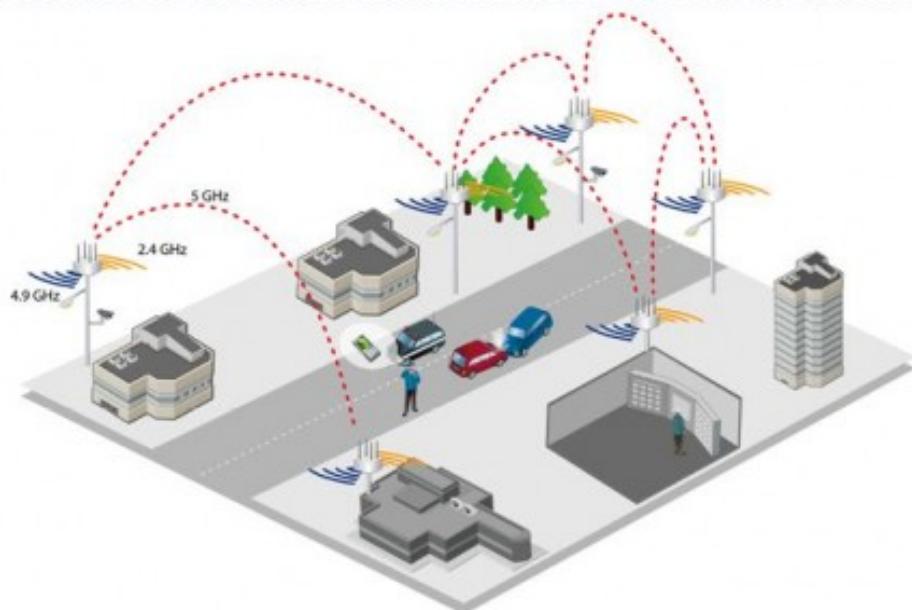


UNIVERSITY OF WATERLOO

FACULTY OF ENGINEERING  
MECHANICAL AND MECHATRONICS ENGINEERING

**VIDEO ENCODING TECHNIQUES FOR  
NETWORKED LOW POWER  
APPLICATIONS**



SELF STUDY REPORT

Prepared by  
SATYA PRAHLADA STHITA PRAJNA KANDARPA  
UW ID 20402024 | USERID *spspkand*  
4B MECHATRONICS ENGINEERING  
31 DECEMBER 2015

41 Pineslope Crescent  
Scarborough, Ontario, Canada  
M1E 4M5

31 December 2015

Professor William Melek,  
Director of Mechatronics Engineering  
Department of Mechanical and Mechatronics Engineering  
University of Waterloo  
Waterloo, Ontario  
N2L 3G1

Dear Sir,

This report, titled "Video encoding techniques for networked low power applications", was prepared as my 4B Work Report for the University of Waterloo. This report is in fulfillment of the course WKRPT 400. The purpose of this report is to evaluate standard video encoding techniques in the context of networked low power video sensor applications and compare their performance with novel encoding techniques developed for distributed sensor networks.

I got exposed to some innovative and novel media processing techniques at a startup I worked at which helped pique my interest in audio visual media processing. This report intends to provide guidance and critical technical evaluation from a software performance perspective to anyone interested in developing low power sensor networks that integrate cameras, audio sensors and mobile communication devices.

The technical analysis conducted by me for this purpose incorporates machine learning techniques to evaluate the standard video encoding technologies to produce optimized encoding parameters that may match the aforementioned specialized distributed video codecs in terms of video transcoding performance. This analysis maybe useful to anyone who would like to avoid the high license costs for the specialized video codecs.

This report was written entirely by me and has not received any previous academic credit at this or any other institution.

