depth maps by using new tools which exploit the redundancies amongst texture and depth views, and pay attention to the unique characteristics of

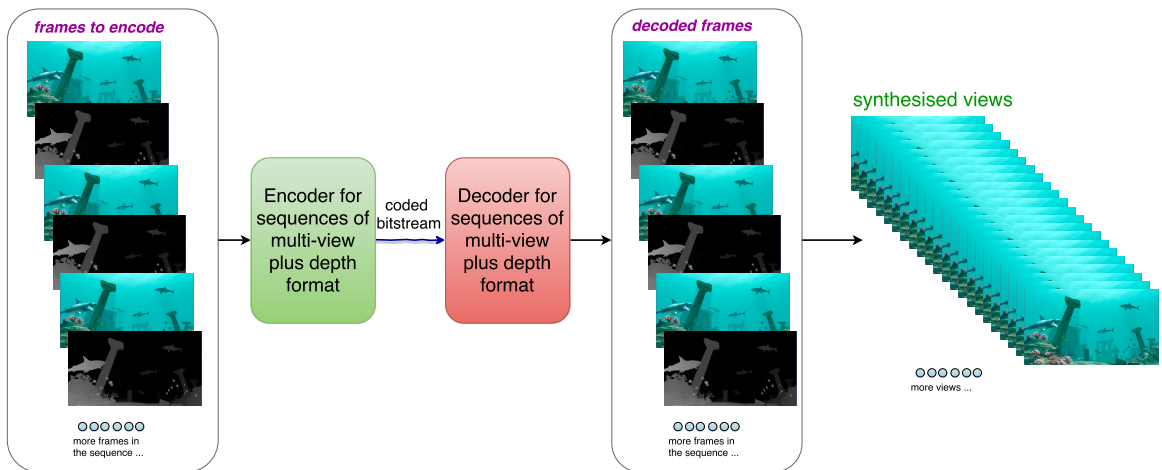


Figure 1.2: System Structure for transmitting videos of Multi-view Plus Depth format.

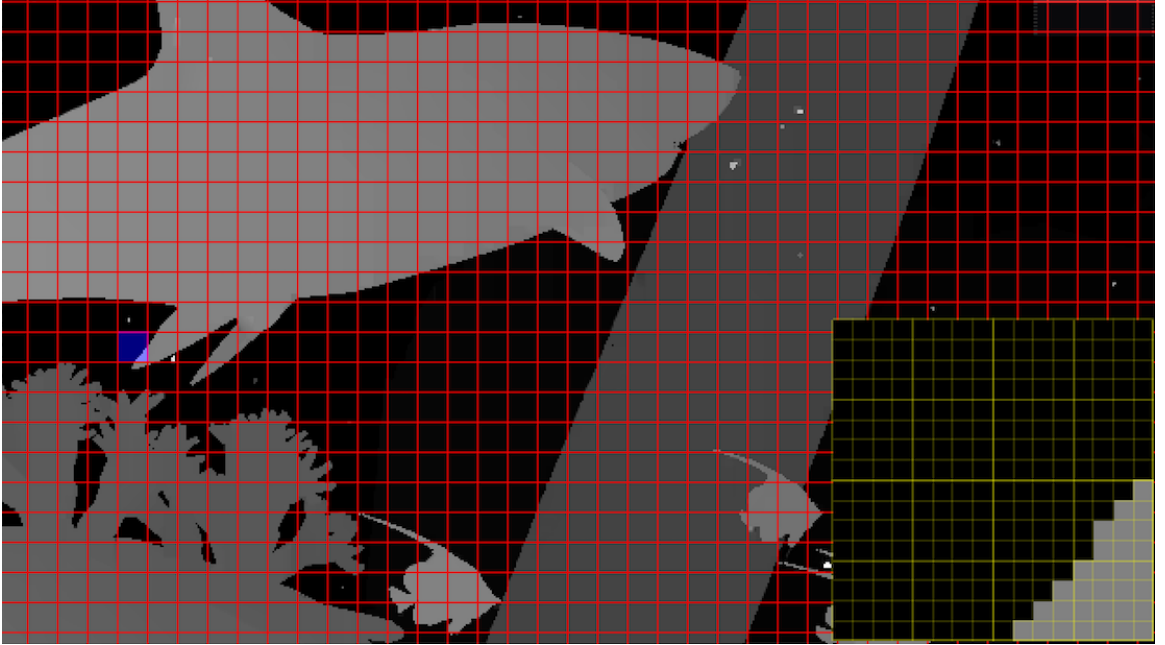


Figure 1.3: Example of wedgelet partition in a block of size 16 by 16 in depth map from Shark video sequence.

the depth maps, such as large homogeneous regions separated by sharp boundaries [11].

