



Faculty of Electrical Engineering and Information Technology
Institute for Media Technology
Audio Visual Technology

MASTER THESIS

Netflix Like Encoding Optimization

Submitted by: Vijaykumar Singh Rana

Major: Media Technology

Advisor: Prof. Dr.-Ing. Alexander Raake

Co-Advisor: M.Sc. Steve Göring

Ilmenau, October 29, 2019

Acknowledgments

optional

Zusammenfassung

maximum of 2400 chars; one paragraph

Abstract

maximum of 2400 chars; one paragraph

Contents

1	Introduction	1
1.1	Motivation & Goals	5
1.2	Thesis Structure	7
2	Fundamentals	9
2.1	Related Work	9
2.2	Two-Pass Encoding	12
2.3	FFMPEG	13
2.4	Video Quality Metric: Video Multi-Method Assessment Fusion (VMAF)	14
3	Design & Implementation	17
3.1	Theoretical System Design	17
3.1.1	Extract Scenes/Shots	18
3.1.2	Generate Resolutions	18
3.1.3	Generate Encodes	18
3.1.4	Video Quality Score Estimation	18
3.1.5	Analysis and Evaluation	19
3.2	Initial Pipeline Design Setup	21
3.3	Pipeline based on Bisection Method Approximation	28
3.3.1	Theoretical Approach with Bisection	30
3.3.2	Generate Resolutions	30
3.3.3	Generate Encodes and Estimate Quality	30
3.3.4	Compute Next bitrate point	31
3.3.5	Quality and Bitrate Threshold	31
3.4	Pipeline based on Bisection Method Approximation with Bitrate Threshold	35
3.4.1	Generate Resolutions	37
3.4.2	Generate Encodes & Quality Estimation	37
3.4.3	Compute Bitrate Points	38
3.4.4	Quality Threshold/Bitrate Threshold	38
3.4.5	Live Plotter	39
4	Analysis & Evaluation	43
5	Conclusion	45

Contents

Bibliography	47
List of Figures	51
List of Tables	53
A AppendixExample	

Chapter 1

Introduction

In today's world, media is being offered to the users on a large scale. That means users around the world are indulged in the streaming services for their favourite shows. These shows are delivered to them based on the available bandwidth. From the production point of view, it must be encoded using schemes that would allow the end-users to have a good quality of experience by the content of the show. YouTube[Yt], Netflix[Neta], Amazon Prime Video[Pri] etc. are the most popular video hosting websites on internet that provides VOD (video on demand). They have evolved their schemes for transcoding videos in a way such that the end-user would receive reasonable quality of videos corresponding to the bandwidth available to them. The service providers that can deliver these video streaming services have to use various schemes for transcoding videos based on the subjective quality for end-users to have a good quality of experience. These video stream providers mainly aim to reduce the bandwidth, consequently reducing the costs and to satisfy the end user so that they can perceive the best possible quality. Dynamic Adaptive Streaming over HTTP (DASH) is one of such schemes that provides streams based on the bandwidth. As the name suggests, it uses a streaming technique that provide media content with quality that can be delivered with available bandwidth.

Dynamic Adaptive Streaming over HTTP (DASH) [ISO14], also known as MPEG-DASH, is an adaptive bitrate streaming technique that enables high quality streaming of media content over the Internet delivered from conventional HTTP web servers. The content is made available at a variety of different bit rates, i.e., alternative segments encoded at different bit rates covering aligned short intervals of playback time. While the content is being played back by an MPEG-DASH client, the client automatically selects from the alternatives the next segment to download and play based on current network conditions. The client automatically selects the segment with the highest bit rate possible that can be downloaded in time for playback without causing stalls or re-buffering events in the playback. Considering video streaming the main goal of a video stream provider is to reduce bandwidth consequently the cost and to satisfy the end user, so that they are able to perceive the best possible quality [Sto11].

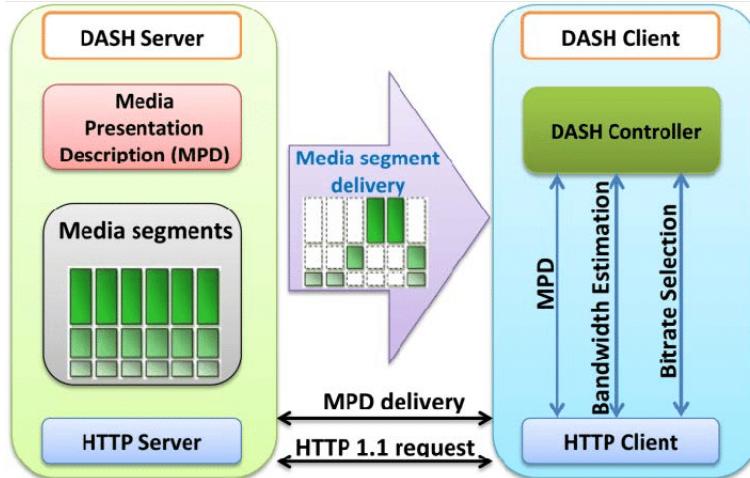


Figure 1.1: The overview of Dynamic Adaptive Streaming over HTTP technique. [VMC17]

Currently, YouTube and Netflix are some of the most popular video hosting websites in the world [Moo+19]. They provide the end-users with streaming service which is affordable and provide the best possible quality content based on the fluctuating network conditions. To provide the high quality content, YouTube uses chunked video playback. The video is split into smaller parts called chunk. Each of the chunks are encoded from low quality to high quality. The suitable quality chunk is selected based on network bandwidth. Because of chunking, the web browsers are able to manage the smaller pieces of the videos, in their dedicated cache for playback. For fluctuating network speeds and optimal user viewing experience, YouTube uses Adaptive Bitrate Streaming (ABS). Figure 1.2 demonstrates how a video stream is selected based on the network condition. When the video starts on YouTube it starts with low quality and then based on the network condition the quality keeps switching. As it can be seen in the Figure 1.2 when the available bandwidth gets higher, the high bitrate stream is then selected. in the next timestamp there appears a network congestion and eventually the quality drops i.e., high bit rate stream is not streamed and based on lower bandwidth low bit rate stream is then selected. Further as the network bandwidth changes, the suitable stream is then selected. Beside network conditions YouTube also acquires screen dimensions and hardware ability of the client device, where in addition an adaption is performed, e.g. a 4K stream would not be played on a Full-HD smartphone's display.

Netflix also uses the similar adaptive streaming approach to provide streaming service to the users. Netflix also uses granular approach of chunking the video stream based on the shots. A sequence of frames shot by a camera until a point is called a shot. A video sequence is sequence of shots alternated based on the actors in the video sequence. A scene on the other hand is a collection of shots focused on an actor(s) or or an object(s)[BR96]. In this thesis, scene and shot are referred to be same.

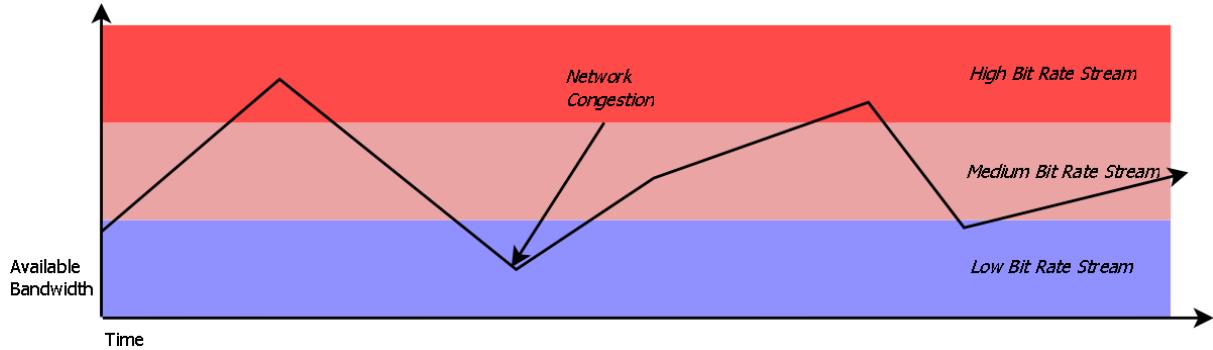


Figure 1.2: Adaptive Streaming [OJ]

because in this thesis a scene cut detection is performed, thus individual shots may not be correctly identified. Besides that, a detected scene segmented from a given video sequence has shots with similar attributes as described in [BR96] i.e., a scene is collection of shots with a focus on same spot, object(s) and actor(s). Therefore, the scene cuts are "better" than only shots, because some shots could be just similar. Netflix uses optimized encoding schemes for encoding a video i.e., Netflix analyzes the scenes for a video sequence individually and optimizes each of them to produce the best quality of content for given set of bitrates. On the top level Netflix uses Per-Title Encode Optimization where a scheme to encode the video optimally is prepared. For every individual video sequence the scheme is generated because the content of the video sequences varies and hence Netflix analyzes and then encodes the video individually. The objective is to deliver better quality of experience for users with bandwidth available to them. The reason for using per title encoding is because the approach of "one-size-fits-all" fixed bitrate ladder is not a good solution [Año+15]. In many cases the video content differs from video to video. So setting a fixed bitrate for certain resolution is not a good solution and is inefficient. Figure ?? show a table with pairs for bitrate and resolution. Now based on this table 5800 kbps should be good enough for generating a stream of FULL HD resolution(1080p). However, certain scenes might contain fluctuating luminance for a region in the frames which leads to camera noise. The camera noise distorts the quality of the the scene that results into blockiness artifacts in the noisy region of the frames. On the other hand for generating a Full HD stream for a video containing cartoon content 1750 Kbps is more than enough, so utilizing 5800 Kbps is not an optimal solution and would lead to bandwidth wastage [Aar+15].

In case of Netflix, the videos are not just encoded using conventional constant bitrate video encoders. Netflix encodes the videos by using multiple coded representations (with multiple bitrates) of a scene for different resolutions and the encoded segments are delivered, based on the request to the end-user[Kat18]. It allows the end-users to stream high quality video with low bitrate and provides the end-users seamless playback of the content. The reason behind achieving this is subjective. From an end-user's perspective, who can afford good network

Bitrate (kbps)	Resolution
235	320x240
375	384x288
560	512x384
750	512x384
1050	640x480
1750	720x480
2350	1280x720
3000	1280x720
4300	1920x1080
5800	1920x1080

Table 1.1: Bitrate-Resolution pairs or Bitrate Ladder [Aar+15]

bandwidth, would not want to have artifacts like blurriness or blockiness in the video content. This requires the development of algorithms that could deal with the issues of artifacts in videos encoded with low bitrates. The content of a given video differs with other videos more or less, as it can have scenes with more of constant background covering large area of frames or it can be a dynamic scenes based on lot of motion involved i.e., rapid change in the frame content over time or it can have fluctuating intensity values for pixels over the frames in the scene. When using constant bitrate technique, the dynamic scenes may utilize the bits efficiently but constant scenes where the region of the frame is not changing too much, the bits allocated for those frames would be more than required, thus leading to over utilization of bits .The worst case scenario would be a scene where it is completely dark which does not require as many number of bits as a scene with frames containing varying pixel intensity. The figure 1.3 demonstrates how the quality varies for different video sources. Certain sources never reach high quality even though the bitrate is high enough as 25000 kbps. The reason for high bitrate requirement arises due to the complexity of content in a given video sequence. IF the given video sequence has high varying temporal information then this leads to coding of more number of frames i.e., the information from both intra-frame (spatial feature) and inter-frame (temporal feature) keeps changing after every few frames. The video compression tends to exploit the redundant information in intra-frame and inter-frame to achieve maximum video compression [GSW00]. Thus high change in the content within fewer frames leads to low compression of the video sequence.

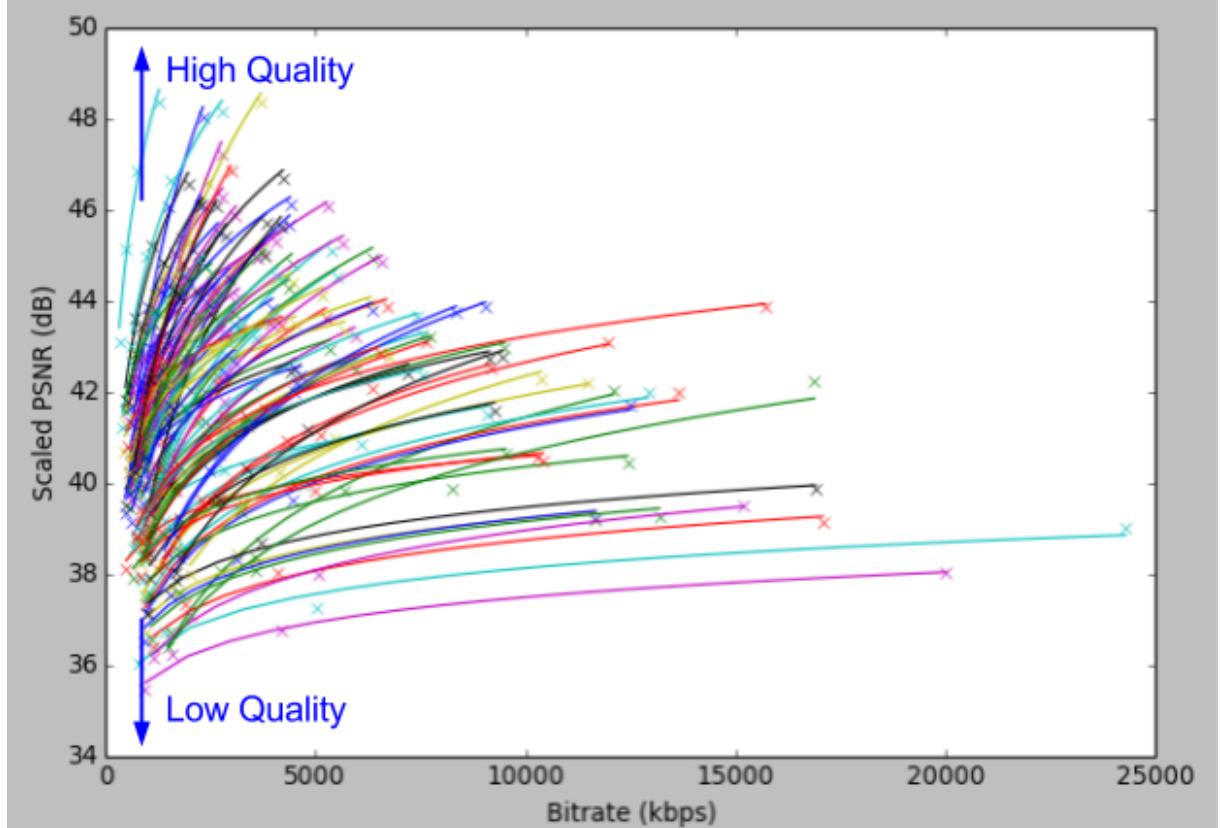


Figure 1.3: 100 randomly sampled sources at 1080p resolution using x264 constant QP (Quantization Parameter) rate control. [Aar+15]

Also the process of achieving best quality goes further into encoding each of chunks based on the scenes. Every video sequence can be subdivided onto multiple scenes. As each of the scenes have different characteristics, so the idea is to encode them individually with required bitrates for different resolutions based on the subjective quality.

Objective video quality metrics like PSNR (Peak Signal to Noise Ratio) [Psn], SSIM (Structural Similarity Metric)[Wan+04], VMAF (Video Multi-Method Assessment Fusion) [Blo18] are used for video quality evaluation. The videos encoded for each quality with all possible bitrates are evaluated. Finally those bitrates are chosen which provides high quality amongst the different resolutions on a rate-distortion or rate-quality curve.