Sincerely,  
Satya Prahlada Sthita Prajna Kandarpa  
ID 20402024

# TABLE OF CONTENTS

<b>List of Tables</b>	<b>iv</b>
<b>List of Figures</b>	<b>v</b>
<b>Summary</b>	<b>vi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	1
<b>2 Objectives</b>	<b>3</b>
<b>3 Video Processing And Transport</b>	<b>4</b>
3.1 Anatomy of a video . . . . .	4
3.1.1 Pre-recorded Videos . . . . .	5
3.1.2 Live/Streaming Videos . . . . .	6
3.2 Video Processing . . . . .	7
3.2.1 Encoding and Decoding . . . . .	7
3.2.2 H.264/AVC encoder . . . . .	12
3.2.3 Transcoding . . . . .	14
<b>4 Distributed Video Coding</b>	<b>17</b>
4.1 Overview . . . . .	17
4.1.1 Slepian-Wolf (SW) Theorem for Lossless Coding . . . . .	17
4.1.2 Wyner-Ziv(WZ) Theory . . . . .	18

4.2	PRISM Encoder . . . . .	18
4.2.1	PRISM encoder performance . . . . .	19
4.3	DISCOVER Encoder . . . . .	19
4.3.1	DISCOVER Encoder Performance . . . . .	19
<b>5</b>	<b>Regression Analysis of Youtube Videos</b>	<b>20</b>
5.1	Dataset Characteristics . . . . .	20
5.2	Methodology for Regression based training . . . . .	22
5.3	Evaluation of Software Packages . . . . .	22
5.4	Regression Training using <i>caret</i> . . . . .	22
5.4.1	Data Splitting and Pre-processing . . . . .	22
5.4.2	Tuning and building models . . . . .	22
5.4.3	Resampling and Model Cross-Validation . . . . .	22
5.4.4	Characterizing performance and variable importance .	22
<b>6</b>	<b>Conclusions and Recommendations</b>	<b>24</b>
6.1	Conclusions . . . . .	24
6.2	Recommendations . . . . .	24
<b>References</b>		<b>25</b>

## LIST OF TABLES

3.1 Some common video containers and compatible video coding formats . . . . .	6
---	---

## LIST OF FIGURES

3.1	Matrix representation of a video and its component frames . . . . .	4
3.2	Location of pixel $I_2(0, 2)$ in a video frame . . . . .	5
3.3	Structure of a video codec . . . . .	8
3.4	Digital video produced through compression techniques . . . . .	8
3.5	Demonstration of temporal redundancy[1] . . . . .	9
3.6	Subsampling chromaticity to achieve data reduction[2] . . . . .	9
3.7	Demonstration of spatial redundancy[1] . . . . .	10
3.8	Structure of H.264/AVC encoder[3] . . . . .	13
4.1	Achievable rates of lossless coding by distributed coding of two statistically dependent random signals[4] . . . . .	17
4.2	Architecture of PRISM coding[4] . . . . .	18
5.1	Histogram of videos by codec type in the youtube dataset . . . . .	20
5.2	Jitter plots of framerates separated by codec type . . . . .	21
5.3	Jitter plots of bitrates separated by codec type . . . . .	21
5.4	Transcoding rate [fps] v/s video duration [min] . . . . .	22
5.5	Transcoding rate [fps] v/s video bitrate [Kbps] . . . . .	22
5.6	Transcoding rate [fps] v/s video framerate [fps] . . . . .	23

# SUMMARY

The main purpose of this report is to give broad insight into the current state of video encoding technologies for various applications. The report introduces the burgeoning world of low power video sensor networks and talks about their applications in fields such as crowd, traffic and home surveillance and audio visual media (from an artistic context).

The report then describes the computational capabilities available for low power video sensors by analyzing the technical specifications for one standard low power video camera with wireless capabilities, the Dakota Ultra-Low Power Day/Night Camera. These technical specifications are then used to come up with viable constraints for video processing benchmarks such as time taken to encode a frame, network bandwidth required for continuous transmission, latency and Quality of Experience (QoE). These constraints may also be thought of as targets for the video processing system to be implemented by a low power camera.

The report then introduces standard video processing techniques from a very low level so that the reader may gain insight into the computations and algorithms that power and drive today's digital media driven world. This is done to give the reader enough background knowledge to properly evaluate the project. Specifically, the structural anatomy of a video (on disk, in memory and during transport), is provided. Then, standard video processing techniques are described, namely encoding, decoding, transcoding and network transport. The report dives pretty deep into mathematical and image processing concepts that power video processing, and thus, a basic knowledge of signal processing techniques and linear algebra is assumed. The mathematical theorems that govern data redundancy reduction and computationally efficient algorithm design are explained from a higher level as a lot of the math was beyond an undergrad engineering student's grasp.