Depth information measures of the distance between the object in the far position and the object in the near position from a static viewpoint, which is expressed in the format of depth map. Instead of presenting depth maps directly to the viewer, views in the medium positions are generated by Depth-Image-Based Rendering (DIBR) technique. The qualities of the depth maps are vital to the DIBR process. Corona artifacts (a.k.a. ringing artifacts) can be discovered in synthesized views if the edge sharpness in depth maps can not be well preserved. Therefore, retaining the edge sharpness in depth map is the key to avoid the artifacts in the synthesized views. In 3D-HEVC, new intra-picture prediction tools and residual coding methods have been applied to preserve the special properties of depth maps. Depth Modelling Mode (DMM) which is one of the new intra-picture prediction tools, is designed to provide much more granularity for encoding the depth maps than the normal angular intra prediction modes. DMM is more capable of approximating the depth maps to be encoded due to the fact that it provides a vast amount of non-rectangle partitions. Figure 1.3 presents an example of the wedgelet partition from the depth map in Shark video sequence. The small block highlighted by blue color amongst the blocks separated by the red grid is magnified at the right-bottom position in Figure 1.3. A straight line is used for the partition in wedgelet mode. Figure 1.4 shows a sample of the contour partition from the same depth map as Figure 1.3. The partition pattern comprises contour lines instead of one single straight line. Wedgelet partition and contour partition for depth

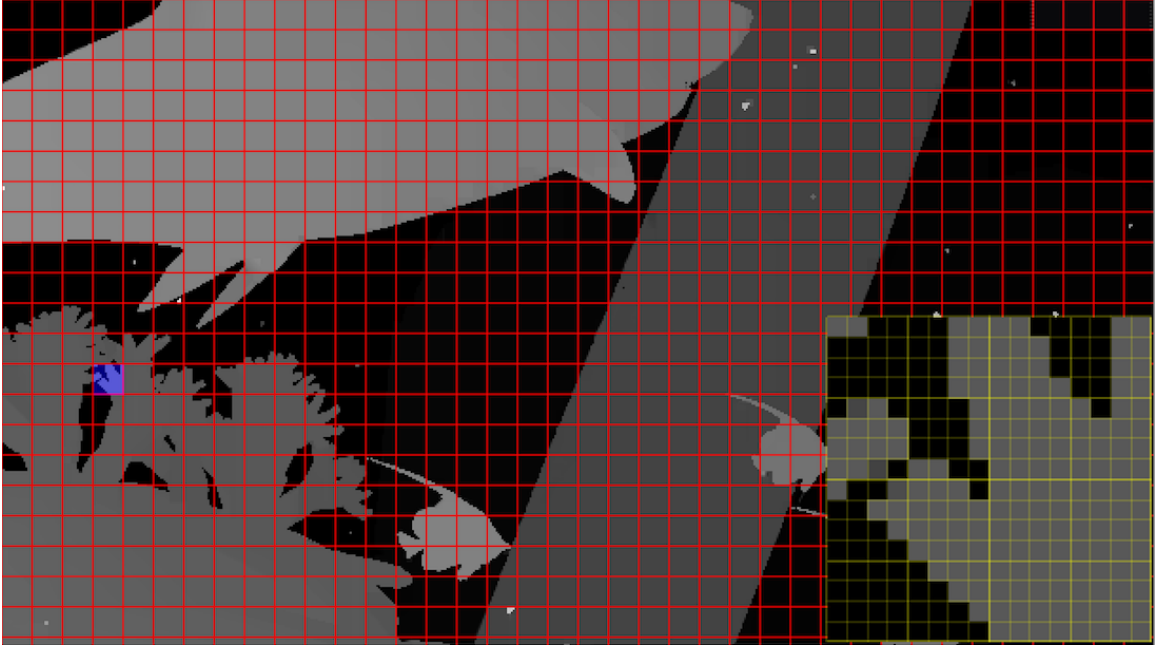


Figure 1.4: Example of contour partition in a block of size 16 by 16 in depth map from Shark video sequence.

maps are enabled by DMM1 and DMM4 separately.

1.1 Motivation

The idea of this work originates from the discovery of the computational complexity of the wedgelet searching process in depth modelling modes. The immense complexity for searching the best wedgelet candidate lead to the strongly marked increase of encoding time. The time consumed for compressing a single depth map in 3D-HEVC encoder is roughly a sixfold increase relevant to the encoding time of a single texture frame wherein the all intra configuration in HTM-16.2 is used. Thus we designed a deep neural network architecture which is trained subsequently for predicting the most probable wedgelet candidates. The learned model achieves 92.2% to 97.3% top-16 accuracy for various block sizes. The inference engine is integrated into the reference software (HTM-16.2) of 3D-HEVC. The learned models reduce roughly half of the wedgelet searching candidates. It provides 64.6% time reduction in average while the BD performance has a negligible decrease comparing with the unmodified 3D-HEVC encoder.

Motivation for Wedgelet Candidates Reduction: The encoding time consumed in HTM-16.2 encoder for each view (including both texture and depth) by default can be observed from the command line outputs. Figure 1.5 shows a piece of command line outputs from the encoding process

Layer	POC	TId	Encoding Time (s)	Layer Id	Picture Order Count (POD)
0	143	0	932224	0	0
1	143	0	60160	0	0
2	143	0	932592	0	0
3	143	0	59616	0	0
4	143	0	932312	0	0
5	143	0	61920	0	0
0	144	0	933808	0	0
1	144	0	62288	0	0
2	144	0	932096	0	0
3	144	0	61496	0	0
4	144	0	931432	0	0
5	144	0	62440	0	0
0	145	0	927656	0	0
1	145	0	64216	0	0
2	145	0	928776	0	0
3	145	0	62608	0	0
4	145	0	929248	0	0
5	145	0	65928	0	0
0	146	0	928136	0	0
1	146	0	62480	0	0
2	146	0	923304	0	0
3	146	0	61344	0	0
4	146	0	927120	0	0
5	146	0	65248	0	0
0	147	0	918384	0	0
1	147	0	64480	0	0
2	147	0	926144	0	0
3	147	0	64000	0	0
4	147	0	923752	0	0
5	147	0	64664	0	0