1.1 Motivation & Goals

The main goal of this thesis is to develop a pipeline for video transcoding using an approximation method as an alternative to the Netflix's brute-force based pipeline technique. The objective is to achieve the reduction in the computation cost for transcoding a given video sequence only for relevant bitrate points for respective resolutions(6 resolutions ranging from

Chapter 1 Introduction

360p to 2160p). Section 1 describes about how Netflix approaches to encoding a video by subdividing it into segments and then encode them individually to provide a good quality of experience from subjective point of view. This pipeline includes segmenting a video as a first step. Next each of these segments are assessed and divided into scenes. These scenes are based on the type of content and have relatively similar kind of frames which makes it easier to encode them with low bitrate compared to encoding a video with more varying content. For this Netflix applies encoding with all possible bitrates for each resolution [Blo17]. The encodes generated are then scored using a objective quality metric method called Video Multi-Method Assessment (VMAF) from Netflix. Then based on the quality scores obtained from a rate distortion plot, the appropriate bitrate is chosen for each resolution. Netflix applies a brute force technique in order to compute quality score for encodes transcoded at each bitrate for all the resolutions. This requires intensive computational resources. The question arises that if there is any method to reduce the overall computation costs and approximate this whole process using optimization techniques. If the number of bitrate points can be reduced to only select number of bitrate points and still be able to obtain the complete rate-quality curve by approximation, the overall computation cost will be reduced and as a result it will speed up the whole process chain.

To have such comparison both of the approaches will be implemented; a) brute force technique of computing quality scores by computing quality for each bitrate for different resolutions to generate a rate-quality curve and b) selecting fewer points from the generated rate-quality curve and then trying to approximate using mathematical equations e.g., a logarithmic curve.

The motivation here is to do the computation for less number of bitrates and try to obtain other points by the means of any approximation method. This will include interpolation and extrapolation of the remaining bitrate points to generate a rate-quality curve that would resemble closely to the rate-quality curve generated using brute-force technique. Then the overall number of encoding steps would be reduced. The next step would be to compare the resulting points after approximation with points from the brute-force approach. Further analysis would provide the error between points obtained using approximation approach and points from the brute-force approach. Based on the error the plausibility of the approach to minimize the overall encoding steps can be analyzed. For a given video sequence the resolutions of 360p, 480p, 720p, 1080p, 1440p and 2160p will be considered and corresponding to each resolution, the computation of the quality scores for all possible coded representation using various bitrates will be done.

1.2 Thesis Structure

In order to provide the information regarding the work done and the contributions in this master thesis, the documentation is structured in the following manner:

- ▷ Chapter 1 addresses about current overview of encoding pipeline used by Netflix and the challenges about setting it up. In addition to that how the challenging aspects of Netflix pipeline could be replaced using approximation methods in order to achieve reasonably similar results which is the main motivation of this thesis topic.
- ▷ Chapter 2 describes different principles of encoding schemes, video quality metrics and tools already available to set up the pipeline for different approaches for transcoding video sequences. Additionally it includes theories about previous work done with regards to video encoding optimization.
- ▷ Chapter 3 focuses on the various encoding pipelines used to generate the rate-quality curves for the video sequences by transcoding them with multiple bitrate points for various resolutions. The approaches being the pipeline for brute-force technique used by Netflix, Bisection approximation technique and finally the approximation technique based on the mathematical equation of logarithmic curve. Apart from that it also addresses the challenges that were encountered in each of the approaches including the pros and cons with respect to the approaches.
- ▷ Chapter 4 is about analysis and evaluation of the results obtained after the implementation of the approximation method of mathematical equation, where the resulting curves describes how plausible is the approximation approach of mathematical equation with respect to brute-force technique i.e., the curve generated for a corresponding resolution will be evaluated using mean squared error to find the plausibility of the approximation approach.
- ▷ Finally Chapter 5 concludes about the approximation approach and comments about overall work done and finally providing some future work recommendations.

Chapter 2

Fundamentals

In this chapter, the first section 2.1 focuses on the previous work done related to this thesis. Following the related work, this section will continue to describe the concepts and methods for various tools used in the implementation phase. The description of the tool would also describe how it is essential to this thesis.

2.1 Related Work

As described in the previous chapter 1 about Netflix encoding videos, firstly using the different encoding strategy for each video and then further using granular approach to optimize each of the scenes for a the given video. From the point of view of adaptive streaming, the problem was not about encoding a single scene for each video but rather how to do it for collection of all the scenes with all the possible bitrate/qualities range that the user would desire the service to receive. This led to development of a framework which considered all the above points from Netflix Technology called Dynamic Optimizer [Kat18]. This framework analyzes an entire video over multiple quality and resolution pairs, finds the best trade off to generate optimal scheme to encode a scene[Man+18].

The idea behind finding the optimal scheme to encode a video sequence is that a video consists mainly of numerous shots, so each shot should be encoded independent of other shots. These shots have little correlation with each other which can be treated an element that can be encoded independently of each other. Each of the shots to be encoded can be considered as visual content which is uniform and homogeneous in sense of the object or the actor the scene is focused on. This means applying constant or fixed quality for an individual shot fits well. From the perception point of view, viewers are not bothered by the coding difference that would occur at the boundaries of adjacent shots i.e., whenever there is a scene or a shot change it does not affect the subjective quality of experience. It provides a good deal to adjust the coding resolution and quantization parameter at such points of shot boundaries. The decision

Chapter 2 Fundamentals

of choosing optimal resolution and quantization parameter to be used for each of the shots to be encoded, is done on the basis of perceptual video quality metrics like SSIM (Structural Similarity) and VMAF (Video Multimethod Assessment Fusion). Dynamic Optimizer uses cloud computing to generate multiple encodings for a given set of frames from a shot for various resolutions and quantization parameters[KG18]. In this thesis topic the idea is similar to setup a pipeline for video encoding. So in the similar way the videos will be split into shots. Each of the shots will be reproduced for various resolutions. Each of those resolutions of the scene will be coded with multiple bitrate values producing encodes for each bitrate value. Each of the encodes will be then tested with VMAF tool for VMAF score (quality score) for each bitrate for a given resolution. For each bitrate and its corresponding VMAF score, a plot will be generated. The plot would be similar to the Figure 2.1, i.e., rate-quality curve with bitrate versus corresponding scaled PSNR values. The labels (A, B, C and D) in the plots represents each of the curves corresponding to different resolution. The curve **A** (in black) is the plot generated using rate-quality pairs of points (R_i, Q_i - set of bitrate points and its corresponding quality scores for a given resolution) for various bitrates for the resolution 720x480. Similarly curve **B** (in green) is generated using rate-quality pair of points, for the resolution of 1280x720 and finally the curve **C** is generated with rate-quality pair of points for 1920x720. The red curve **D** is the convex hull obtained over all the curves from the resolutions and its R_i, Q_i points. This curve is used as the basis to select the points for a given bitrate. Although the quality parameters is generated using traditional metric of scaled PSNR in the figure, however in this thesis the rate-quality curve will include quality scores generated using VMAF tool. As discussed in the previous section, the next step will include approximation technique where a fewer points for a curve will be selected and then the rest of the points will be approximated.

The transcoding procedure used by Netflix requires to encode all video sequences with all possible bitrate settings. Netflix uses transcoding pipeline in the cloud based services to transcode the videos and to speed up the processing. Transcoding each video with unique set of settings where each of its segments are transcoded with various bitrates settings requires large amount of computational and storage resources. The authors from [KZS15] discussed about how large number of transcoded video segments are never used which leads to wasteful management of storage and computation. They have described about a possible solution that when a user requests a segment of a video, the bitrates with which the segment is transcoded can be predicted using a Markov model predictor[OR02]. In addtion to that they have also described that based on the Netflix study of codec and bit-rate combinations for a single segment can result up to 120 transcoding operations. For a given video of a regular television show of about one hour, assuming having scenes on an average duration of 10 seconds will have 360 video segments. Applying 120 transcoding operations for each of those scene or a video segment is not feasible approach. They have shown that based on the the next

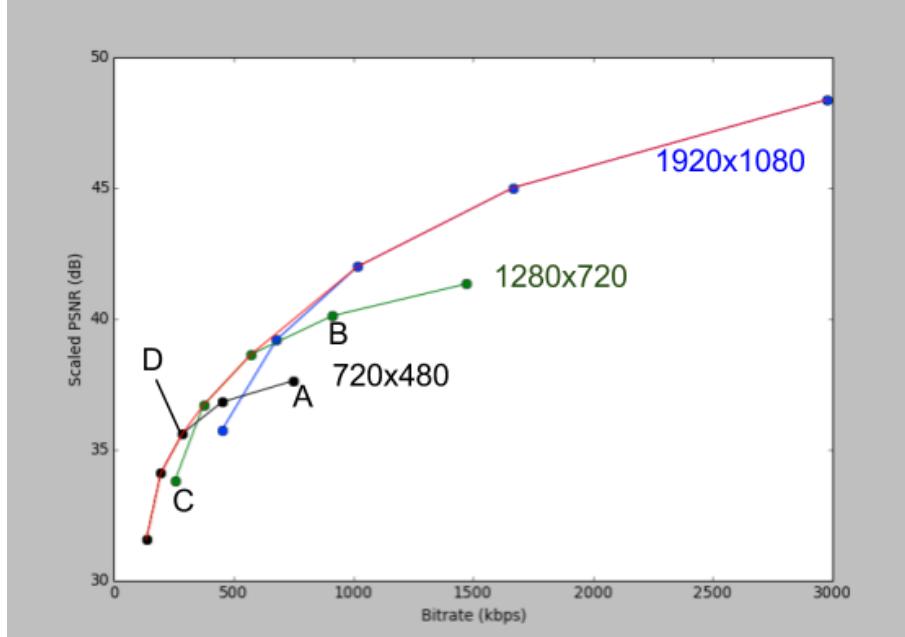


Figure 2.1: Example encodes showing individual R-D curves and convex hull.[Aar+15]

video segment requested by the user, the corresponding bitrate with which the segment is to be encoded, can be predicted. The main objective of the authors being to transcode video segment with only select settings of bitrate and quality. In this thesis the motivation is quite similar where instead of transcoding a given video segment with multiple bitrates the idea is to transcode for fewer bitrate points for a given resolution and then try to approximate remaining bitrate points based on the mathematical approximation of curve generated using the few select bitrate points. Therefore reducing the computational complexity of trancoding pipeline.

The authors from [Sat+19] presented an alternative to per-scene video encoding done by Netflix. They used the CRF (Constant rate factor) encoding scheme for 3 resolutions and for each of those resolutions they computed 4 different (R_i, Q_i) points. For generating quality points for each of the encodes, they used a full reference metric named PEVQ (Perceptual Evaluation of Video Quality) from OPTICOM [Pev]. The authors explained about how the function of (R_i, Q_i) can be used to model the prediction of PEVQ quality scores by using the combination of logistic and inverse exponential functions. Using this approach they try to compute all the quality points for all bitrates for a given resolution and finally getting all computations for a given video. Their approach also takes into consideration the tradeoffs between bitrate and quality. In this thesis topic the approach is quite similar. The mathematical approximation would be done using a logarithmic curve instead which is same as an inverse exponential function. The logarithmic curve will be fitted with respect to fewer points selected from (R_i, Q_i) pairs for every given resolution.

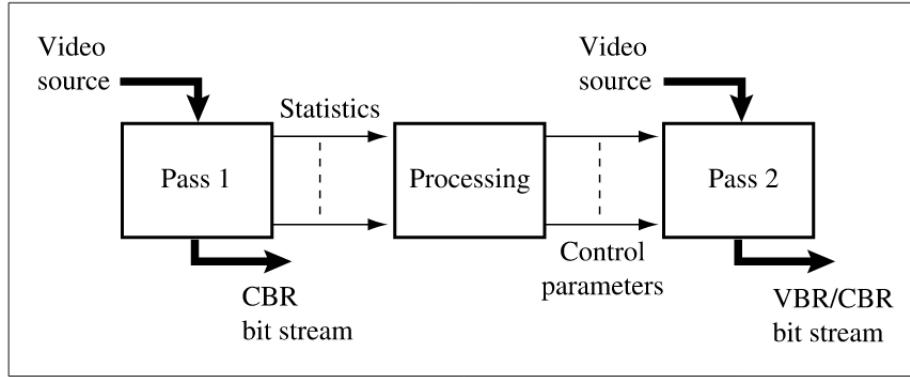


Figure 2.2: Schematic for Two-Pass Video Encoding System.[WRG99]

2.2 Two-Pass Encoding

The pipeline used for the brute force technique by Netflix requires to transcode video at variable bitrates for every resolution of a shot (or a scene) from a given video. The transcoding scheme that is used in this thesis topic is Two-Pass Video Encoding. The two-pass encoding scheme is a for non-real time applications that enables processing a video in two passes. The first pass usually analyzes the given video sequence and generates the information regarding statistics about each element of the video. The results obtained from the first pass is then used to optimally encode the video better in the second pass. This kind of video encoding uses variable bitrate (VBR) to encode the different elements of the video unlike the CBR encoding where the whole video is encoded and compressed using a constant bitrate (CBR). In the CBR encoding the rate remains constant but the quality of the the video sequence keeps varying based on the varying content[MRG99]. To have a constant quality video sequence it has to be code using variable bitrate which can be achieved using two-pass encoding. During the first pass the video sequence is encoded using the conventional CBR and at the same time the statistics data about the coding the elements of the video is are obtained. The first pass statistics are processed to produce the control parameters for the second pass. The second pass is then used to generate VBR compressed stream of the the given video sequence using the first pass statistics data. The control parameters which are generated between the first pass and the second pass are used to distribute the bits appropriately to the different elements of the video such that at the end the video is encoded with variable rate and a constant quality.[WRG99]

Figure 2.2 shows the overview of a two pas encoding system. As shown in the figure, the first pass is responsible for applying constant rate CBR encoding to the given video sequence. The statistics obtained from this pass are then processed after the video sequence. The processing of the statistics is done after gathering the data from the complete video sequence in the first

Command Tool	Usage
ffmpeg	is used for converting audio or video formats
ffplay	is a command line media player that uses Simple Directmedia Layer(SDL) and FFMPEG libraries
ffprobe	displays text based media information in formats like XML, CSV etc

Table 2.1: FFMPEG Command Line Tools

pass. So the processing for statistical data cannot be done before. Based on the statistical data processed after the first pass, the control parameters are computed. In the second the video is encoded using the computed control parameters for each of the elemental parts of the video. Thus the final video encode is optimized CBR sequence or VBR stream. [WRG99]

So to transcode a given video sequence requires a tool called FFMPEG [Ffmb]. The following section describes the overview about the FFMPEG tool.

2.3 FFMPEG

FFMPEG is an open source cross platform tool which allows for handling audio and video data by the means on encoding, decoding, transcoding, multiplexing, demultiplexing etc. This tool allows for many options for customizing video data by the means of its capabilities[Ffmb].

FFMPEG provides various command line tools, which is described in the table 2.1.

FFMPEG also provides libraries for various applications described in the table 2.2.

Libraries	Usage
libswresample	provides functions to resample audio
libavresample	provides functions to resample audio from Libav project [tea]
libavcodec	provides native FFmpeg audio/video encoders and decoders
libavformat	provides demuxers and muxers for audio/video container formats
libavutil	includes hash functions like SHA-1, LZO decompressor and Base64 encoder/decoder.
libpostproc	provides functions for old H.263 video postprocessing
libswscale	provides functions for video image scaling and colorspace/pixelformat conversion
libavfilter	is the substitute for vhook which allows the video/audio to be modified or examined between the decoder and the encoder

Table 2.2: FFMPEG Libraries

After the video is processed and encoded, the quality estimation has to be done. There are numerous schemes and metrics available to compute the quality of a video based on the quality of a video. The focus of this thesis for the video quality estimation is based on a full-reference metric introduced by Netflix and later made open-source framework. The metric is called Video Multi-method Assessment Fusion or VMAF. The following section describes it in details.

is an objective full-reference video quality metric developed by Netflix in cooperation with the University of Southern California and the Laboratory for Image and Video Engineering (LIVE) at The University of Texas at Austin.

2.4 Video Quality Metric: Video Multi-Method Assessment Fusion (VMAF)

VMAF is developed by Netflix in collaboration with the researchers from University of South Carolina University of Texas at Austin. The tool uses a machine learning approach to fuse

2.4 Video Quality Metric: Video Multi-Method Assessment Fusion (VMAF)

together the elementary video quality metrics. The elementary metrics are the highly efficient objective video quality metrics based on the human perception. The first set of elementary metrics are chosen based on the image quality metric called Visual Information Fidelity (VIF). This metric provides five different quality scores for each resolution of a given image i.e., they rescale a given image 5 times and calculate VIF on all the scaled versions. The second set of elementary metrics are chosen based on quantification of structural loss which occurs due to blurriness and encoding compression called Detail Loss Measure (DLM) [LMN12]. VIF and DLM are both image based quality metric where it computes the quality based on the spatial features. Finally the feature based on temporal information is obtained based on the temporal change in the information i.e., between the consecutive frames, the difference in the pixel is calculated. For this, the luminance component for each of the corresponding pixel is taken into consideration. This is same as the TI (Temporal Indicator)[Win12].