A comprehensive analysis of standard video processing techniques including encoding, decoding and transcoding is presented from the perspective of

their applicability to a network distributed video collection system. The various protocols available that enable video transport over a network aren't covered in much detail, however. The aforementioned distributed system has a server that acts as the centralized repository of video streams from each of the video sensors in the network. The process occurs sequentially starting from the capture of raw frames by a physical sensor, to the raw frames being encoded to a bit-stream by a video codec, the transportation of this bit-stream over a network to the centralized server and the final processing task, which involves using multi-view coding techniques to generate a multi-dimensional representation of all the videos.

The standard video codec covered by this report is the ubiquitous H.264/HVC compression standard developed by the ITU-T Video Coding Experts Group (VCEG) together with the ISO/IEC JTC1 Moving Picture Experts Group (MPEG). It is one of the most widely used video standards with applications including Internet streaming, Blu-ray disks and HDTV broadcasts over satellite and cable networks. The advantages offered by this video codec include providing imperceptible quality loss at lower bitrates than other standards and perhaps, the most important one, its ability to integrate well with existing video encoding and transmission infrastructure across a wide range of applications.

The report then presents an overview of the technologies behind the new video encoding technique known as Distributed Video Coding (DVC). The main advantages offered by video codecs that implement DVC include offloading encoding complexity to the decoder, which, for the purposes of this report happens on the centralized server. DVC is very flexible in that it allows user configurable distribution of coding complexity between the encoder and decoder based on application requirements. Two experimental video codecs implemented by researchers, namely the PRISM codec and DISCOVER codec are presented and their performance, characterized by the results of a few research papers, is analyzed based on the objectives defined for the low power video sensor.

The report then analyses a publicly available dataset of the transcoding performance of a few thousand youtube videos. This dataset includes char-

acteristics such as input codec, frame-rate, size and output codec, frame-rate, bitrate and most importantly, memory and cpu time taken to transcode the input video to output video. This dataset is used to train two classes of statistical regression models, namely linear regression and non-linear regression models, to be able to predict the cpu time and memory required for transcoding. The videos in the dataset use 4 different codecs - flv, h264, mpeg4 and vp6. The transcoding was performed using the most widely used open source audio/video processing library, FFmpeg.

The linear regression models trained using the dataset include Partial Least Squares Regression (PLS), Elastic Nets (ENET), Multiple Linear Regression (LM) and Robust Linear Models (RLM). The second class of non-linear regression models covered include k-Nearest Neighbors (kNN), Neural Networks (NN), Model Averaged Neural Networks (avNNet), Multivariate Adaptive Regression Splines (MARS) and finally, Support Vector Machines (SVM). All of the analyses are implemented using the statistical programming language, R using the packages *caret* and *AppliedPredictiveModelling*. The results from all the models are evaluated using two metrics, namely  $R^2$  and Root Mean Squared Error, which are the statistical performance metrics for regression models.

It is to be noted that the transcoding performance in the youtube dataset was measured on a fairly standard server computer with a high amount of processing power. However, the rationale behind training multiple models using this dataset is to be able to predict transcoding times for standard codecs irrespective of the computational power of the machine on which the transcoding operation is being run. This was done since there was no way to measure the computational encoding performance on an actual low power video sensor.

In conclusion, this report agrees on the advantages offered by the DVC based codecs and recommends anyone interested in such application to research them and try to build new products based on these codecs to become future proof.

# 1 | INTRODUCTION

## 1.1 | BACKGROUND

The deployment of high-speed, wired and wireless networks such as 802.16, 802.16a, and 802.11b/g and the explosion of digital camera equipped cellular phones has already provided basic infrastructure for supporting communications in high data-rate wireless video sensor networks. These networks can find their way into many real-time applications needing video-based active monitoring of telemetry data in such diverse indoor and outdoor environments as hospitals, hotels, parking lots, highways, airports, and international borders. Typical video sensor networks are made up of multiple cameras with varying degrees of spatially and temporally overlapping coverage, generating correlated signals that need to be processed, compressed, and exchanged in a loss-prone wireless environment to facilitate real-time decisions. However, the sheer volume of visual data involved, with video signals ranging from a few hundreds of kilobits per second to a few megabits per second and more, poses new and unique challenges. There are numerous challenges to be addressed in order to make the second generation of broadband enabled wireless sensor video networks to take hold.