Figure 1.5: An example showing a piece of the command line outputs during the encoding process for Shark sequence.

of Shark sequence. The numbers in red blocks stands for the encoding time of certain views, while the corresponding layer Id and Picture Order Count (POD) are in the green blocks. A repetitive pattern of the encoding time for each view can be observed every six numbers vertically. A simple calculation using six numbers within the top-most red block, $(90+69+69)/(90+69+69+14*3) \approx 0.84$, shows that approximately 84% of the total encoding time is busy with encoding the depth maps. Similarly, it is reported in [12] that the coding for depth map consumes near 86% of total 3D-HEVC encoding time. A trial of time profiling for 3D-HEVC encoder is performed using Instruments which is available on macOS. After encoding the Newspaper sequence for more than one hour, Figure 1.6 clearly shows 97.8% time is used to compress the CUs recursively. The first recursive xCompressCU function (denoted as XC1 thereafter) is for CUs of size 64x64, the second recursive xCompressCU

Weight	Symbol Name
66.42 min 100.0%	TAppEncoder (11026)
66.42 min 100.0%	Main Thread 0x2a7867
66.42 min 99.9%	main TAppEncoder
66.42 min 99.9%	TAppEncTop::encode() TAppEncoder
66.38 min 99.9%	TEncTop::encode(bool, TComPicYuv*, TComPicYuv*, InputColourSpaceConversion, TComList<TComPicYuv*> &, std::__1::list<AccessUnit> &) TAppEncoder
66.38 min 99.9%	TEncGOP::compressPicInGOP(int, int, TComList<TComPic*> &, TComList<TComPicYuv*> &, std::__1::list<AccessUnit>, std::__1::all) TAppEncoder
66.29 min 99.8%	TEncSlice::compressSlice(TComPic*, bool, bool) TAppEncoder
65.04 min 97.9%	TEncCu::compressCtu(TComDataCU*) TAppEncoder
65.02 min 97.8%	TEncCu::xCompressCU(TComDataCU* &, TComDataCU* &, unsigned int, PartSize) TAppEncoder
58.61 min 88.2%	TEncCu::xCompressCU(TComDataCU* &, TComDataCU* &, unsigned int, PartSize) TAppEncoder
42.71 min 64.3%	TEncCu::xCompressCU(TComDataCU* &, TComDataCU* &, unsigned int, PartSize) TAppEncoder
28.46 min 42.8%	TEncCu::xCompressCU(TComDataCU* &, TComDataCU* &, unsigned int, PartSize) TAppEncoder
13.53 min 20.3%	TEncCu::xCheckRdCostIntra(TComDataCU* &, TComDataCU* &, PartSize, bool) TAppEncoder
25.24 s 0.6%	TEncCu::xCheckRdCostDIS(TComDataCU* &, TComDataCU* &, PartSize) TAppEncoder
7.76 s 0.1%	TComRdCost::setRenModelData(TComDataCU const* &, unsigned int, short const* &, int, int) TAppEncoder
2.69 s 0.0%	TComDataCU::initSubCU(TComDataCU* &, unsigned int, unsigned int, int) TAppEncoder
1.96 s 0.0%	TComDataCU::initEstData(unsigned int, int, bool) TAppEncoder
1.04 s 0.0%	TComDataCU::xPartFrom(TComDataCU* &, unsigned int, unsigned int) TAppEncoder

Figure 1.6: A screen capture of the time profiling information for Newspaper sequence.

Table 1.3: The summary of the time percentages occupied by DMM1 searching in the process for compressing CUs

size of CU	Recursive xCompressCU Function	Time percentage of DMM1 searching process
32 by 32	XC2	30.0%
16 by 16	XC3	25.6%
8 by 8	XC4	18.8%

Table 1.4: The summary of the time percentages occupied by VSO in DMM1 searching

size of CU	process	Time percentage of VSO in DMM1 searching
32 by 32	VSO in DMM1 searching from XC2	80.1%
16 by 16	VSO in DMM1 searching from XC3	83.7%
8 by 8	VSO in DMM1 searching from XC4	78.8%

(denoted as XC2 thereafter) is targeting CUs of size 32x32, the third one (denoted as XC3 thereafter) is dedicated to CUs of size 16x16, and the last one (denoted as XC4 thereafter) is bound to CUs of size 8x8. It is observed that the most time consuming part during the process of compressing the depth CUs is DMM1 searching. The DMM1 searching time percentages are summarised in Table 1.3 wherein the summary for XC1 is omitted since DMM1 is not applicable to CUs of size 64 by 64 in HTM-16.2. [t] The major reason leading to the time consuming property of DMM1 searching is the View Synthesis Optimization (VSO) Method for improving quality of synthesized views [13], wherein the Synthesized View Distortion Change (SVDC) is computed. The time percentages of the VSO processes in DMM1 searching are summarised in Table 1.4. In HTM-16.2, many wedgelet candidates are evaluated using the VSO which has a high computational complexity. Evaluating less wedgelet candidates will help to relieve the burden of heavy computation required by VSO, thereby certain time reduction can be achieved.