The figure 2.3 shows VMAF system diagram. The authors from [KG18] stated about the studies describing how VMAF quality score comes really close to human perception for video quality in terms of adaptive streaming where artifacts are caused mainly by scaling and compression. Netflix has made the VMAF project as an open-source project which means researchers and developers can contribute to this project for further improvement and new use-cases[KG18]. As described in previous paragraph the elementary quality metrics which are computed and described using VIF and DLM. All these features are then fused together through a Support Vector Machine (SVM) regressor [CV95]. The regressor produces a video quality score in the range [1 to 100], where 0 means the lowest quality to 100 which means nearly perfect quality. The regressor uses the parameters based on the collection of subjective testing based on full reference technique i.e., the users rated the quality of reference and the distorted videos respectively. [KG18]

VMAF is a full reference quality metric. As compression and scaling several artifacts result in the video so VMAF uses the existing perceptual quality scores from multiple quality assessment algorithms to estimate the quality of a given video. As it is a full reference metric, the original source and encoded/compressed source both are taken in to consideration. For multiple encoded versions of a source video the comparison is done using a scaled VMAF score. [Ras17][KG18].

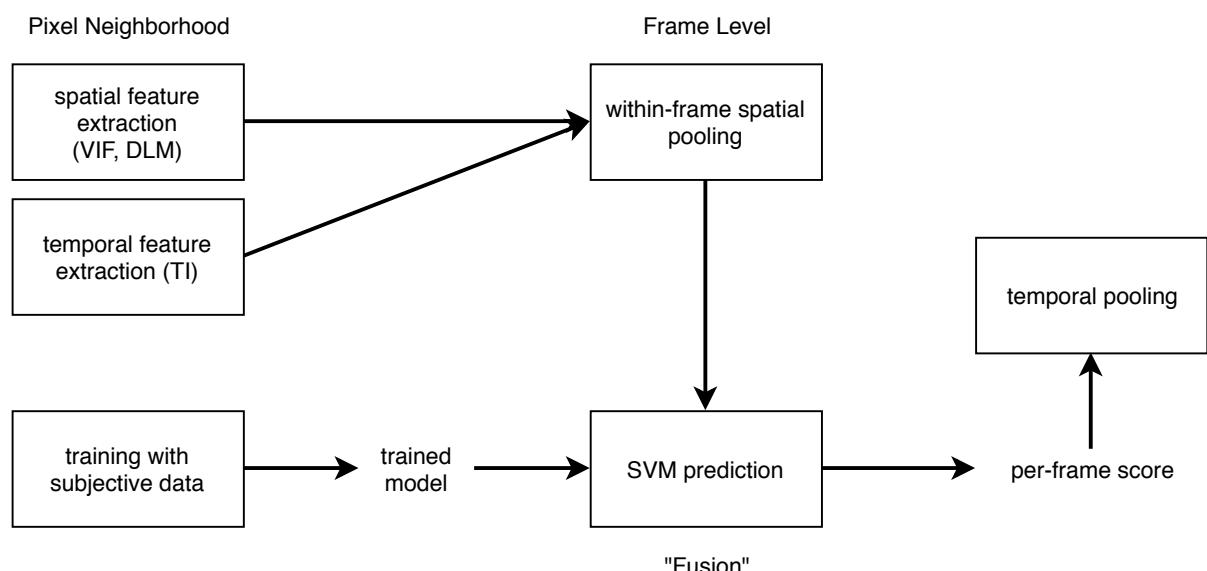


Figure 2.3: High Level VMAF System Diagram [KG18]

Chapter 3

Design & Implementation

3.1 Theoretical System Design

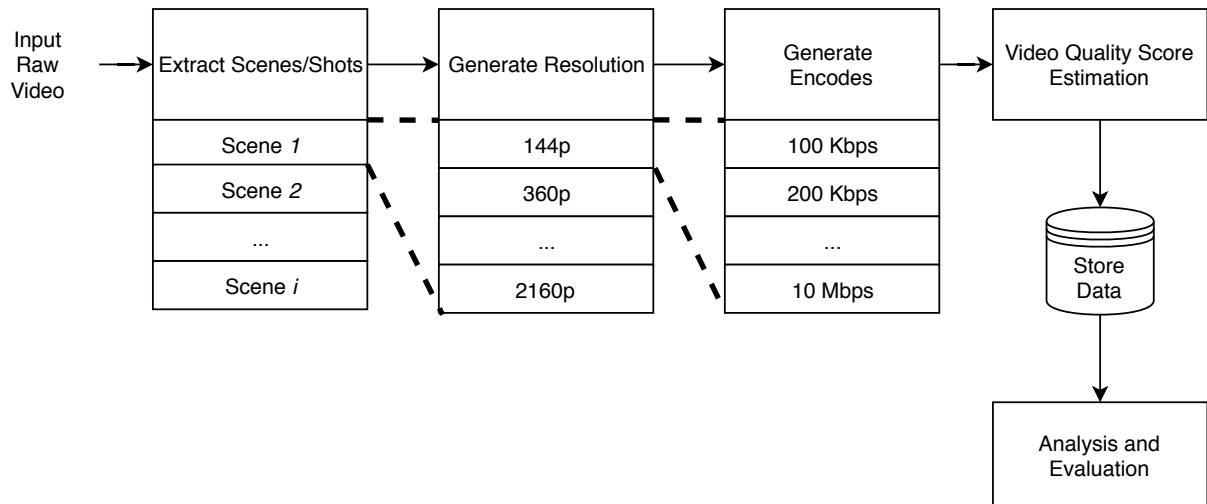


Figure 3.1: Overview of Theoretical Video Encoding Pipeline

The initial idea is to implement set of algorithms for various operations to transcode a given video sequence. Figure 3.1 provides the general overview of the encoding pipeline. As described in previous chapters about how Netflix uses various transcoding operations about how a given video sequence is transcoded with all possible bitrate settings for each resolutions. Similarly in the figure 3.1 each of the blocks represent the operations that takes place for transcoding a video.

The following subsections will describe the pipeline in detail and how each of the blocks are responsible and for what operations.

3.1.1 Extract Scenes/Shots

A raw video is the input to the pipeline and at first operation of extracting the scenes is performed. In this phase a given video sequence is analyzed to find all the possible scenes. The scenes are detected by the means of change in correlation of adjacent frames. Using a threshold value to find the difference in the correlation, the scene or shot boundaries can be identified. There are various approaches available which goes from simple fade in/fade out changes to content based detection to find the scene boundaries in a given video sequence.

3.1.2 Generate Resolutions

In the previous phase, a given video sequence is subdivided into multiple scenes or shots. Now in this phase for each of the scenes, different resolutions are generated. For every scene Scene_i , its different resolutions are generated. The resolution settings will range from the lowest resolution of 144p upto 1440p (QHD). Please note that for the highest quality of 2160p is not produced by down-scaling but instead the original scene will be taken into consideration.

3.1.3 Generate Encodes

This is in general overview about transcoding phase describing how each video is dealt with in the pipeline. After down-scaling the scenes from previous phase to various resolutions, each of those resolutions are then subjected to transcode with bitrate settings ranging from 100kbps to 10Mbps. For each of the resolutions, several bitrate transcoded encodes are generated. So to summarize, a given video sequence is subdivided into X number of chunks called the scenes or shots, then for each of the scenes, Y number of resolutions are generated and then finally for each of the resolution, Z number of encodes are generated. So, at the end, for a given video sequence there will be a total of $X \cdot Y \cdot Z$ number of encodes are generated.

3.1.4 Video Quality Score Estimation

After the encodes are generated, each of them have to be analyzed using a video quality metric. The metric to be used in this thesis is VMAF from Netflix which is a full reference metric. The original split scene will be used as reference for all its resolutions with multiple encode generated using various bitrate settings i.e., for each encoded version, an upscaling to the native resolution of the reference video is performed. For a source video scene with the resolution of 1080p used as reference for its encoded version with resolution 480p, the VMAF calculation will be done by upsampling the encode to 1080p to match the source video scenes

resolution [Blo18]. For each and every encode, a score value is generated. The resultant scores are aggregated based on encodes corresponding to each resolution. Finally the results are stored for further evaluation.

3.1.5 Analysis and Evaluation

The scores which were generated from the encodes are kept separately for each of the resolutions. This is done to obtain a rate-quality curve. Each rate-quality curve corresponds to each resolution which would describe the change in quality as the bitrate increases. It is anticipated that as the bitrate increases the quality score would initially increase by large steps of quality score and keep increasing until a point from which the quality score steps starts to decrease as the bitrate increases. The curve would resemble to the Figure 2.1 from section 2. In the same manner, the rate-quality curve will be generated for all the other resolutions.

In this thesis, the pipeline was assumed to take into consideration a single scene for each of the different video sequences. The video sequences were raw videos with resolution 2160p. However the framework is not limited to a single scene, it can be used for all the detected scenes. The main reason for processing only a single scene was to test the idea of shot based encodes used by Netflix [Man+18]. As each shot has its own background and foreground properties it can be encoded based on the complexity of its content optimally as described by Netflix Tech blog in [Man+18]. So the first objective was to extract scenes or shots as described in the figure 3.1.

As described before the scene detection can be done based on threshold between adjacent frames or content based difference in adjacent frames. Some of the tools available for detecting scenes in a given video sequence are as follows:

- ▷ FFmpeg blackframe filter [Ffma]
- ▷ Shotdetect [Can]
- ▷ Matlab Scene Change Detection [Mat]
- ▷ PySceneDetect [**pyscenedetect**]

FFmpeg blackframe filter is a video scene detection tool based on threshold technique where the intensities of adjacent frames are correlated to find the scene or shot boundaries. The drawback is that this tool is limited to threshold only. If there are consecutive scenes which are dark then this tool may fail as it was observed in pre-test that this tool does not always work properly.

On the other hand *Shotdetect* is another video scene detection tool which finds scene for a given video sequence based on content change only. In this thesis, a tool suitable for scene

Chapter 3 Design & Implementation

detection should be based on approaches of both threshold mode and content mode. This tool takes content mode into consideration only. This tool works on the basis of change in the motion in the consecutive frames. The change in attributes of RGB (Red, Green and Blue) and HSV (Hue, Saturation and Value) intensity is also taken into consideration for finding scene or shot boundaries [Joh].

Matlab Scene Change Detection is yet another sophisticated video segmenting tool to find scenes or shots in a given video sequence. This tool is based on a model which uses edge detection techniques to obtain features for adjacent features. The frame is then processed block-wise to compare features in the corresponding blocks of adjacent frames. If the change in corresponding blocks for adjacent frame exceeds the specified threshold for the blocks, it is then able to determine the scene boundary.

Finally **PySceneDetect**, a Python [Pyt] based command line tool which provides various operations for scene detection in video sequences. It is an open source project being contributed by other users on Github [Bre19]. PyScenedetect uses various features based detection methods to find scenes in a given scene. The methods are based on various techniques of threshold only technique similar to ffmpeg blackfilter to advanced content aware detection of a scene or shot. The tool utilizes both open source computer vision libraries of OpenCV [Ope] and FFMPEG libraries.

In this thesis, *PySceneDetect* is used for finding scene or shot boundaries for any given video. The reason is that it allows for scene detection based content and threshold both. Secondly it is open source project whereas *Matlab Scene Change Detect tool* requires purchasing the Matlab and Simulink toolkits. Figure 3.2 shows the result in a tabular form when an input video sequence is used with PySceneDetect. The figure shows a comma separated values (csv) type of representation of the data extracted. Each row in the table shows a tuple of attributes for a single scene. E.g., in the figure 3.2 the first row is the heading for the attributes represented in each of the columns for the corresponding scene. The next row shows the attributes of first scene in each of the columns. The first column describes the *Scene Number* followed by the column *Start Frame*, which shows the start of the frame for the given scene. The value 0 is used for indexing the first frame and it increments frame by frame. The next two columns *Start Timecode* and *Start Time(seconds)* shows the starting time in timecode and seconds format. The next column *End Frame* shows the frame index where the scene ended followed by next two columns of *End Timecode* and *End Time(seconds)* same as start time formats of timecode and seconds. After the end time column follows the column named *Length (frame)* which shows the number of frames for the scene i.e., difference between the values at the column *End Frame* and column *Start Frame* for the given scene. Similarly the time difference is shown in both format of timecode and seconds in the columns *Length (timecode)* and

3.2 Initial Pipeline Design Setup

Scene Number	Start Frame	Start Timecode	Start Time (seconds)	End Frame	End Timecode	End Time (seconds)	Length (frames)	Length (timecode)	Length (seconds)
1		00:00:00.000	0	210	00:00:08.400	8.4	210	00:00:08.400	8.4
2	210	00:00:08.400	8.4	324	00:00:12.960	12.96	114	00:00:04.560	4.56
3	324	00:00:12.960	12.96	402	00:00:16.080	16.08	78	00:00:03.120	3.12
4	402	00:00:16.080	16.08	543	00:00:21.720	21.72	141	00:00:05.640	5.64
5	543	00:00:21.720	21.72	590	00:00:23.600	23.6	47	00:00:01.880	1.88
6	590	00:00:23.600	23.6	644	00:00:25.760	25.76	54	00:00:02.160	2.16
7	644	00:00:25.760	25.76	696	00:00:27.840	27.84	52	00:00:02.080	2.08
8	696	00:00:27.840	27.84	806	00:00:32.240	32.24	110	00:00:04.400	4.4
9	806	00:00:32.240	32.24	1051	00:00:42.040	42.04	245	00:00:09.800	9.8
10	1051	00:00:42.040	42.04	1093	00:00:43.720	43.72	42	00:00:01.680	1.68
11	1093	00:00:43.720	43.72	1301	00:00:52.040	52.04	208	00:00:08.320	8.32
12	1301	00:00:52.040	52.04	1405	00:00:56.200	56.2	104	00:00:04.160	4.16
13	1405	00:00:56.200	56.2	1624	00:01:04.960	64.96	219	00:00:08.760	8.76
14	1624	00:01:04.960	64.96	1756	00:01:10.240	70.24	132	00:00:05.280	5.28
15	1756	00:01:10.240	70.24	1829	00:01:13.160	73.16	73	00:00:02.920	2.92
16	1829	00:01:13.160	73.16	1894	00:01:15.760	75.76	65	00:00:02.600	2.6
17	1894	00:01:15.760	75.76	1979	00:01:19.160	79.16	85	00:00:03.400	3.4
18	1979	00:01:19.160	79.16	2467	00:01:38.680	98.68	488	00:00:19.520	19.52
19	2467	00:01:38.680	98.68	2535	00:01:41.400	101.4	68	00:00:02.720	2.72
20	2535	00:01:41.400	101.4	2708	00:01:48.320	108.32	173	00:00:06.920	6.92
21	2708	00:01:48.320	108.32	2828	00:01:53.120	113.12	120	00:00:04.800	4.8
22	2828	00:01:53.120	113.12	2865	00:01:54.600	114.6	37	00:00:01.480	1.48
23	2865	00:01:54.600	114.6	2890	00:01:55.600	115.6	25	00:00:01.000	1
24	2890	00:01:55.600	115.6	2963	00:01:58.520	118.52	73	00:00:02.920	2.92
25	2963	00:01:58.520	118.52	3008	00:02:00.320	120.32	45	00:00:01.800	1.8
26	3008	00:02:00.320	120.32	3046	00:02:01.840	121.84	38	00:00:01.520	1.52
27	3046	00:02:01.840	121.84	3211	00:02:08.440	128.44	165	00:00:06.600	6.6
28	3211	00:02:08.440	128.44	3251	00:02:10.040	130.04	40	00:00:01.600	1.6
29	3251	00:02:10.040	130.04	3320	00:02:12.800	132.8	69	00:00:02.760	2.76
30	3320	00:02:12.800	132.8	3456	00:02:18.240	138.24	136	00:00:05.440	5.44
31	3456	00:02:18.240	138.24	3541	00:02:21.640	141.64	85	00:00:03.400	3.4
32	3541	00:02:21.640	141.64	3612	00:02:24.480	144.48	71	00:00:02.840	2.84
33	3612	00:02:24.480	144.48	3674	00:02:26.960	146.96	62	00:00:02.480	2.48
34	3674	00:02:26.960	146.96	3700	00:02:28.000	148	26	00:00:01.040	1.04
35	3700	00:02:28.000	148	3730	00:02:29.200	149.2	30	00:00:01.200	1.2
36	3730	00:02:29.200	149.2	3803	00:02:32.120	152.12	73	00:00:02.920	2.92
37	3803	00:02:32.120	152.12	3912	00:02:36.480	156.48	109	00:00:04.360	4.36
38	3912	00:02:36.480	156.48	3955	00:02:38.200	158.2	43	00:00:01.720	1.72
39	3955	00:02:38.200	158.2	3994	00:02:39.760	159.76	39	00:00:01.560	1.56
40	3994	00:02:39.760	159.76	4029	00:02:41.160	161.16	35	00:00:01.400	1.4
41	4029	00:02:41.160	161.16	4098	00:02:43.920	163.92	69	00:00:02.760	2.76
42	4098	00:02:43.920	163.92	4152	00:02:46.080	166.08	54	00:00:02.160	2.16
43	4152	00:02:46.080	166.08	4320	00:02:52.800	172.8	168	00:00:06.720	6.72

Figure 3.2: Scene Information in CSV Format of a Random Video Sequence using PySceneDetect

column *Length (seconds)*. The scenes are detected and placed in the order as it gets detected in a given video sequence. After finalising the scene detection tool, initial setup was designed.