A broadband network of wireless video sensors is subjected to three principal constraints:

1. Limited processing capabilities and diverse display resolutions due in part to inexpensive device designs and limited battery power. These call for lightweight signal processing and compression algorithms at the individual sensor nodes and an architecture that can adapt to the differing processing capabilities of the encoding and decoding nodes.
2. Limited power/energy budget requiring careful management for maximizing network lifetime, the quality of the acquired data, and the ac-

curacy of the decisions. Communication is often the dominant power-consuming activity. Power management requires efficient compression algorithms that maximize the power utilization per bit communicated and controlled dormancy cycles in inter-sensor communication that preclude frequent intersensor communication. This motivates the need for distributed coding and processing.

3. Information loss that is endemic to the harsh, loss-prone, wireless communication environment. This calls for robust coding algorithms, communication and networking protocols, and architectures that are immune to single points of failure. It is important to proactively build in robustness considerations into the architectural foundation rather than as after-thought bandage fixes.

The technologies that can make this vision a reality are within the reach of the general consumer and before them, entrepreneurs and electronics enthusiasts who would like to drive this revolution. Motivated by the enormous impact these technologies could potentially have, a survey of the general state of video encoding technologies revealed a growing interest among the Research and Development community in the field of Distributed Video Coding.

With the above constraints, the traditional views of video coding and transmission as being confined to a downlink scenario (such as television broadcast or download from a video server) need to be relaxed. In the prevalent video coding architectures such as MPEG-x and H.26x, video encoding is the primary computationally intensive task with the complexity dominated by the motion-search operation. Conventional video decoding, on the other hand, has significantly lower complexity. This skewed, somewhat rigid, complexity compartmentalization conflicts with the heterogeneous processing capability requirements of video sensor networks where the encoding units might be able to do only lightweight processing but the relay or decoding units might be more capable. The prevalent video coding architectures are also built upon the principle of (deterministic) predictive coding from which they derive their compression efficiency.

## 2 | OBJECTIVES

## 3 | VIDEO PROCESSING AND TRANSPORT

### 3.1 | ANATOMY OF A VIDEO

Digital videos are ubiquitous in the era of endless streaming/download/playback services such as Youtube, Netflix and VLC.

Digital video is an ordered sequence of digital images, known as frames, played in succession at a given rate, usually represented as a framerate (frames per second or fps ).

$$I_1 = \begin{bmatrix} 0 & 1 & 2 & 2 & 2 & 2 & 3 & 5 & 7 & 7 \\ 0 & 0 & 1 & 2 & 2 & 3 & 5 & 7 & 7 & 7 \\ 0 & 0 & 0 & 1 & 2 & 3 & 5 & 7 & 7 & 7 \\ 0 & 0 & 0 & 2 & 2 & 3 & 6 & 7 & 7 & 7 \\ 0 & 0 & 0 & 2 & 2 & 3 & 7 & 7 & 7 & 7 \\ 0 & 0 & 0 & 2 & 2 & 3 & 7 & 7 & 7 & 7 \\ 0 & 0 & 0 & 0 & 1 & 2 & 5 & 6 & 7 & 7 \\ 0 & 0 & 0 & 0 & 1 & 2 & 3 & 5 & 6 & 7 \\ 0 & 0 & 0 & 0 & 1 & 2 & 3 & 5 & 6 & 7 \\ 0 & 0 & 0 & 1 & 1 & 1 & 2 & 4 & 5 & 7 \end{bmatrix}, I_2 = \begin{bmatrix} 0 & 1 & 2 & 2 & 2 & 2 & 3 & 5 & 7 & 7 \\ 0 & 0 & 0 & 2 & 2 & 3 & 5 & 7 & 7 & 7 \\ 0 & 0 & 0 & 0 & 2 & 3 & 5 & 7 & 7 & 7 \\ 0 & 0 & 0 & 0 & 0 & 2 & 3 & 6 & 7 & 7 \\ 0 & 0 & 0 & 0 & 0 & 0 & 2 & 3 & 7 & 7 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 5 & 6 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 3 & 5 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 3 & 5 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 2 & 4 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 2 & 3 \end{bmatrix}$$