Motivation for Using Deep Learning: Deep learning such as Multi-Layer Perceptron (MLP) is a subfield of representation learning, which is in turn a major subset of machine learning [14]. Machine learning such as the support vector learning [15] is applied to many methods in the domain of Artificial Intelligence (AI). Deep learning based on back propagation training has been found hard to proceed in the late 1980s [16], however, starting from 2012, it kicks off the glorious comeback. The deep Convolutional Neural Network (CNN) has won the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) from 2012 to 2015 with the CNN architecture of the winner going deeper and deeper year by year. The great achievements attract attentions from people all over the world and make deep learning the hottest topic in our daily lives. Inspired by the fact that supervised deep learning can learn multiple layers of abstract representations in the visual recognition tasks, it should be applicable to recognize the angular modes of the intra-picture prediction in the 3D-HEVC. The final DMM1 candidate selected in depth map coding is essentially determined by the angle pattern

of the depth blocks. If we can make use of deep learning to predict the most probable angles of the target pixel block, a vast amount of angular modes and DMM1 wedgelet candidates can be naturally skipped by which the time saving can be achieved without decreasing the coding performance.

Motivated by the discussions above, we adopt deep learning approach with deep convolutional neural network to accelerate the depth map coding in 3D-HEVC.

1.2 Contribution and Dissertation Outline

We accelerate the depth map coding in 3D-HEVC leveraging the power of deep learning. The contributions of the dissertation are:

- A deep convolutional neural network with 32 layers comprising ResNet units [17] has been designed and trained for recognizing the angular directions of the blocks from intra-picture prediction in 3D-HEVC encoder. The learned models have high top-k precisions which work well on the tasks of recognizing intra angular patterns in 3D-HEVC.
- A way of integrating the learned model into the HTM-16.2 encoder has been suggested. By making use of Bazel [18] to compile the encoder binary, the data level parallelism (instead of concurrency) functionality in CPU as well as the parallel architecture in GPU are fully utilized for efficient computations of matrix operations.
- An algorithm, illustrated in Figure 1.7 on page 9, for fast depth map coding based on the predictions from learned deep models has been proposed and implemented. The simulation results show that the proposed algorithm is capable of reducing 64.6% time in wedgelet searching during 3D-HEVC encoding process while the BD performance only has a trivial decrease.

The first two contributions lay the foundation for the third one, which is the main objective of this work: to accelerate the depth map encoding process in 3D-HEVC.

Chapter 2 supplies the background of video coding history, video coding standards, and deep learning using artificial neural network. Prior arts in video coding and deep learning are surveyed in this chapter.

Chapter 3 describes the methodology which has been implemented to collect the data to be used for deep learning. The pre-processing steps for the data are provided in details along with the reasons behind the scene. We also visualize lots of collected data to help with the understanding of their properties.

Chapter 4 presents the designed deep convolutional neural network which has been adopted in the deep learning process. Discussions on choosing proper hyper-parameters for the devised neural network are given. The stopping criteria are presented with the training results.

Chapter 5 provides evaluation results for the learned models. Block resizing with different approaches are compared using the accuracy of prediction.