3.2 Initial Pipeline Design Setup

The motivation with this design was to test if encoding a single resolution with different bitrate settings would yield out the quality scores similar to one of the rate-quality curves from figure 2.1. The encodes were generated using bitrates ranging from 25 kbps till 10 Mbps. A total of 67 points in between 25 Kbps to 10 Mbps were chosen. The initial 29 points ranged between 25 Kbps to 750 Kbps with interval of 25 Kbps. The rest of the points ranging from 750 Kbps till 10 Mbps had an interval of 250 Kbps. The reason for having bitrate steps was that the Rate-Quality curve tends to rise quickly in the lower bitrate region and then flattens in higher bitrate region.

The figure 3.3 shows the overview of the video encoding pipeline used for the approach for encoding a single resolution of 720p and its key frames are shown in figure 3.4. For its inter-frame and intra-frame characteristics a plot for spatial and temporal information was generated

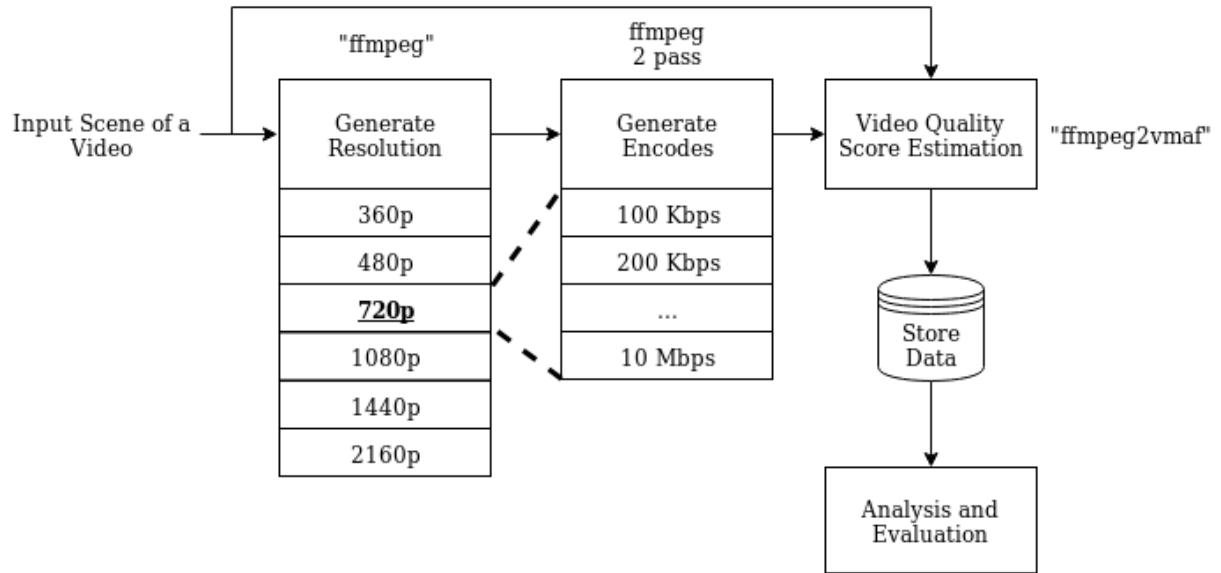


Figure 3.3: Video Encoding Pipeline for only 720p Resolution



Figure 3.4: Key Frames of the Hongkong Scene

as shown in the figure 3.5. First a scene is provided as an input to the pipeline and the first block generates its different resolutions from 360p, 480p, 720p, 1080p, 1440p till 2160p i.e., in total 6 different resolutions of the scene. Then the 720p resolution of the scene is forwarded to next block where the scene is transcoded using several bitrate points as described above for 67 different bitrate points.

The encodes generated are then passed to the next block *Video Quality Score Estimation*. Each of the encodes are analyzed with respect to original scene to estimate the quality. The metric as already described is VMAF from Netflix. This open source project provides various command line tools for video quality estimation. One of the command line tools which is used here is "ffmpeg2vmaf". This command line tool is provided with the original video and encoded video as arguments. The tool compares the encoded video with the original video and generates results in various format such as json, text and xml. One of the sample output can be seen in the figure 3.6. The only result that was used in this thesis is the VMAF_score, the second last entry in figure 3.6. The other generated values are the features extracted from

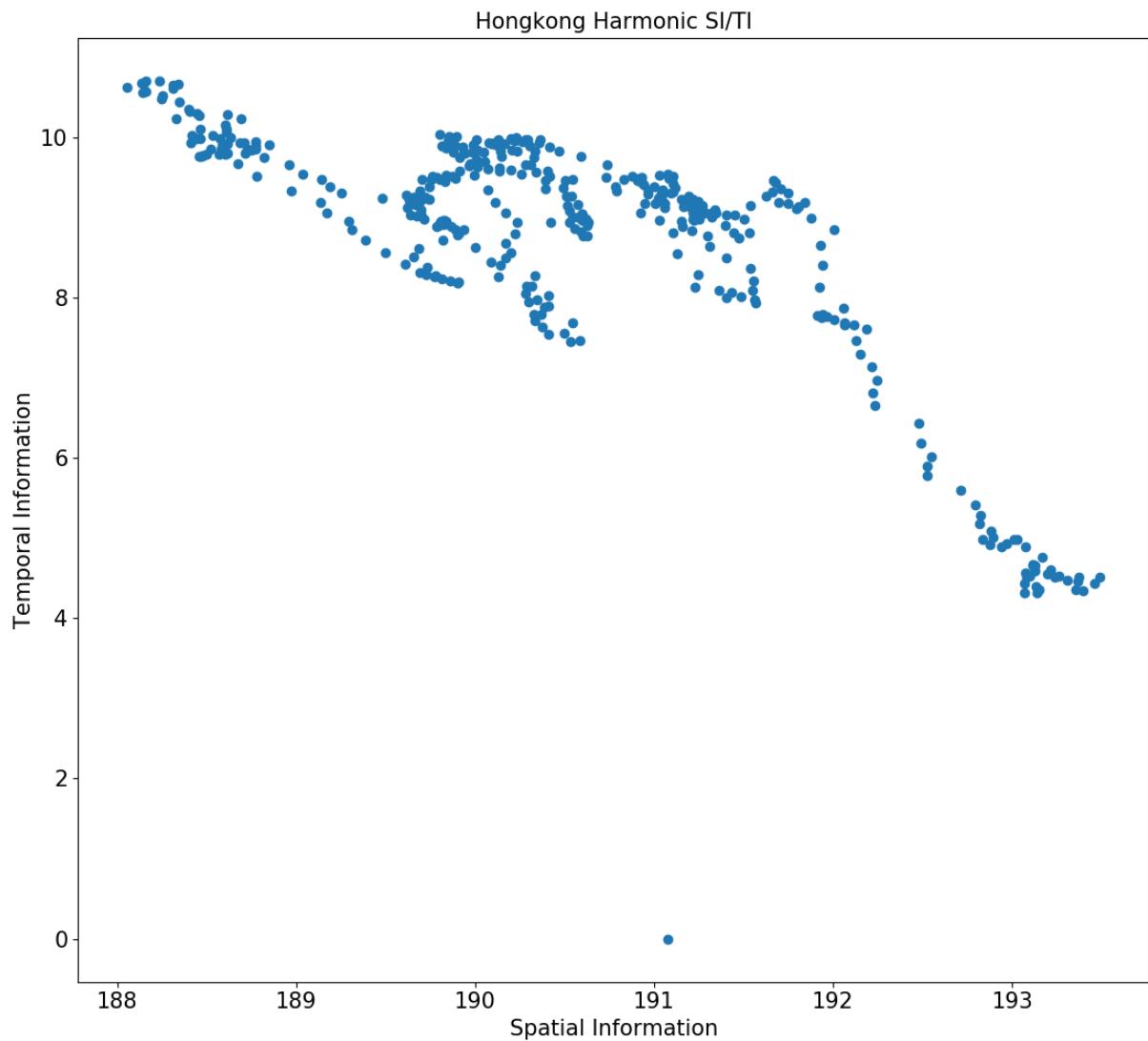


Figure 3.5: Plot for Information versus Temporal Information Plot for Hongkong Scene

```

"aggregate": {
    "VMAF_feature_adm2_score": 0.93458780776205741,
    "VMAF_feature_motion2_score": 3.8953518541666665,
    "VMAF_feature_vif_scale0_score": 0.36342081156994926,
    "VMAF_feature_vif_scale1_score": 0.76664738784617292,
    "VMAF_feature_vif_scale2_score": 0.86285338927816291,
    "VMAF_feature_vif_scale3_score": 0.91597186913930484,
    "VMAF_score": 76.699271371151269,
    "method": "mean"
}

```

Figure 3.6: Sample VMAF Output[Netb]

the encode after comparing with the original scene per frame and also per frame VMAF scores are generated.

The Figure 3.7 shows the schematic of the plot generated with quality scores generated using VMAF metric against the bitrates used for encoding the scene. The curve rises quickly in the region of bitrates from 25 Kbps to 2 Mbps, then it starts to flatten. It shows that quality does not increase as much as it increases in lower bitrate points.

The rate-quality curve in figure 3.7 turned out to be as expected. Next step was to setup the pipeline for all the resolutions of the video. The figure 3.8 shows the pipeline used in determining rate quality curve for all resolutions. Each of resolutions generated are subjected to 67 different bit rate points. For each of the encodes, all VMAF calculations are done using "ffmpeg2vmaf" tool. The results for each of the encodes are separated based on the resolutions. Encodes and their scores were grouped together while analyzing the quality of the encodes. After the evaluation of the scores computed for the corresponding encodes, the plot was generated.

The figure 3.9 shows rate vs quality curve for all the resolutions. The observation at curves of different resolution shows the characteristics about the intersection points between curves that can be used to generate convex hull above all the curves for individual resolutions. The curve generated using VMAF scores for each bitrate points corresponding to each of the resolutions are way too much detailed. In order to have clear picture, some of the extra points can be ignored and plot only the remaining points of rate-quality pair corresponding to each resolution. Consecutively this would provide much clear curves. For a given bitrate, if there exists quality points for more than one resolutions then the low resolution quality points must be ignored because higher resolution with high quality for a given bitrate would offer more **quality of experience in the subjective sense for a viewer.** After trimming down the irrelevant points of rate-quality pairs and then setting up the plot again would allow to understand it better. **The resulting** rate-quality points in the plot looks like figure 3.10.