**Figure 3.1** – Matrix representation of a video and its component frames

A grayscale video, represented by  $V$ , is a sequence of images

$$V = I_1, I_2, \dots, I_n, n = \text{number of frames in the video}$$

, and

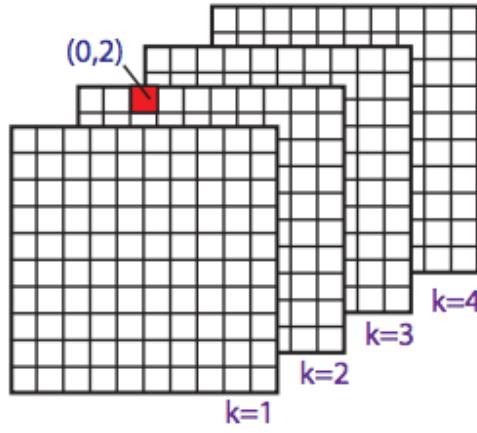
$$I_k \mid k = 1 \dots n$$

is the matrix representation of an image of dimension  $a \times b$ . Please refer to Fig. 3.1 for a visual representation of the matrix images. Each image,  $I_k$  consists of grayscale(brightness or intensity) values from a finite set  $C$  of size  $c$ , where

$$C = \{x \mid x = 0, 1, 2, \dots, N_c - 1\}$$

. A pixel is the basic unit of processing in images. Its location in a video maybe denoted by the 3-D co-ordinates

$$(k, m, n) \text{where } (k, m, n) = (\text{frame number, row number, column number})$$



**Figure 3.2** – Location of pixel  $I_2(0, 2)$  in a video frame

The 3-D co-ordinates may also be represented by  $I_k(m, n) \in C$ . A visual representation of a pixel  $I_2(0, 2)$  is shown in Fig. 3.2[1].

Videos with color information use a similar representation with an additional color component. Pixels in color video frames may be represented by

$$P(I_k, C_t, m, n)$$

where  $C_t$  represents the color component numbered  $t$ . For example, an RGB image can have three possible values for  $t$ , i.e.,  $t \in (1, 2, 3)$ . So,  $P(I_k, C_t, m, n)$  represents the value of the color component  $C_t$  for the pixel with frame co-ordinates  $(m, n)$  and frame  $I_k$  [5].

Videos can be said to have two main representations in digital media - Pre-recorded videos and live/streamable videos. Disk based videos are playable files that may be stored on a personal computing device or a cloud server. These are binary representations of the video data, obtained by compressing raw video frames to achieve optimal spatiotemporal data representation, i.e., reduce redundant data in frames using a combination of motion tracking, Fourier or Discrete Cosine Transforms, Quantization and Variable Length Encoding [2].

### 3.1.1 | PRE-RECORDED VIDEOS

Disk based video file formats may contain uncompressed video footage (RAW format) or encoded video footage (MP4, AVI, etc. formats). Most consumer

focused video file formats consist of the following components:

## || CONTAINER

The container stores the video and/or audio data using separate encoding formats for video and audio. Popular container types include Matroska(MKV), FLV, Ogg, AVI, etc. It is to be noted that container selection constrains the available video encoding formats. The following table lists a few popular containers and their supported encoder formats.

Name	File Extension	Container	Coding Formats
MPEG-4(MP4)	.mp4	MPEG-4 Part 12	H.264
Matroska	.mkv	Matroska	Any
Flash Video(FLV)	.flv	FLV	H.264, VP6

**Table 3.1** – Some common video containers and compatible video coding formats

## || VIDEO CODING(ENCODING) FORMAT

A video coding format (or sometimes video compression format) is a content representation format for storage or transmission of digital video content (such as in a data file or bitstream). Examples of video coding formats include MPEG-2 Part 2, MPEG-4 Part 2, H.264 (MPEG-4 Part 10), HEVC, Theora, Dirac, RealVideo RV40, VP8, and VP9.

### 3.1.2 | LIVE/STREAMING VIDEOS

Streamable videos are defined as multimedia that is constantly received by and presented to an end user while being delivered by the provider. Streaming refers to the delivery method of the video, rather than the video itself, and is an alternative to downloading a full video file. This report deals specifically with live streaming videos, which involves a source media type, a screen recorder in this case, an encoder to digitize the content, and a transport medium, usually one of HTTP, RTSP or RTP.