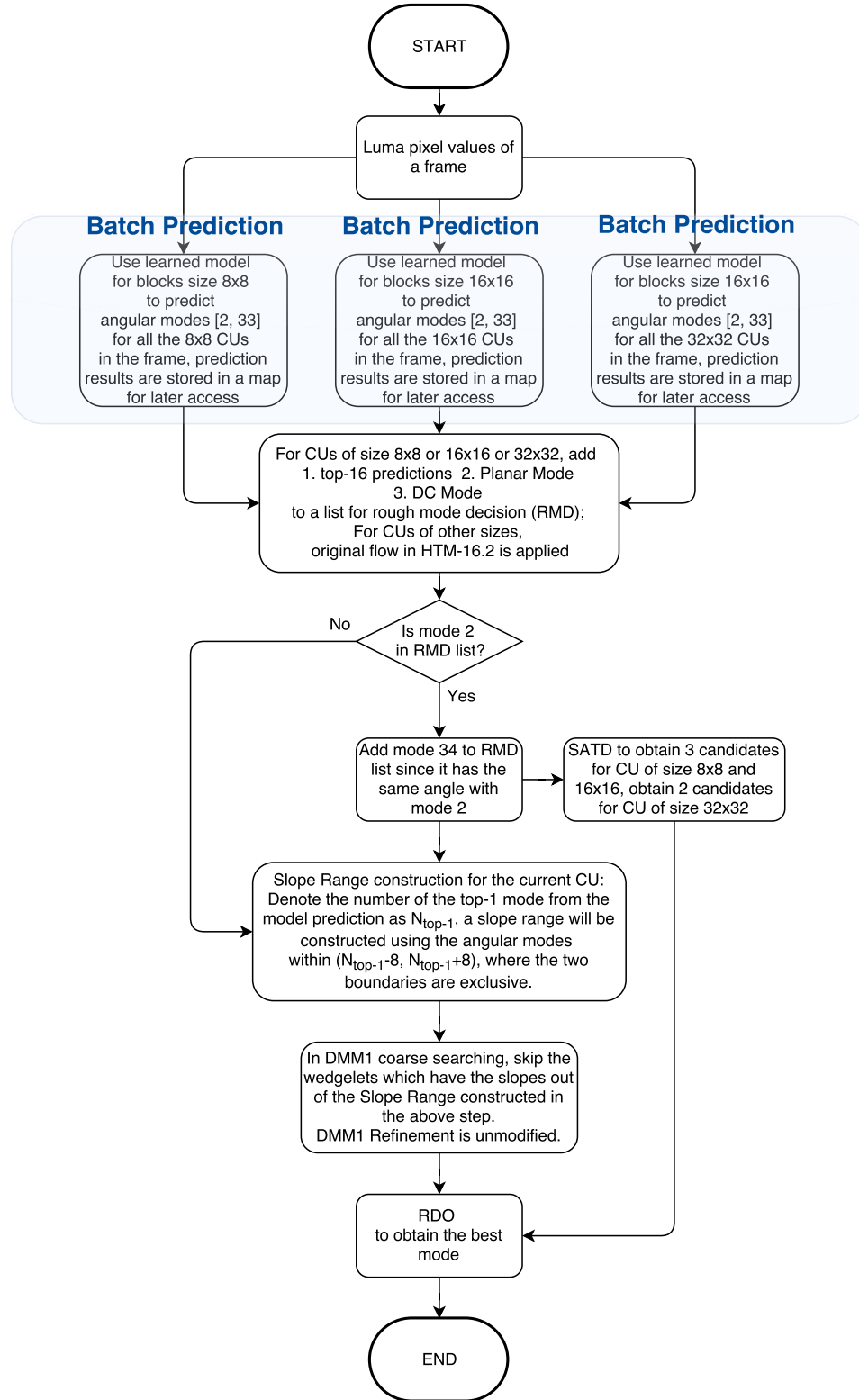


Figure 1.7: Flowchart for proposed fast depth coding algorithm.

Chapter 6 shows the methods been used for integrating the learned model into the 3D-HEVC encoder. Advantages and drawbacks of the integration are discussed. Simulation results comparing with the original HTM-16.2 are given in this chapter.

Chapter 7 concludes the thesis and discusses the future work for fast depth coding using deep learning.

Chapter 2

Background

To start with, we bring up what is video coding, why it is needed, and its challenges. Next we discuss what is deep learning, the history of deep learning and how it works for vision tasks. Furthermore, we introduce how we plan to apply deep learning to optimize video coding tasks and why it should work. In the end, a survey of related works in video coding and deep learning is given.

2.1 Video Coding

Video playback is the most straightforward way for human to perceive dynamic scenes that exist across a time series. More than half of the neurons in human brain are born to process the visual information which is supplied by human eyes. It becomes effortless for human to understand things presented by the video playback instead of a long paragraph of words. Videos are made up of consecutive sets of image frames, which in turn are made up of pixel matrices. Visual information of a cosmic scale is first stored by various methods then delivered during a period of video playback.

In 1950s, video tapes were employed to store the videos. Video tape is able to serve for about eight to twelve years before the video quality starts to degrade. In 1970s, laser disc appeared in the US market as an alternative of video tapes. Start from laser disc, the video storage started its new era in digital world. In 1990s, DVDs were released after laser disc. Data is stored in spiralling tracks on the disc. A laser beam can be utilized to read the data. In addition, hard drives, flash drives and SD cards were also starting to become popular in the late 90s. Nowadays, the cloud storage is very common in daily lives. It is capable of storing data on the servers which are accessible from any devices via internet connections.

Although so many formats are available for video storage, they share a common feature: the more storage you use, the more cost it will be. Let's take the cloud storage as an example. Google cloud is one of the most popular cloud services in our daily lives. It provides cloud storage with a price of \$0.026 per GB/month [19] (this price is observed on 21 Nov 2017, it may change in the

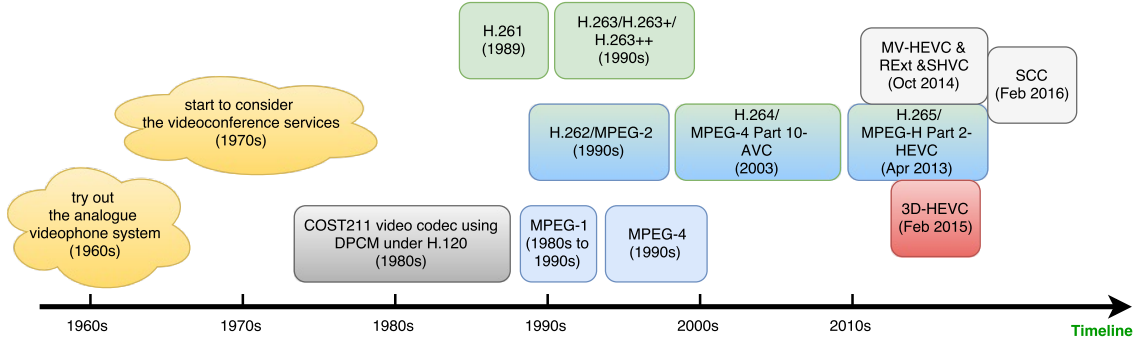


Figure 2.1: The brief history of the video coding standards

future). If a 4K video with a resolution of 4096×2160 , at 120 frames per second, 8 bits for each of the RGB component, needs to be stored without any compression in Google cloud, we need to pay a monthly fee: $(4096 * 2160 * 120 * 60 * 90 * 3 * 0.026) / (1024 * 1024 * 1024) \approx 416.47$ \$. Without doubt, this figure is relatively not acceptable for just storing the video. High compression is needed to store the videos in a practical way.

From the other perspective, let us take the bandwidth into consideration. To deliver the uncompressed 4K video which has been mentioned in the previous paragraph, we need a bandwidth of: $(4096 * 2160 * 120 * 3) / (1024 * 1024 * 1024) \approx 2.97$ Gigabytes per second. The maximum bandwidth of Wireless 802.11ac, which is one of the common internet access technologies, is 1.3 Gigabytes per second [20]. Apparently, the wireless connection is not able to deliver such kind of 4K videos. High compression is desired to deliver the video through the internet.