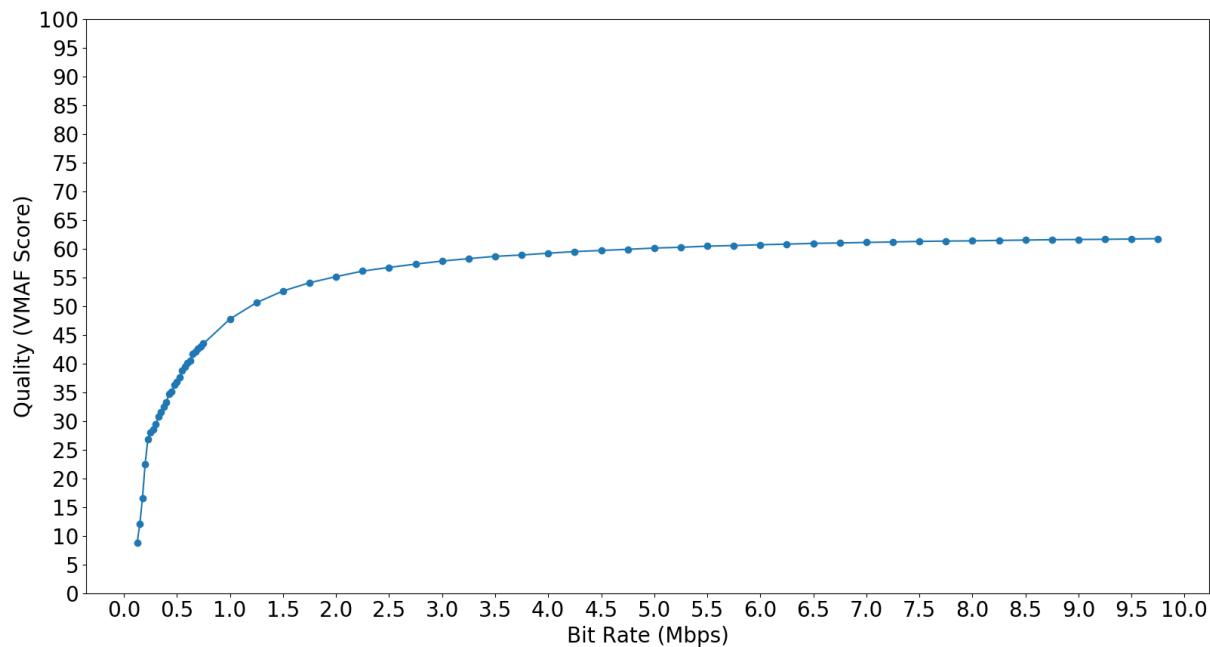


Figure 3.7: Plot for Rate-Quality Curve at 720p resolution for Hongkong scene

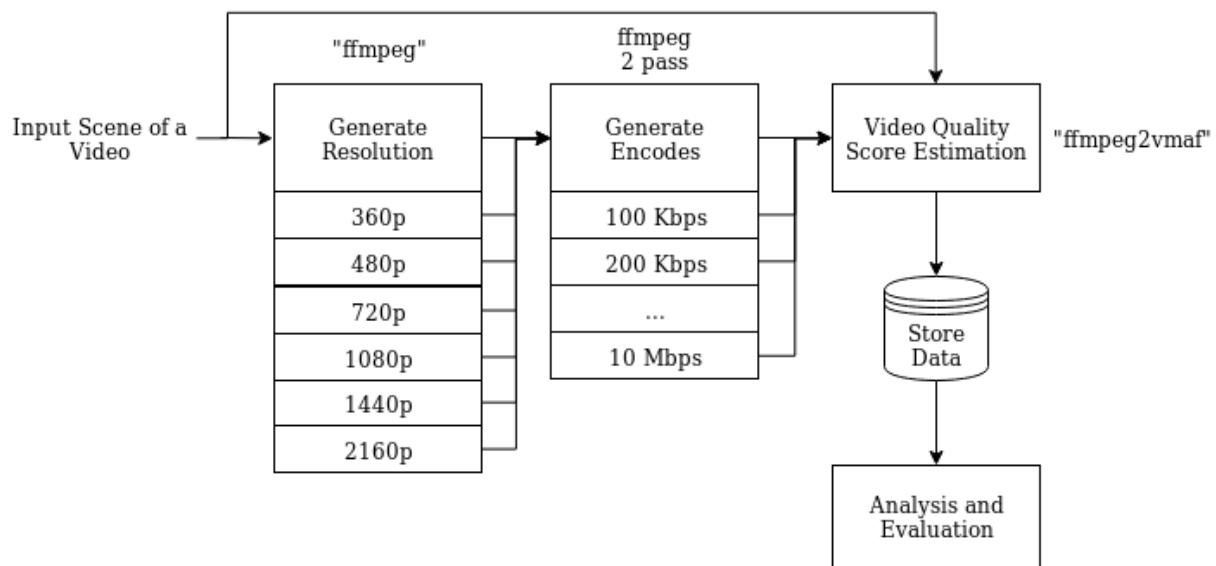


Figure 3.8: Video Encoding Pipeline for all Resolution

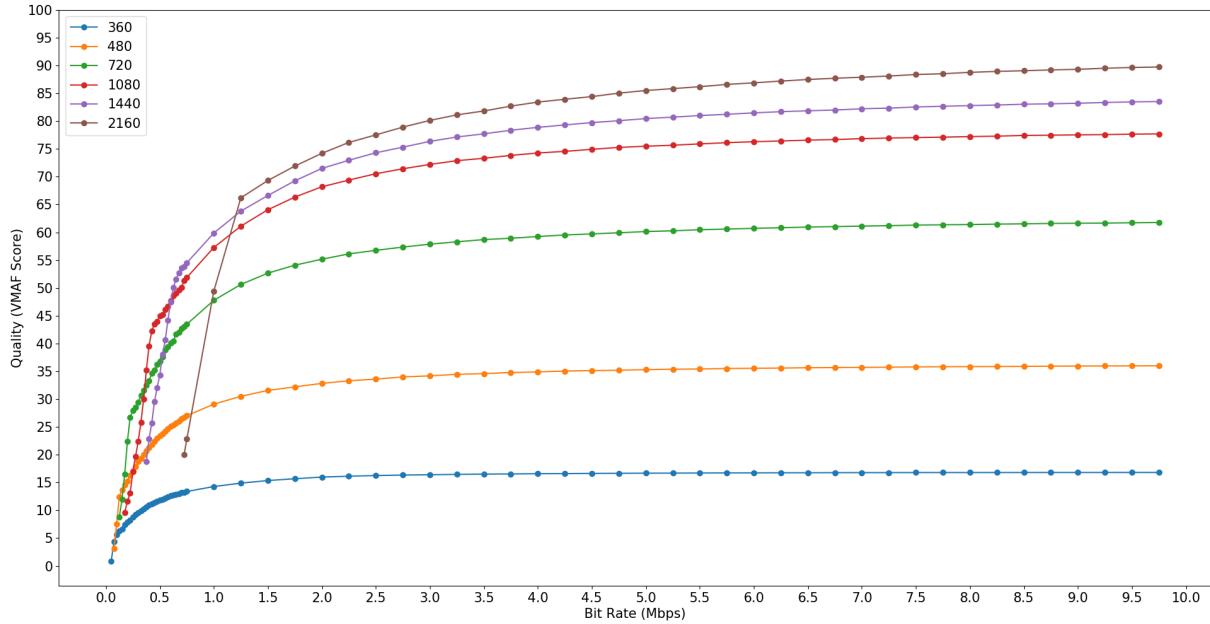


Figure 3.9: Plot for Rate-Quality Curves for Hongkong scene at all resolutions over different Bitrates

Analyzing the plot in the figure 3.10 shows that the curves for low resolutions i.e., 360p, 480p and 720p the frequency rate-quality pairs is too high when it comes to change in quality. The change in quality for those curve is not as rapid as it is for the curves of higher resolutions i.e., 1080p, 1440p and 2160p. This analysis shows that the curves for the higher resolutions were able to achieve quality higher VMAF Score than 60 at the 1.5 Mbps. The reason behind it is that the scene that is being analyzed here had very smooth and slow movements with most of the uniform regions across the frames and it would be nice to know which sequence it is. The figure 3.5 shows the plot with spatial Information ranging in a very short range of approximately 188 to 193, which is comparatively much higher than temporal information range of approximately 4 to 7. Such video scenes where the movement is least and background remains more or less constant they reach higher quality at lower bitrates as well. For dynamic videos the curves might take higher bitrates ranging in 5-10 Mbps to reach satisfactory quality for the viewers. The rate-quality curves tend to have much higher slope as the bitrate increases initially for videos with low amount of motion than the rate-quality curve of dynamic scenes where there is sudden motions and content change. In this thesis there were several other videos used with unique attributes. Some of the videos had dynamic motion in it and after the analysis when the rate-quality curve was plotted. The curves did not flatten as much as the curves for smooth video scenes. The further sections of this thesis will show the proof for the claim being made about relationship of the characteristics of rate-quality curve with the different scenes. Take the figure 3.7 into consideration, for the bitrate points below 2.5 Mbps, there was reasonable quality difference after each bitrate points from the start but after 2.5 Mbps the curve tends to flatten, which signifies that the quality will not improve as much after

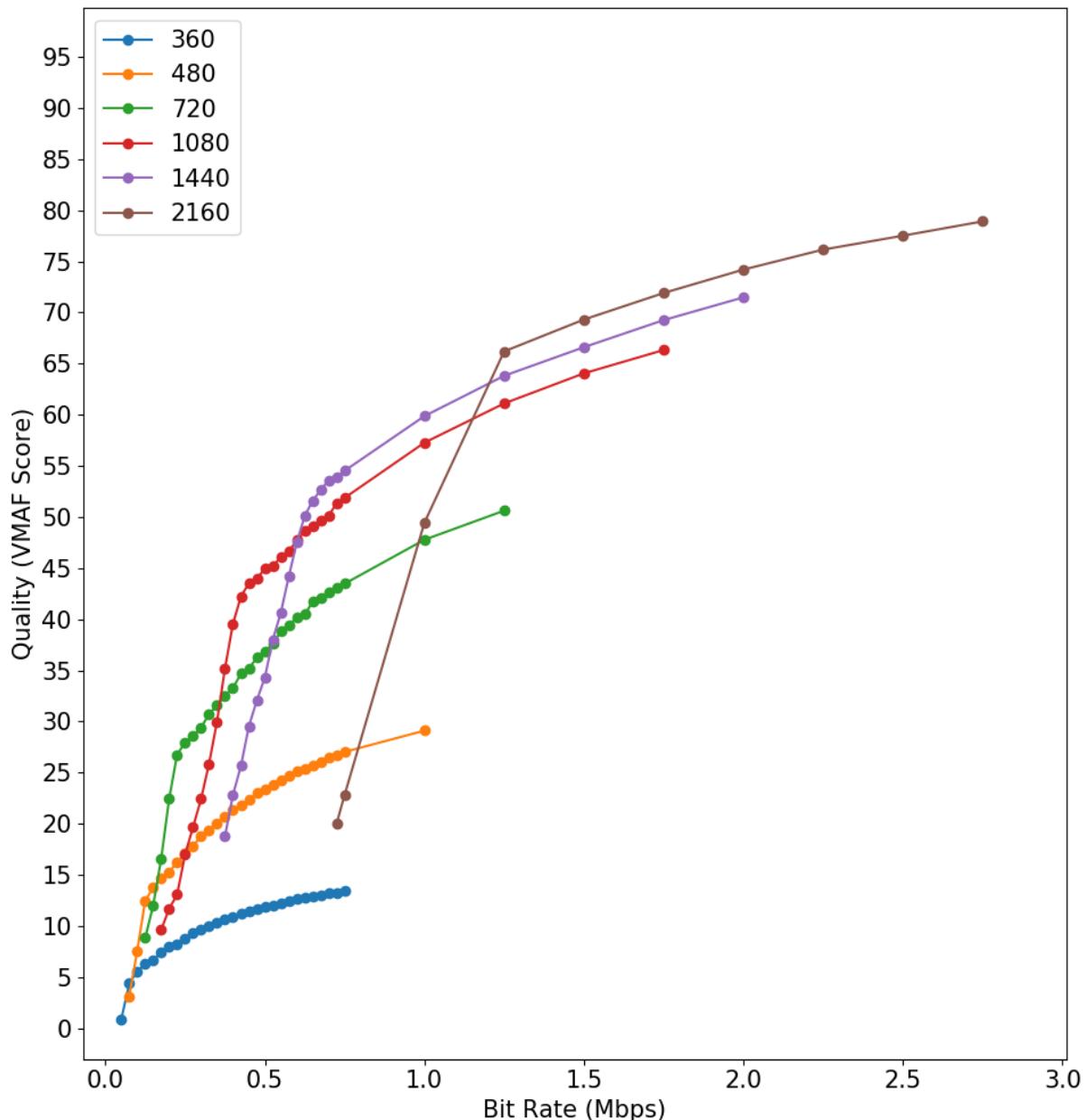


Figure 3.10: Plot for rate-Quality Curves for Hongkong scene at all resolutions over "Selected" Bitrates

2.5 Mbps bitrate. Now this point may vary from video to video depending upon the content of the video.

To summarize about the initial design setup and its results, it yielded smooth plots for individual resolutions and looked promising. This approach could be used to analyze other scenes or shots as well. But this approach has a downside. The number of bitrate points required for each of the resolutions for computing VMAF score quality can not be anticipated. The bitrate points for the resolutions would also vary in case of different kind of video scenes. Another reason is the computational resource limitation because the process of transcoding requires intensive computations. **Netflix uses pipeline in the cloud based on parallel computing. Their cloud uses the Linux EC2 instances which provides powerful computing resources[Blo17].** For this thesis, the computational resources were limited so transcoding a scene for 6 resolutions and then computing 67 encodes for each of the resolutions reaches very high computational complexity. The problem of computational cost had to be reduced in order to find a feasible solution that can provide enough bitrate points to do the transcoding. Secondly the bitrate points must be recurring where the quality for those bitrate points really have reasonable difference. This required to develop a new approach which could divide and conquer the bitrates points which are relevant. So the next section would describe about the new approach and how it performs for different video scenes.

3.3 Pipeline based on Bisection Method Approximation

As described in the previous section 3.2, about the pipeline and its limitation, a new approach was needed which could search for the relevant bitrate points used in transcoding for quality score estimation. Secondly, the resolution of bitrate points can be reduced i.e., in the previous approach the bitrate points in the lower bitrate had a 25 Kbps resolution. So to avoid the long computation hours, the high resolution of bitrate points has to be traded off. The ideas for selecting the bitrate points where the quality increment is higher can be based on numerical analysis. Some numerical analysis methods use the approaches to find the roots for a given mathematical equation of a curve by choosing the initial set of points on the curve to find the next suitable point. With this approach, when a new point is found the process is reiterated but now using the new obtained point as part of the method. This iteration keeps on going until the root(s) are evaluated. The approach of numerical analysis used in this thesis is called Bisection Method [Cor77].

The bisection method is also known as Interval Halving in other terms. The bisection method starts with set of two points on the curve between which a root can be found. The concept is that, for a set of two points on x_a and x_b axis on a curve $f(x)$ being a continuous curve, if

3.3 Pipeline based on Bisection Method Approximation

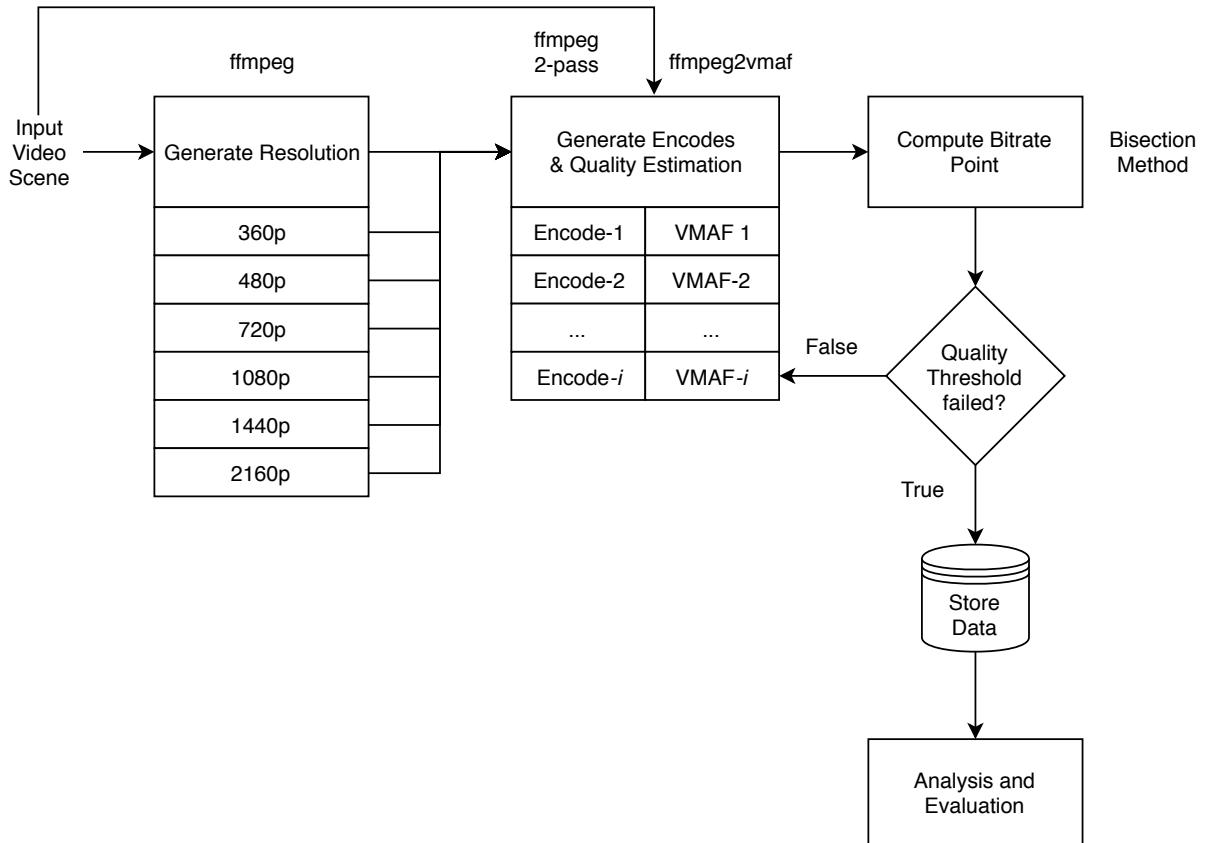


Figure 3.11: Video Encoding Pipeline based on Bisection Method Approximation

the there is change of sign on the given set of two x points then there exists a root between the given set of two values i.e., if $f(x_a) \cdot f(x_b) < 0$ then there exists a root between the points x_a and x_b . The curve must be plotted in prior to analyze the curve visually to observe the change of sign. The method is iterated over the update for the initial set of points. To find the root for $f(x)$ between points x_a and x_b such that $f(x_a) \cdot f(x_b) < 0$ the following instructions iterated:

- ▷ Repeat the following step till the tolerance condition is met:
 - ◊ Find a mid-point value x_c between given two points x_a and x_b i.e., $x_c = (x_a + x_b)/2$
 - ◊ Check if $f(x_a) \cdot f(x_c) < 0$ or if $f(x_b) \cdot f(x_c) < 0$
 - ◊ If $f(x_a) \cdot f(x_c) < 0$ then update $x_b = x_c$ else if $f(x_b) \cdot f(x_c) < 0$ then update $x_a = x_c$
 - ◊ Repeat the above steps until a tolerance condition is met i.e., $|x_a - x_b| < 2 \cdot tolerance$
- ▷ Once the tolerance condition fails then stop. The latest x_c point is the approximate root with error less than $|x_a - x_b| < 2$

3.3.1 Theoretical Approach with Bisection

The new pipeline has used the similar approach of Bisection method for approximating the bitrate points based on the quality scores obtained on those points. The pipeline described in the figure 3.11 is based on the numerical approach of root finding using Bisection Method. The approach describes the method to obtain the optimal bitrate points where the occurrence of the quality difference between the selected bitrate points is higher. However in this method instead of finding the roots, the next bitrate point is found using the quality scores for the given bitrate points. The theoretical approach can be described by the following points:

- ◊ Repeat the following step until the tolerance condition is met;
 - ◊ Select initial set of two bitrate points R_1 and R_2
 - ◊ Transcode a given resolution of the scene at the given two bitrate points
 - ◊ Compute VMAF quality scores Q_1 and Q_2 at the initial bitrate points
 - ◊ Compute a third bitrate point such that $R_3 = (R_1 + R_2)/2$ and its corresponding VMAF quality score Q_3
 - ◊ Check if the quality difference $Q_3 - Q_1 > Q_2 - Q_3$, then update Q_2 with the value Q_3 else update Q_1 with the value Q_3
 - ◊ Repeat until the quality difference remains greater than the threshold value i.e., for updated bitrate points R_1 and R_2 with their corresponding quality score Q_1 and Q_2 if the quality difference $|Q_1 - Q_2|$ is less than the threshold value
- ▷ Once the threshold condition fails, stop the iteration

The Bisection Method based video encoding pipeline is shown in the figure 3.11 which consists of the blocks described in detail in the following subsections.