Despite the fact that raw videos usually contain a large amount of data, a lot of redundancies exist. For every video sequence, two types of redundancies are ubiquitous: Spatial Redundancy and Temporal Redundancy. Video coding technologies are taking advantages of those redundancies to achieve the efficient compression for video data. Many of the useful video coding technologies have been adopted by the international video coding standards, such as MPEG-4, H.264, H.265, etc.

Figure 2.1 shows the brief history of the video coding standards. In 1980s, the COST211 video codec, built on top of Differential Pulse Code Modulation (DPCM), was standardized under H.120 standard by CCITT (now known as ITU-T). In late 1989, the H.261 was completed and its success marked a milestone for video coding at low bit rate with fairly good quality [21]. The Motion Picture Experts Group (MPEG) kicked off the exploration of video storage, such as CD-ROMs. Their objective was to achieve a competitive performance with cassette recorders in terms of compression of videos which have rich motions. The framework of H.261 had been used to start the codec design of MPEG-1. MPEG-2 was one generation after the MPEG-1. It featured higher capabilities when handling videos with high bit rates and high resolutions. In MPEG-2, the encoder is allowed to make its own decision on the the number of bi-directionally predicted pictures according to a suitable

coding delay. ITU-T found this technique applicable to telecommunication applications, as a result MPEG-2 has been adopted as H.262 for telecommunications. Right after the MPEG-2 standard, MPEG-3 was designed mainly for coding of high definition videos. However, MPEG-3 was discarded due to the versatility of MPEG-2, which can be used to encode videos of any resolutions. In the late 1998, MPEG-4 was introduced as a way of defining compression of both audio and visual digital data. Later on MPEG-4 was divided into several parts during its continuously evolving. Among its sub-parts, MPEG-4 part 10 (a.k.a. Advanced Video Coding) is mainly for the video compression. With the rising popularity of the high definition videos, the new standard termed High Efficiency Video Coding (HEVC) for compressing videos in a more efficient way comparing with previous standards, such as H.264/AVC, has emerged under the efforts from the Joint Collaborative Team on Video Coding (JCT-VC). In the meanwhile, five extensions of the HEVC standard, comprising Format Range Extension (RExt), Scalability Extension (SHVC), Multi-view Extension (MV-HEVC), 3D Extension (3D-HEVC), Screen Content Coding Extension (SCC), have been finalized from 2014 to 2016 to fulfill extra requirements in various video coding scenarios.

In this work, we focus on the depth map coding in 3D-HEVC. The 35 angular modes and depth modeling modes have been embraced in the depth map coding tools in 3D-HEVC. The DMM1 mode introduces an huge increase for the encoding time of 3D videos. Acceleration of the depth map coding is needed.

2.2 Deep Learning

Deep learning is an approach of representation learning (a.k.a. feature learning), which is essentially a method to learn from data. Numerous layers of computational units together with appropriate activating mechanism comprise the basic architecture for deep learning. Multitudinous data sets are needed for those computational architectures to learn data abstractions for tasks such as image classification, speech recognition, object detection, etc. Each layer learns a level of abstraction from the data sets using back-propagation algorithm [22]. Making use of those learned abstractions, the computational architectures are able to solve complex problems which are typically non-linear and normally hard to solve by using specific rules that are designed in advance.

Deep learning has been attracting wide attention from all over the world in recent years, not only because of the great achievements it has made in various application scenarios, but also due to the promise of an intelligent future it gives. Such a learning methodology makes people believe it is possible for the formation of wise machines that they have long dreamed to possess. The growing data accessibility provides rich examples for deep computational architectures to adjust their internal weights and bias until their predictions have low error rate. On the other hand, the computational devices are relatively affordable than in the previous years by the society, with the help of which, accelerations of learning processes has been achieved, hence a bunch of time consuming deep learning

architectures can be tried within acceptable periods.

In the ILSVRC-2012 competition [23], AlexNet [24] received the championship with the 15.3% top-5 error rate, compared to 26.2% achieved by the runner-up. Such a large margin of error rate claimed a breakthrough in object recognition history. It kicked off a blistering pace of trying out deep learning by both academia and industry, which in turn led to an increase of the convolutional neural networks' submissions to ILSVRC-2013, in which ZF Net [25] was the winner. It fine-tuned the architecture of AlexNet based on the gorgeous visualizations of trained models. Both AlexNet and ZF Net are of the same structure which is built up by simply stacking computational layers while GoogLeNet [26] is composed of Inception modules. This new architecture was the most successful candidate in ILSVRC-2014. It has not only set the new height of object recognition but also started to optimize the computational resources of the network by design. It consists of 22 layers, which was deeper than all the previous networks in ILSVRC. However, it is still not deep enough. In ILSVRC-2015, Residual Neural Network (ResNet) [17] with 152 layers won the championships in all the five main tracks. ResNet introduced a brand new notion into the neural network architecture named identity mapping. The shortcut connection in the identity mapping prevents the degradation of training accuracy when the network goes deeper. Besides, the converging speed of ResNet is faster than the network built up with Inception modules when both are of the similar size.

Despite the fact that neural networks built up from Inception modules converge slower than those built up from ResNet modules, it is still worth it for a brief review of the valuable insights residing in the Inception networks. A typical incarnation of the first generation of Inception networks is named GoogLeNet [26]. It was intricately carved with a responsibility to win computer vision tasks in ILSVRC-2014, on which it performed better than all the other deep neural network architectures. There exist philosophical reflections which are intend to serve as guidelines for the construction of Inception networks. Two major downsides of a enlarged neural network have been discussed in [26]. One is the higher chances of overfitting while the other is the strikingly increased requirements of computational resources with the enlarged network size. For handling those drawbacks, based on the new ideas which were introduced in [27] about how to construct the reasonable architecture of neural networks, new experiments orienting sparse network structure have been tried out. One year later after GoogLeNet hold the championship of ILSVRC-2014, a method named Batch Normalization [28] has been proposed by Google researchers to accelerate and ease the training of deep neural networks. The core idea behind Batch Normalization is to normalize the inputs to each layer for every batch of training data. More importantly, based on the observation that the normalization process essentially is matrix multiplications followed by adding biases, the Batch Normalization is implemented as additional layers which makes it part of the network architecture. This fairly novel method started a new chapter for the training of deep neural networks. With the adoption of Batch Normalization, higher learning rates no longer impede the convergence of the deep networks, oppositely faster training speed is brought to scene which can achieve a better accuracy of prediction with considerably less

time. Additionally, in some cases, it can even replace the Dropout [29] which is an effective method to prevent overfitting. The incorporation of Batch Normalization into the first generation of Inception network architecture led to the formation of Inception-v2, which improved the best accuracy on ImageNet classification with less training steps. In the same year, Inception-v3 [30] joined the show, the objective of which was to effectively leverage the power of additional computation by factorizing to smaller size convolutions and regularizing the classifier layer with the estimation of minor effect of label-dropout in the training process. The network architectures were scaled up in Inception-v3, which consequently imposed higher requirements of available computational resources.

Chapter 3

Prepare the Data for Deep Learning

With the rising popularity of the high definition videos, the new standard termed High Efficiency Video Coding (HEVC) for compressing videos in a more efficient way comparing with previous standards, such as H.264/AVC, has emerged under the efforts from the Joint Collaborative Team on Video Coding (JCT-VC). In the meanwhile, five extensions of the HEVC standard, comprising Format Range Extension (RExt), Scalability Extension (SHVC), Multi-view Extension (MV-HEVC), 3D Extension (3D-HEVC), Screen Content Coding Extension (SCC), have been finalized from 2014 to 2016 to support fulfill extra requirements in various scenarios.

3.1 Video Coding

Chapter 4

Train the Deep Model for Prediction

With the rising popularity of the high definition videos, the new standard termed High Efficiency Video Coding (HEVC) for compressing videos in a more efficient way comparing with previous standards, such as H.264/AVC, has emerged under the efforts from the Joint Collaborative Team on Video Coding (JCT-VC). In the meanwhile, five extensions of the HEVC standard, comprising Format Range Extension (RExt), Scalability Extension (SHVC), Multi-view Extension (MV-HEVC), 3D Extension (3D-HEVC), Screen Content Coding Extension (SCC), have been finalized from 2014 to 2016 to support fulfill extra requirements in various scenarios.

4.1 Video Cnbnbnbnboding

Chapter 5

Evaluate the Learned Deep Model

With the rising popularity of the high definition videos, the new standard termed High Efficiency Video Coding (HEVC) for compressing videos in a more efficient way comparing with previous standards, such as H.264/AVC, has emerged under the efforts from the Joint Collaborative Team on Video Coding (JCT-VC). In the meanwhile, five extensions of the HEVC standard, comprising Format Range Extension (RExt), Scalability Extension (SHVC), Multi-view Extension (MV-HEVC), 3D Extension (3D-HEVC), Screen Content Coding Extension (SCC), have been finalized from 2014 to 2016 to support fulfill extra requirements in various scenarios.

5.1 Video Codhgvingg

Chapter 6

Employ the Learned Deep Model

With the rising popularity of the high definition videos, the new standard termed High Efficiency Video Coding (HEVC) for compressing videos in a more efficient way comparing with previous standards, such as H.264/AVC, has emerged under the efforts from the Joint Collaborative Team on Video Coding (JCT-VC). In the meanwhile, five extensions of the HEVC standard, comprising Format Range Extension (RExt), Scalability Extension (SHVC), Multi-view Extension (MV-HEVC), 3D Extension (3D-HEVC), Screen Content Coding Extension (SCC), have been finalized from 2014 to 2016 to support fulfill extra requirements in various scenarios.

6.1 Video Cobnbvnbvnbvnding

Chapter 7

Conclusion

With the rising popularity of the high definition videos, the new standard termed High Efficiency Video Coding (HEVC) for compressing videos in a more efficient way comparing with previous standards, such as H.264/AVC, has emerged under the efforts from the Joint Collaborative Team on Video Coding (JCT-VC). In the meanwhile, five extensions of the HEVC standard, comprising Format Range Extension (RExt), Scalability Extension (SHVC), Multi-view Extension (MV-HEVC), 3D Extension (3D-HEVC), Screen Content Coding Extension (SCC), have been finalized from 2014 to 2016 to support fulfill extra requirements in various scenarios.