3.3.2 Generate Resolutions

The approach here is same as described in the previous pipeline approach of 3.2 for generating a total of 6 resolution versions of the input scene.

3.3.3 Generate Encodes and Estimate Quality

This block is different from the previous approach. Each of the resolutions of the input scene is processed and transcoded at two initial bitrate points of 0.1 Mbps and 10 Mbps as described

3.3 Pipeline based on Bisection Method Approximation

in 3.3.1. Based on the bitrate points computed from the next block *Compute Next Bitrate Point*, its transcoding and quality score of the encode is estimated.

3.3.4 Compute Next bitrate point

After transcoding the corresponding VMAF, quality scores are evaluated at the points. The next bitrate point is computed by taking the mid-point of the initial bitrate points 0.1 Mbps and 10 Mbps i.e., $(0.1 + 10)/2 = 5.05$. Then the VMAF quality score is evaluated at this new bitrate point of 5.05 Mbps. The quality differences for both set of bitrate points of 0.1 Mbps 5.05 Mbps and 5.05 Mbps 10 Mbps is calculated. If the quality difference for the points 0.1 Mbps and 5.05 Mbps is greater than the quality difference for the points 5.05 Mbps and 10 Mbps then the next set bitrate points are updated as 0.1 Mbps and 5.05 Mbps, else the next set of points are updated as 5.05 Mbps and 10 Mbps. The computation of the new bitrate points updating keeps on iterating until the threshold condition fails. The next block is responsible for decision making for the iteration to continue.

3.3.5 Quality and Bitrate Threshold

For the initial approach only quality threshold was used. If the quality difference between subsequent bitrate points, say at 0.1 Mbps and 5.05 Mbps or the quality difference for points at 5.05 Mbps and 10 Mbps is greater than the quality threshold then the iteration continues to compute new bitrate point. The block *Generate Encodes and Estimate Quality* transcodes the scene with the new computed bitrate point and consequently the VMAF score for the encode is estimated. But when the any of quality differences fails to be greater than the quality threshold, then the iteration stops. The generated encodes and its VMAF score is then forwarded to further analysis. The whole process goes on for every resolution version of a given scene. The encodes and their respective VMAF quality scores are grouped together on the basis of the resolutions. The data obtained is then forwarded for further analysis.

The video scene that can be used for testing the video encoding pipeline based on bisection method is *Skateboarding Scene*. The keyframes for the scene are shown in the figure 3.12. The video scene is about a person skateboarding in a skate-park. As the person moves around in this scene the spatial information varies over time as the position of the person changes after every few frames but the major part of skate-park remains almost similar. The characteristics about spatial and temporal information variation can be seen in the figure 3.13. The plot signifies about the video scene behaviour and how it keeps changing over time. The Spatial Information varies in the approximate range from 100 to 220.



Figure 3.12: Key Frames of the Skateboarding Scene in the sequence as 1) top-left, 2)top-right, 3)bottom-left and 4)bottom-right

After applying the new pipeline based on bisection method on *Skateboarding Scene* resulted into an expected rate-quality curve. The generated rate-quality curve can be shown in the figure 3.14. The generated plot shows the same characteristics of a rate-quality curve. The overall computation processing was reduced compared to the initial pipeline approach from 3.2.

After trimming down the unwanted bitrate points the curve can be plotted again with the relevant bitrate points as shown in the figure 3.16.

Observation

The performance of the bisection method based encoding pipeline was not significantly different than to the performance of initial approach of encoding pipeline 3.2. The observation while transcoding the videos based on bisection method was that it worked good for lower resolutions i.e., number of bitrate points generated with this approach, were comparatively less for the resolutions 360p, 480p and 720p with 12, 13 and 25 encodes respectively, on the other hand the number of encodes generated were greater than or equal to 45 encodes for higher resolutions of 1080p, 1440p and 2160p. The quality difference between subsequent bitrate points was used as the threshold for the pipeline. So based on the encodes generated it was observed that the threshold worked effectively with lower resolutions. But for higher resolution the threshold could was not so effective. The bisection algorithm went on to compute encodes with bitrate

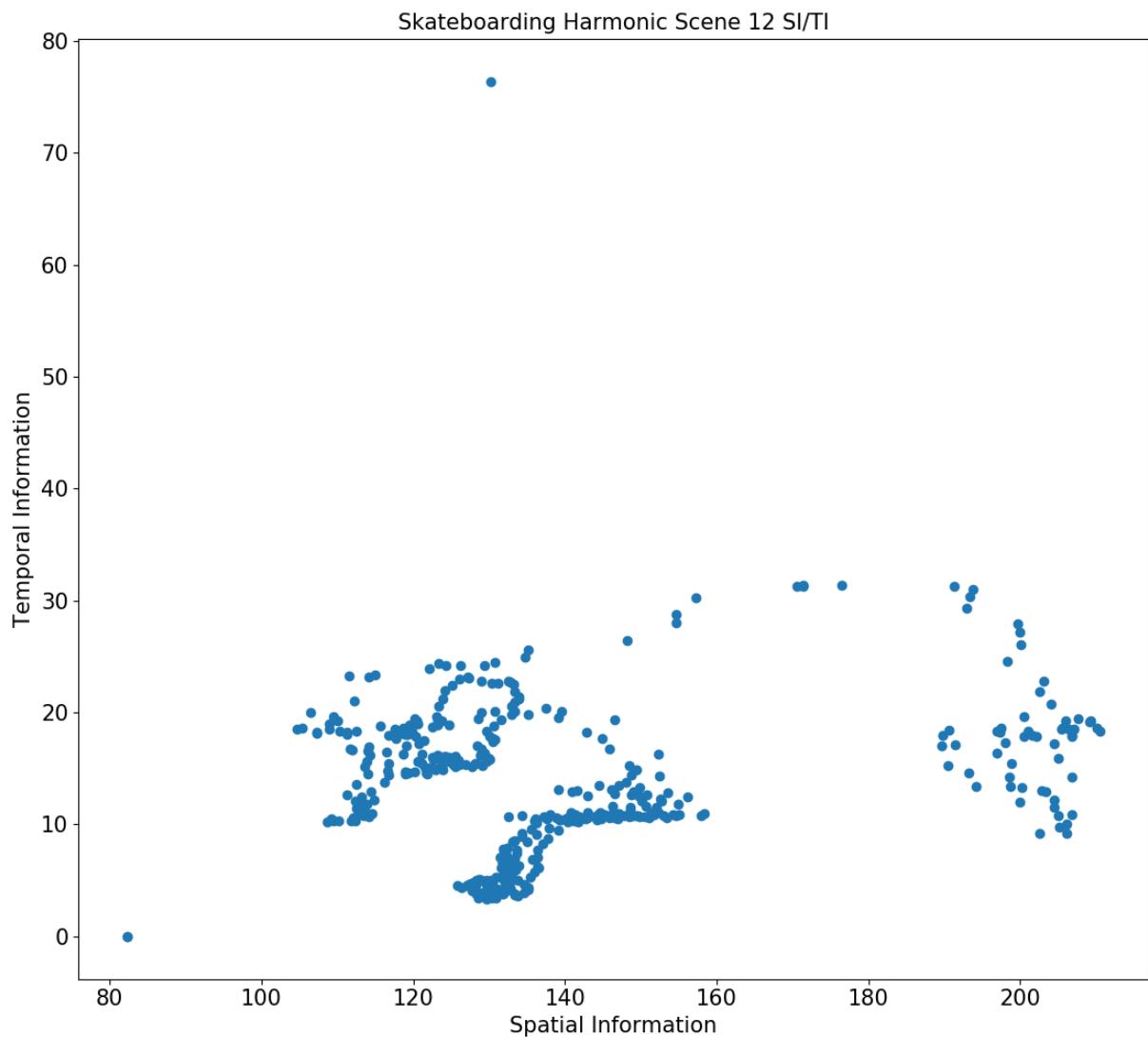


Figure 3.13: Plot for **Spatial Information** versus **Temporal Information** for the Skateboarding Scene

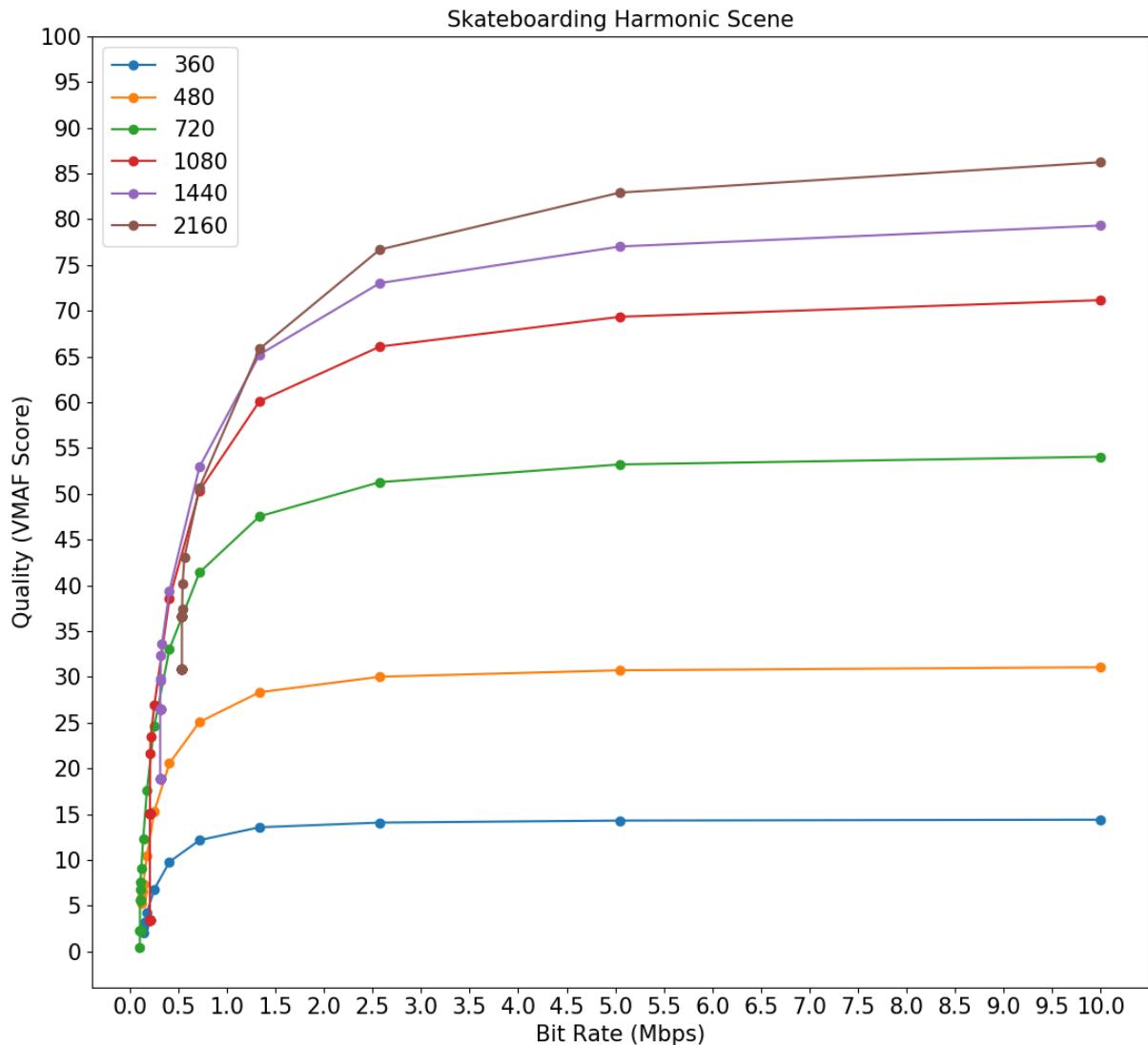


Figure 3.14: Plot for Rate-Quality curves for all resolutions for Skateboarding Scene

3.4 Pipeline based on Bisection Method Approximation with Bitrate Threshold

```
skateboarding_harmonic_scene_12_480_0.1193359375M.abc
skateboarding_harmonic_scene_12_480_0.124169921875M.abc
skateboarding_harmonic_scene_12_480_0.12900390625M.abc
skateboarding_harmonic_scene_12_480_0.138671875M.abc
skateboarding_harmonic_scene_12_480_0.17734375M.abc
skateboarding_harmonic_scene_12_480_0.1M.abc
skateboarding_harmonic_scene_12_480_0.2546875M.abc
skateboarding_harmonic_scene_12_480_0.409375M.abc
skateboarding_harmonic_scene_12_480_0.71875M.abc
skateboarding_harmonic_scene_12_480_10.0M.abc
skateboarding_harmonic_scene_12_480_1.3375M.abc
skateboarding_harmonic_scene_12_480_2.575M.abc
skateboarding_harmonic_scene_12_480_5.05M.abc
```

Figure 3.15: List of encodes generated using Bisection Algorithm for Skateboarding Scene at 720p

values going up to twelfth decimal points on Mbps scale. The figure 3.15 shows the list of encodes generated for *Skateboarding Scene* for the resolution of 720p. Observing the plot in the figure 3.14 closely, does not show large number of rate-quality points, the reason is most of the points which were computed using bisection algorithm were getting close to the quality threshold value and the bitrate points were so close that they look as one point in the plot. In the case of resolution 1080p for *Skateboarding Scene* there were 20 bitrate points computed at approximately 0.207 Mbps, secondly for 1440p resolution, there were roughly 30 bitrate points close to each other at approximate bitrate of 0.314 Mbps and finally for 2160p resolution, there were close to 39 points at approximate bitrate of 0.5379 Mbps. This problem has to be rectified by means of change ion threshold value but different resolutions have different varying quality score steps as the bitrate increases. So the threshold on quality can takes care of low resolution versions of the scene, there is need of threshold such that the video scene encoding is done with minimum relevant bitrate points.