7.1 Video Codvbnvnbvcnbvnbng

Bibliography

- [1] Web Page, 2017. [Online]. Available: <http://hidefnj.com/video.html>.
- [2] C. E. Shannon, *The mathematical theory of communication*. Urbana: Urbana : University of Illinois Press, 1949.
- [3] “Itu-t recommendation database,” 2017. [Online]. Available: <http://www.itu.int/ITU-T/recommendations/rec.aspx?rec=12905&lang=en>.
- [4] “Jct-vc - joint collaborative team on video coding,” 2017. [Online]. Available: <http://www.itu.int/en/ITU-T/studygroups/2013-2016/16/Pages/video/jctvc.aspx>.
- [5] I. E. Richardson, *The H.264 Advanced Video Compression Standard*, 2nd ed., ser. H.264 Advanced Video Compression Standard 2e. Hoboken: Hoboken : Wiley, 2010.
- [6] A. Vetro, T. Wiegand, and G. J. Sullivan, “Overview of the stereo and multiview video coding extensions of the h.264/mpeg-4 avc standard,” *Proceedings of the IEEE*, vol. 99, no. 4, pp. 626–642, 2011, ISSN: 0018-9219. DOI: 10.1109/JPROC.2010.2098830.
- [7] J. Konrad and M. Halle, “3-d displays and signal processing,” *IEEE Signal Processing Magazine*, vol. 24, no. 6, pp. 97–111, 2007, ISSN: 1053-5888. DOI: 10.1109/msp.2007.905706.
- [8] I. Sexton and P. Surman, “Stereoscopic and autostereoscopic display systems,” *Signal Processing Magazine, IEEE*, vol. 16, no. 3, pp. 85–99, 1999, ISSN: 1053-5888. DOI: 10.1109/79.768575.
- [9] Mu, amp, X, K. Ller, P. Merkle, and T. Wiegand, “3-d video representation using depth maps,” *Proceedings of the IEEE*, vol. 99, no. 4, pp. 643–656, 2011, ISSN: 0018-9219. DOI: 10.1109/JPROC.2010.2091090.
- [10] “Jct-3v - joint collaborative team on 3d video coding extension development,” 2017. [Online]. Available: <http://www.itu.int/en/ITU-T/studygroups/2013-2016/16/Pages/video/jct3v.aspx>.
- [11] G. Tech, K. Ying Chen, J.-R. Muller, A. Ohm, A. Vetro, and A. Ye-Kui Wang, “Overview of the multiview and 3d extensions of high efficiency video coding,” *Circuits and Systems for Video Technology, IEEE Transactions on*, vol. 26, no. 1, pp. 35–49, 2016, ISSN: 1051-8215. DOI: 10.1109/TCSVT.2015.2477935.

- [12] H.-B. Zhang, C.-H. Fu, Y.-L. Chan, S.-H. Tsang, and W.-C. Siu, "Probability-based depth intra mode skipping strategy and novel vso metric for dmm decision in 3d-hevc," *Circuits and Systems for Video Technology, IEEE Transactions on*, vol. PP, no. 99, pp. 1–1, 2016, ISSN: 1051-8215. DOI: 10.1109/TCSVT.2016.2612693.
- [13] H. Dou, Y.-L. Chan, K.-B. Jia, and W.-C. Siu, "Segment-based view synthesis optimization scheme in 3d-hevc," *Journal of Visual Communication and Image Representation*, vol. 42, pp. 104–111, 2017, ISSN: 1047-3203. DOI: 10.1016/j.jvcir.2016.11.012.
- [14] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*, ser. Adaptive computation and machine learning. Cambridge, Massachusetts: The MIT Press, 2016, xxii, 775 pages, ISBN: 9780262035613 0262035618.
- [15] B. Schölkopf, C. J. C. Burges, and A. J. Smola, *Advances in kernel methods : support vector learning*. Cambridge, Mass. ; London: Cambridge, Mass. ; London : MIT Press, 1999.
- [16] J. Schmidhuber, "Deep learning in neural networks: An overview," *Neural networks : the official journal of the International Neural Network Society*, vol. 61, p. 85, 2015.
- [17] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," 2015.
- [18] Web Page, 2017. [Online]. Available: <https://www.bazel.build/>.
- [19] Web Page, 2017. [Online]. Available: <https://cloud.google.com/storage/>.
- [20] C. Y. Chou, "Advances in grid and pervasive computing: First international conference, gpc 2006," in. 2006, ISBN: 3540338098. [Online]. Available: [https://en.wikipedia.org/wiki/Bandwidth_\(computing\)#Network_bandwidth_capacity](https://en.wikipedia.org/wiki/Bandwidth_(computing)#Network_bandwidth_capacity).
- [21] M. Ghanbari, *Video coding: an introduction to standard codecs*. London: Institution of Electrical Engineers, 1999.
- [22] Y. A. LeCun, L. Bottou, G. B. Orr, and K. R. Müller, "Efficient backprop," *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 7700, pp. 9–48, 2012, ISSN: 03029743. DOI: 10.1007/978-3-642-35289-8-3.
- [23] Web Page, 2017. [Online]. Available: <http://www.image-net.org/challenges/LSVRC/>.
- [24] A. Krizhevsky, I. Sutskever, and G. Hinton, "Imagenet classification with deep convolutional neural networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017, ISSN: 0001-0782. DOI: 10.1145/3065386.
- [25] Generic, 2014. DOI: 10.1007/978-3-319-10590-1_53.
- [26] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," 2014.

- [27] S. Arora, A. Bhaskara, R. Ge, and T. Ma, “Provable bounds for learning some deep representations,” *CoRR*, vol. abs/1310.6343, 2013. [Online]. Available: <http://arxiv.org/abs/1310.6343>.
- [28] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” 2015.
- [29] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: A simple way to prevent neural networks from overfitting,” *J. Mach. Learn. Res.*, vol. 15, pp. 1929–1958, 2014, ISSN: 1532-4435.
- [30] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” 2015.