3.4 Pipeline based on Bisection Method Approximation with Bitrate Threshold

The new pipeline is based on the same approach of Bisection Method described in 3.3.1 but with a new threshold combined with the existing threshold value. As observed before the bitrare interval halving let to very minute differences in the subsequent bitrate points which is more than required to transcode at such points and estimate the VMAF quality score. The decision making of bisection of initial set of points still depends upon the quality difference for encodes generated for those points but the threshold at which halving of the bitrate points should stop, has to depend on the bitrate difference. So the quality estimation will go on until a reasonable minimum bitrate difference is reached. This way the bitrate points resolution will not go more than the threshold value and make the algorithm computation altogether much

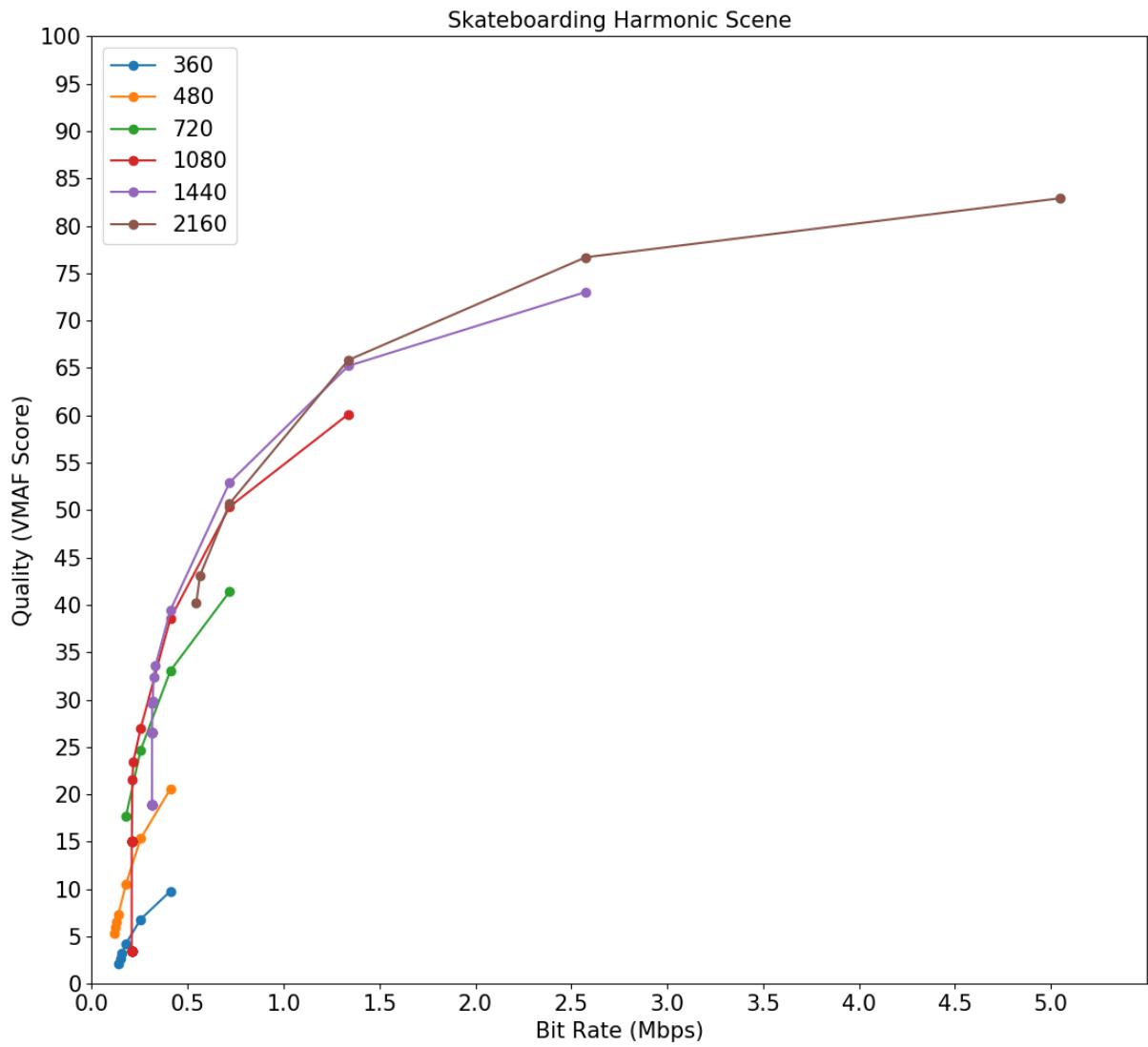


Figure 3.16: Plot for Rate-Quality curves for all resolutions for Skateboarding Scene with "Selected" Bitrate points

3.4 Pipeline based on Bisection Method Approximation with Bitrate Threshold

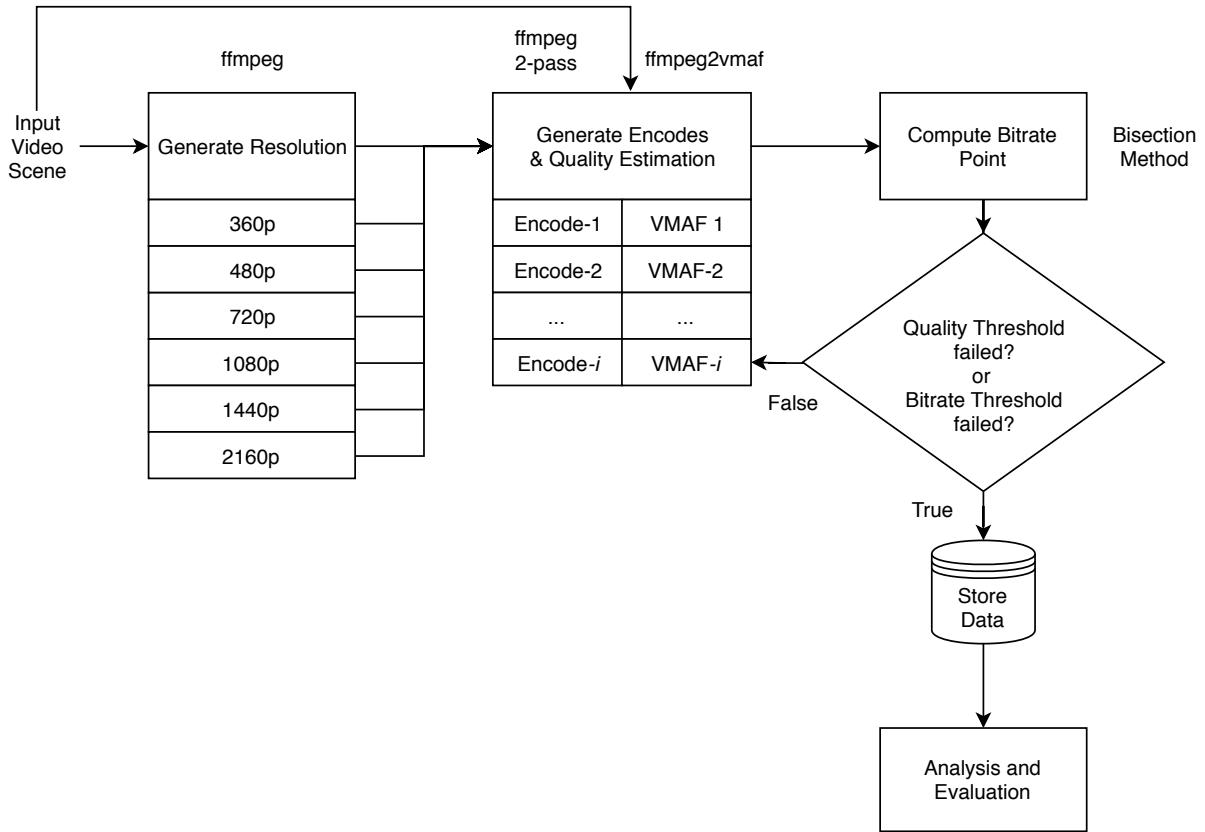


Figure 3.17: Video Encoding Pipeline based on Bisection Method Approximation with "Additional" Bitrate Threshold

faster. This approach will also avoid the bitrate points upto high decimal point values as it was in the case of previous approach where the bitrate values went up to twelfth decimal points. This approach again has the trade-off between computation time and bitrate points interval i.e., approaches with low computation time will have higher interval between bitrate points and vice versa. This new pipeline is shown in the figure 3.17. The blocks in the pipeline for figure 3.17 are described as follows:

3.4.1 Generate Resolutions

This block has the same function as it had before in previous approaches of encoding pipelines in sections 3.2 and 3.3.1 respectively. It generates the a total of 6 resolution versions of the given video scene as an input to the pipeline.

3.4.2 Generate Encodes & Quality Estimation

This block is responsible for generating the encodes at two bitrate points of 0.1 Mbps and 10 Mbps respectively using two-pass video encoding. Once the encodes are generated, their

quality scores are computed with comparison to the original scene using VMAF command line tool. The next set of points are decided by the next block of *Compute Bitrate Points* using the approach of Bisection Method approximation and decision is made by the diamond block of threshold value check.

3.4.3 Compute Bitrate Points

The next next point is calculated by the computing the mid-point of the current set of bitrate points used previously for transcoding and then VMAF quality score estimation for the encodes. This block keeps generating the new bitrate points and feeds it to the *Generate Encodes Quality Estimation* until the decision from next diamond block is to stop the iteration.

3.4.4 Quality Threshold/Bitrate Threshold

The iterative algorithm needs a condition or two to break the loop and stop the iterative computation. The conditions for threshold of iteration were previously based on the quality difference between two subsequent bitrate points which became inefficient for higher resolutions. The need for optimizing this process of computing new bitrate points was required which will avoid irrelevant and very close bitrate points. Hence the need of a new threshold was required. The new threshold value deals with the difference of the current set of two bitrate points which will be used to compute the next bitrate point for transcoding and quality estimation. If the difference between the given bitrate points drops below threshold then the iteration stops.

To test the method of video encoding using the new pipeline with quality and bitrate threshold, one of the [video scenes titled *American Football Scene*](#) was used. The keyframes for the video are shown in the figure 3.18.

If the key-frames are observed closely, the second key frame looks blurry due to quick motion of the video scene. The characteristics can be analyzed by plotting the intra-frame and inter-frame characteristics to display spatial and temporal information frame by frame for the scene. The figure 3.19 shows the plot spatial and temporal information for all the frames. Each dot represents spatial and temporal information for each frame. Based on the plot, it can be observed how the spatial and temporal information is varying over the frames. The spatial information has varied in range from approximately 75 to 275 whereas the temporal information has varied from in the approximate range of 0 to 45. The video describes a pass being made by one of the players to another player across the field i.e., the video had a lot motion.

3.4 Pipeline based on Bisection Method Approximation with Bitrate Threshold



Figure 3.18: Key Frames for American Football Scene with sequence; 1) top-left, 2) top-right and 3)bottom

Observation

This pipeline worked better than the previous approaches. The bitrate points computed using this approach did not get as close as it were close in last approach. After applying the new pipeline to the video scene of *American Football Scene*, the plot for rate-quality curves are generated. The plot for rate-quality curves are shown in the figure 3.20. **The characteristics of the *American Football Scene* differs with the characteristics of previous video scenes.** The reason for the difference is due to a lot of motion in the give video scene. Take the rate-quality curve in the figure 3.20 for the resolution of 360p (Blue colored curve), it can be noticed that the the VMAF quality score reaches maximum reasonable quality just before the bitrate of 3 Mbps. This phenomena occurs when the video characteristics has a high rate of change of content over time i.e., the frame information varies over time rapidly, as it was also observed in the plot 3.19 how the information is varying over the frames in the video scene.

Similarly other curves for the resolutions of 480p, 720p, 1080p, 1440p, and 2160p required higher bitrates to reach the maximum achievable quality scores for respective resolutions. The camera motion generates the additional dynamic change in motion over the frames for the given scene. The given video can be analyzed better by trimming down the irrelevant points similar to the previous videos. The resulting plot looks like, as shown in the figure 3.21.

3.4.5 Live Plotter

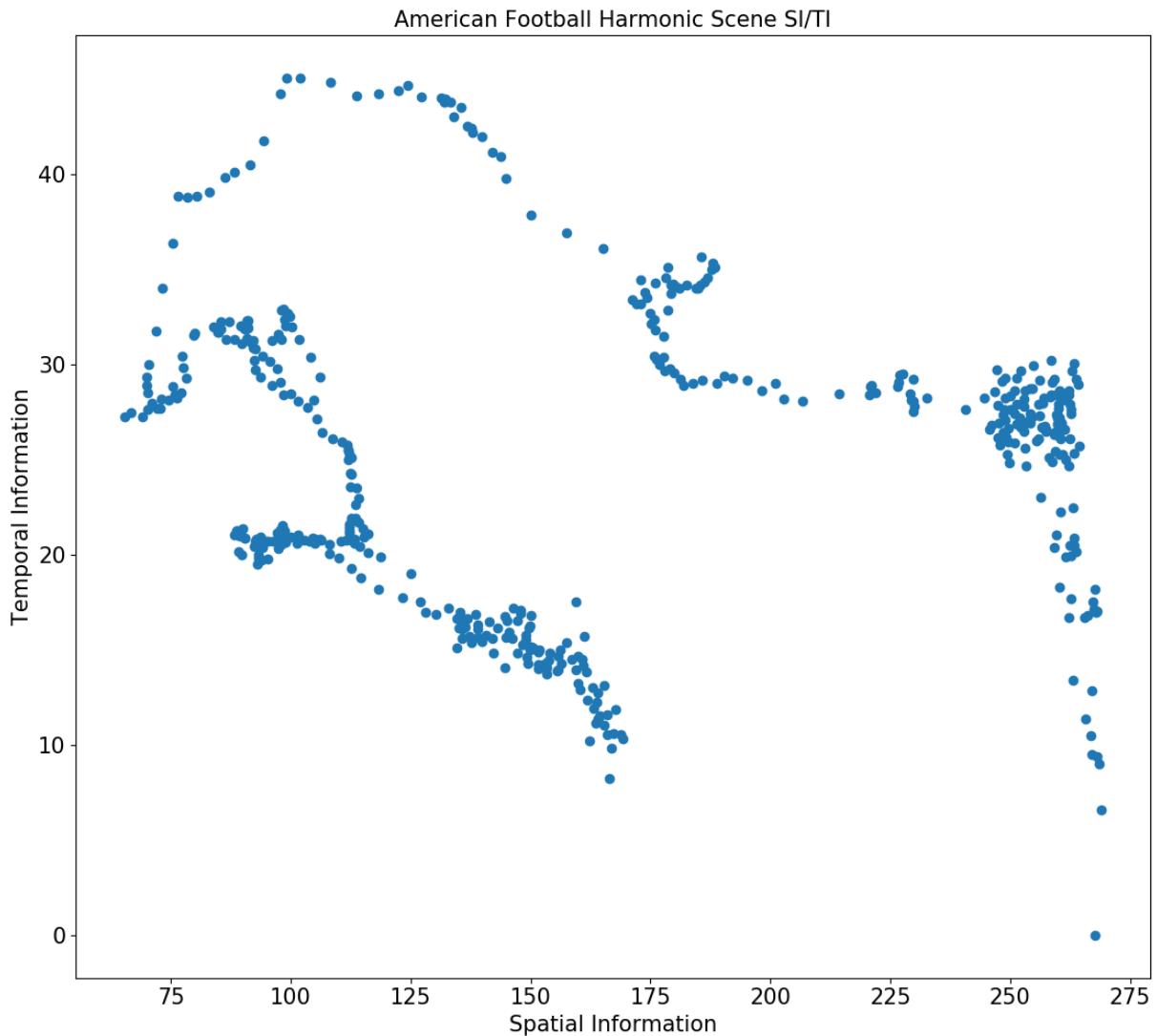


Figure 3.19: Plot for Spatial Information versus Temporal Information for American Football Scene

3.4 Pipeline based on Bisection Method Approximation with Bitrate Threshold

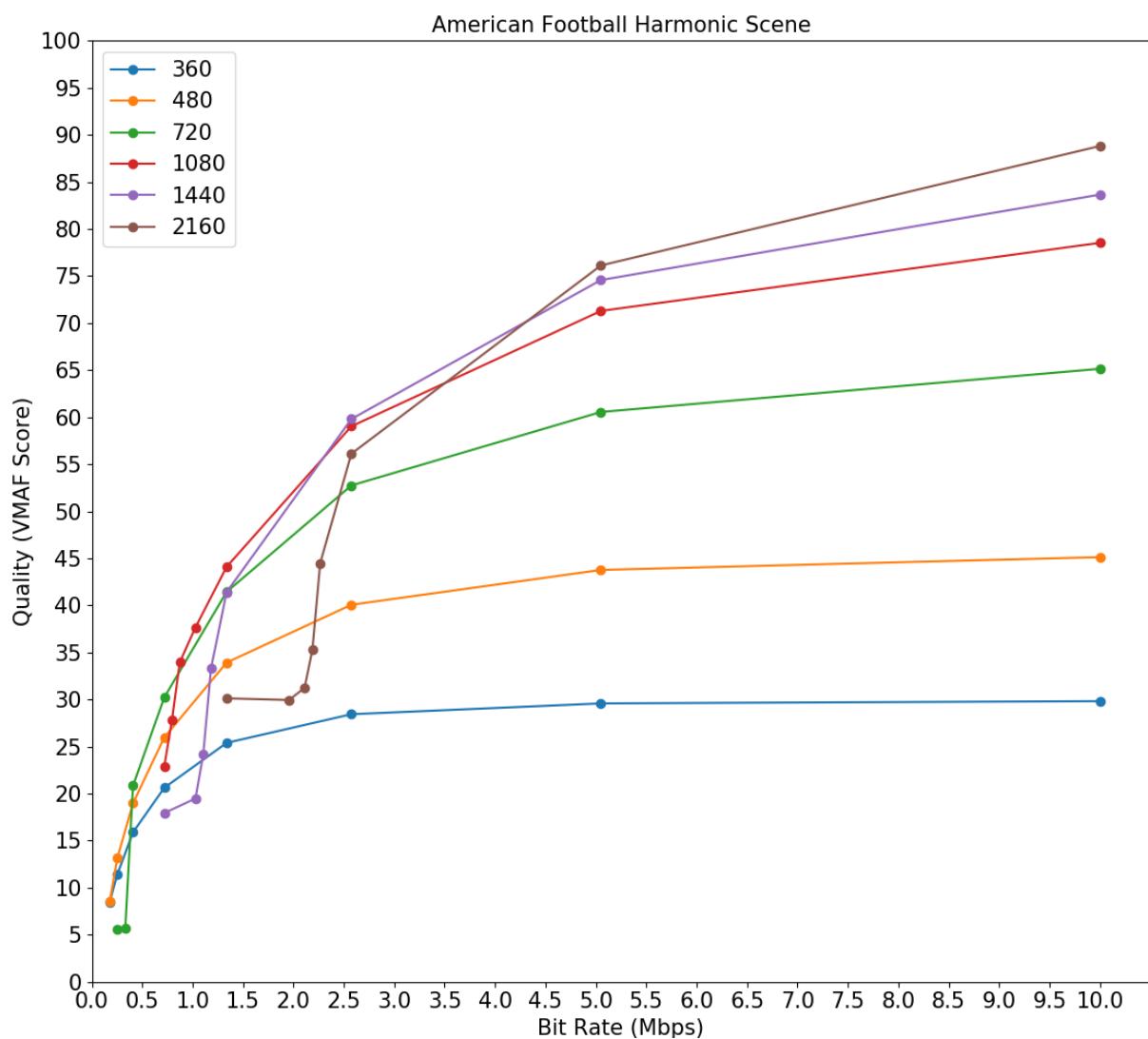


Figure 3.20: Plot for Rate-Quality curves for all resolutions for American Football Scene

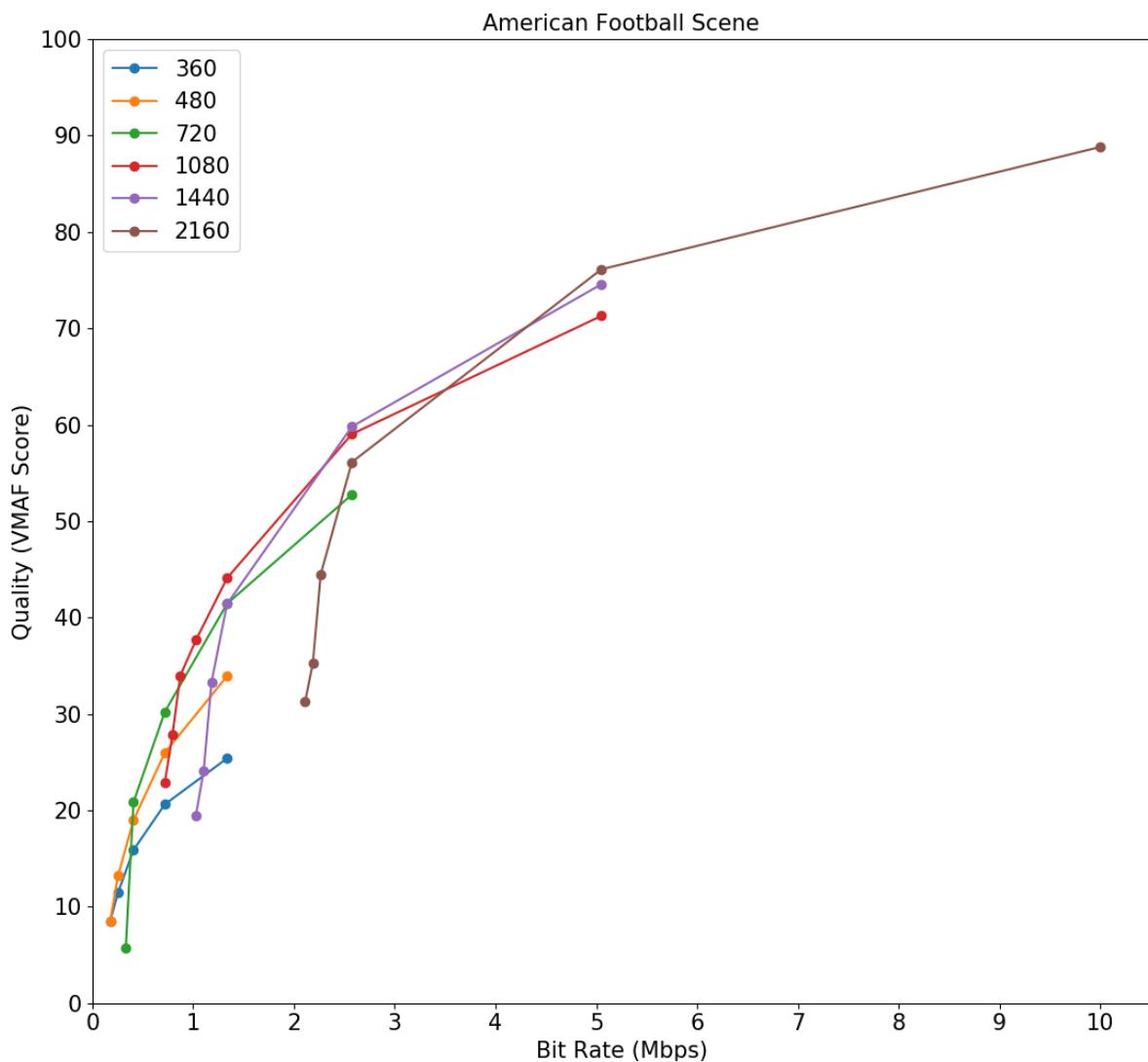


Figure 3.21: Plot for Rate-Quality curves for all resolutions for American Football Scene with only "Selected" Bitrate Points

Chapter 4

Analysis & Evaluation

Measurement results / analysis / discussion: 1/3

- ▷ whatever you have done, you must comment it, compare it to other systems, evaluate it, using e.g. subjective tests
- ▷ usually, adequate graphs help to show the benefits of your approach
- ▷ caution: each result/graph must be discussed! what's the reason for this peak or why have you observed this effect

Chapter 5

Conclusion

Conclusion: 1 page

- ▷ summarize again what your thesis did, but now emphasize more the results, and comparisons
- ▷ write conclusions that can be drawn from the results found and the discussion presented in the paper
- ▷ future work (be very brief, explain what, but not much how)

Bibliography

- [Aar+15] Anne Aaron et al. “Per-title encode optimization”. In: *The Netflix Techblog* (2015).
- [Año+15] Javier Añorga et al. “YouTube’s DASH implementation analysis”. In: *19th International Conference on Circuits, Systems, Communications and Computers (CSCC)*. 2015, pp. 61–66.
- [Blo17] Netflix Technology Blog. *High Quality Video Encoding at Scale*. 2017. URL: <https://medium.com/netflix-techblog/high-quality-video-encoding-at-scale-d159db052746>.
- [Blo18] Netflix Technology Blog. *VMAF: The Journey Continues*. 2018. URL: <https://medium.com/netflix-techblog/vmaf-the-journey-continues-44b51ee9ed12>.
- [BR96] John S Boreczky and Lawrence A Rowe. “Comparison of video shot boundary detection techniques”. In: *Journal of Electronic Imaging* 5.2 (1996), pp. 122–129.
- [Bre19] Breakthrough. *Breakthrough/PySceneDetect*. 2019. URL: <https://github.com/Breakthrough/PySceneDetect>.
- [Can] Canonical. URL: <http://manpages.ubuntu.com/manpages/bionic/man1/shotdetect.1.html>.
- [Cor77] George Corliss. “Which root does the bisection algorithm find?” In: *Siam Review* 19.2 (1977), pp. 325–327.
- [CV95] Corinna Cortes and Vladimir Vapnik. “Support-vector networks”. In: *Machine learning* 20.3 (1995), pp. 273–297.
- [Ffma] URL: <https://ffmpeg.org/ffmpeg-filters.html#blackframe>.
- [Ffmb] *FFmpeg*. URL: <https://ffmpeg.org/>.
- [GSW00] David Gibson, Michael Spann, and Sandra Woolley. “Comparative study of compression methodologies for digital angiogram video”. In: *ISPACS 2000* (2000), pp. 5–8.

Bibliography

- [ISO14] ISO ISO. "ISO/IEC 23009-1: 2014: Information technology—Dynamic adaptive streaming over HTTP (DASH)—Part 1: Media presentation description and segment formats". In: *Geneva, Switzerland: International Organization for Standardization* (2014).
- [Joh] Johmathe. *johmathe/shotdetect*. URL: <https://github.com/johmathe/shotdetect/blob/master/description-pak>.
- [Kat18] Ioannis Katsavounidis. *Dynamic optimizer - a perceptual video encoding optimization framework*. 2018. URL: <https://medium.com/netflix-techblog/dynamic-optimizer-a-perceptual-video-encoding-optimization-framework-e19f1e3a277f>.
- [KG18] Ioannis Katsavounidis and Liwei Guo. "Video codec comparison using the dynamic optimizer framework". In: *Applications of Digital Image Processing XLI*. Vol. 10752. International Society for Optics and Photonics. 2018, 107520Q.
- [KZS15] Dilip Kumar Krishnappa, Michael Zink, and Ramesh K Sitaraman. "Optimizing the video transcoding workflow in content delivery networks". In: *Proceedings of the 6th ACM Multimedia Systems Conference*. ACM. 2015, pp. 37–48.
- [LMN12] Songnan Li, Lin Ma, and King Ng Ngan. "Full-reference video quality assessment by decoupling detail losses and additive impairments". In: *IEEE Transactions on Circuits and Systems for Video Technology* 22.7 (2012), pp. 1100–1112.
- [Man+18] Megha Manohara et al. *Optimized shot-based encodes: Now Streaming!* 2018. URL: <https://medium.com/netflix-techblog/optimized-shot-based-encodes-now-streaming-4b9464204830>.
- [Mat] *Scene Change Detection*. URL: <https://de.mathworks.com/help/vision/examples/scene-change-detection.html>.
- [Moo+19] Ben Moore et al. *The Best Video Streaming Services for 2019*. 2019. URL: <https://www.pcmag.com/roundup/336650/the-best-video-streaming-services>.
- [MRG99] Nader Mohsenian, Rajesh Rajagopalan, and Cesar A Gonzales. "Single-pass constant-and variable-bit-rate MPEG-2 video compression". In: *IBM Journal of Research and Development* 43.4 (1999), pp. 489–509.
- [Neta] NETFLIX. <https://www.netflix.com/>.
- [Netb] Netflix. *Netflix/vmaf*. URL: https://github.com/Netflix/vmaf/blob/master/resource/doc/VMAF_Python_library.md.

- [OJ] Martin Ombura Jr. *How YouTube handles streaming 4,000,000,000+ daily videos without a hitch.* <https://medium.com/@martinomburajr/how-youtube-handles-streaming-4-000-000-000-daily-videos-without-a-hitch-8542741e957a>.
- [Ope] *OpenCV.* 2019. URL: <https://opencv.org/>.
- [OR02] James Oly and Daniel A Reed. “Markov model prediction of I/O requests for scientific applications”. In: *Proceedings of the 16th international conference on Supercomputing*. ACM. 2002, pp. 147–155.
- [Pev] *the Standard for Perceptual Evaluation of Video Quality.* URL: <http://www.pevq.com/>.
- [Pri] *PrimeVideo.* <https://www.primevideo.com/>.
- [Psn] *Peak signal-to-noise ratio.* 2019. URL: https://en.wikipedia.org/wiki/Peak_signal-to-noise_ratio.
- [Pyt] *Welcome to Python.org.* URL: <https://www.python.org/>.
- [Ras17] Reza Rassool. “VMAF reproducibility: Validating a perceptual practical video quality metric”. In: *2017 IEEE International Symposium on Broadband Multimedia Systems and Broadcasting (BMSB)*. IEEE. 2017, pp. 1–2.
- [Sat+19] Shahid Mahmood Satti et al. “Low Complexity” Smart” Per-Scene Video Encoding”. In: *2019 Eleventh International Conference on Quality of Multimedia Experience (QoMEX)*. IEEE. 2019, pp. 1–3.
- [Sto11] Thomas Stockhammer. “Dynamic adaptive streaming over HTTP– standards and design principles”. In: *Proceedings of the second annual ACM conference on Multimedia systems*. ACM. 2011, pp. 133–144.
- [tea] Libav team. *Libav - Open source audio and video processing tools.* <https://www.libav.org/>.
- [VMC17] Van-Huy Vu, Ibrahim Mashal, and Tein-Yaw Chung. “A Novel Bandwidth Estimation Method Based on MACD for DASH.” In: *KSII Transactions on Internet & Information Systems* 11.3 (2017).
- [Wan+04] Zhou Wang et al. “Image quality assessment: from error visibility to structural similarity”. In: *IEEE transactions on image processing* 13.4 (2004), pp. 600–612.
- [Win12] Stefan Winkler. “Analysis of public image and video databases for quality assessment”. In: *IEEE Journal of Selected Topics in Signal Processing* 6.6 (2012), pp. 616–625.

Bibliography

- [WRG99] Peter H Westerink, Rajesh Rajagopalan, and Cesar A Gonzales. “Two-pass MPEG-2 variable-bit-rate encoding”. In: *IBM Journal of Research and Development* 43.4 (1999), pp. 471–488.
- [Yt] *YouTube*. <https://www.youtube.com/>.

List of Figures

1.1	The overview of Dynamic Adaptive Streaming over HTTP technique. [VMC17]	2
1.2	Adaptive Streaming [OJ]	3
1.3	100 randomly sampled sources at 1080p resolution using x264 constant QP (Quantization Parameter) rate control. [Aar+15]	5
2.1	Example encodes showing individual R-D curves and convex hull.[Aar+15] . . .	11
2.2	Schematic for Two-Pass Video Encoding System.[WRG99]	12
2.3	High Level VMAF System Diagram [KG18]	16
3.1	Overview of Theoretical Video Encoding Pipeline	17
3.2	Scene Information in CSV Format of a Random Video Sequence using PySceneDetect	21
3.3	Video Encoding Pipeline for only 720p Resolution	22
3.4	Key Frames of the Hongkong Scene	22
3.5	Plot for Information versus Temporal Information Plot for Hongkong Scene .	23
3.6	Sample VMAF Output[Netb]	24
3.7	Plot for Rate-Quality Curve at 720p resolution for Hongkong scene	25
3.8	Video Encoding Pipeline for all Resolution	25
3.9	Plot for Rate-Quality Curves for Hongkong scene at all resolutions over different Bitrates	26
3.10	Plot for rate-Quality Curves for Hongkong scene at all resolutions over "Se- lected" Bitrates	27
3.11	Video Encoding Pipeline based on Bisection Method Approximation	29
3.12	Key Frames of the Skateboarding Scene in the sequence as 1) top-left, 2)top- right, 3)bottom-left and 4)bottom-right	32
3.13	Plot for Spatial Information versus Temporal Information for the Skateboarding Scene	33
3.14	Plot for Rate-Quality curves for all resolutions for Skateboarding Scene . . .	34
3.15	List of encodes generated using Bisection Algorithm for Skateboarding Scene at 720p	35

List of Figures

3.16 Plot for Rate-Quality curves for all resolutions for Skateboarding Scene with "Selected" Bitrate points	36
3.17 Video Encoding Pipeline based on Bisection Method Approximation with "Additional" Bitrate Threshold	37
3.18 Key Frames for American Football Scene with sequence; 1) top-left, 2) top-right and 3)bottom	39
3.19 Plot for Spatial Information versus Temporal Information for American Football Scene	40
3.20 Plot for Rate-Quality curves for all resolutions for American Football Scene	41
3.21 Plot for Rate-Quality curves for all resolutions for American Football Scene with only "Selected" Bitrate Points	42

List of Tables

1.1	Bitrate-Resolution pairs or Bitrate Ladder [Aar+15]	4
2.1	FFMPEG Command Line Tools	13
2.2	FFMPEG Libraries	14

Appendix A

AppendixExample

Declaration

I declare that the work is entirely my own and was produced with no assistance from third parties.

I certify that the work has not been submitted in the same or any similar form for assessment to any other examining body and all references, direct and indirect, are indicated as such and have been cited accordingly. Ilmenau,

Todo list

- optional
- maximum of 2400 chars; one paragraph i
- maximum of 2400 chars; one paragraph i