The main difference between downloadable and streamable videos is speed with which the end user may start watching the video. In case of downloadable videos (files), the user has to wait till the entire file has downloaded to be able to start playing the video. Streamable videos, however, make use of video codecs that are tailored to give the option of beginning playback from any position. They can make this happen by using multiple frame types, frame prediction methods. One frame type in particular, known as a keyframe, enables this resume capability of video from any position because of its decoupled nature from preceding and succeeding frames. Keyframes are implemented differently by video codecs but their essential function stays the same across all codecs. H.264, has a structure that enables interoperability during the video decode process with older video standards. The operational specifics of H.264 are explained in detail in Section [3.2.2](#).

## 3.2 | VIDEO PROCESSING

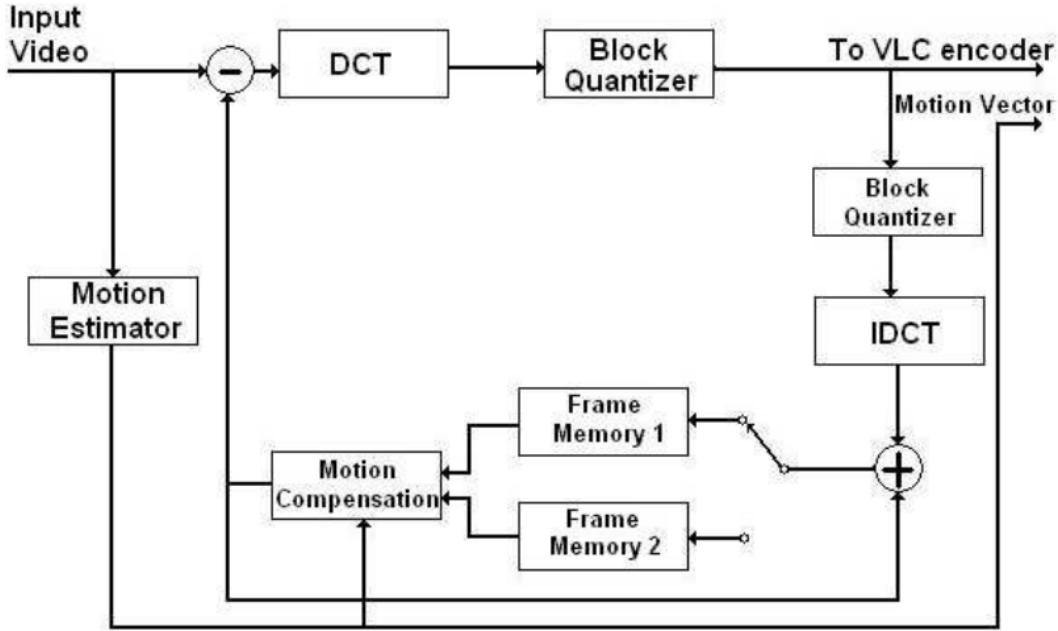
Video processing consists of three main processes - Encoding, Decoding and Transcoding

### 3.2.1 | ENCODING AND DECODING

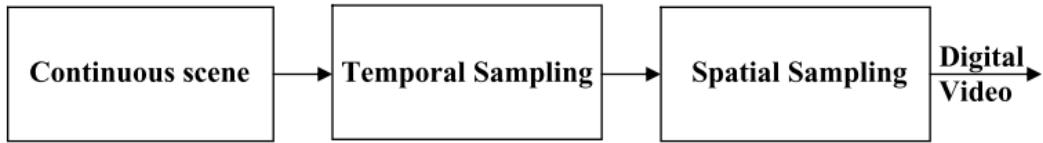
Encoding involves the analysis of uncompressed video files (RAW format) to remove redundant and/or visually indiscernible data and generate a bitstream representation of the video. This bitstream representation may then be used to generate files and/or streamable videos. Decoding involves recovering a playable (streamable) video from the bitstream generated by the encoding process. A *video codec* is a program that can perform both encoding and decoding of a video or bitstream respectively.

The general structure of a video codec is shown in Fig. [3.3](#). It is important to note that the network transport stage of a streamable video involves sending this bitstream representation of a video to an end-user's browser or video playback application like VLC or Quicktime. Decoding occurs in the end-user application via available software or hardware codecs.

This process of data reduction is called video compression or encoding.



**Figure 3.3** – Structure of a video codec



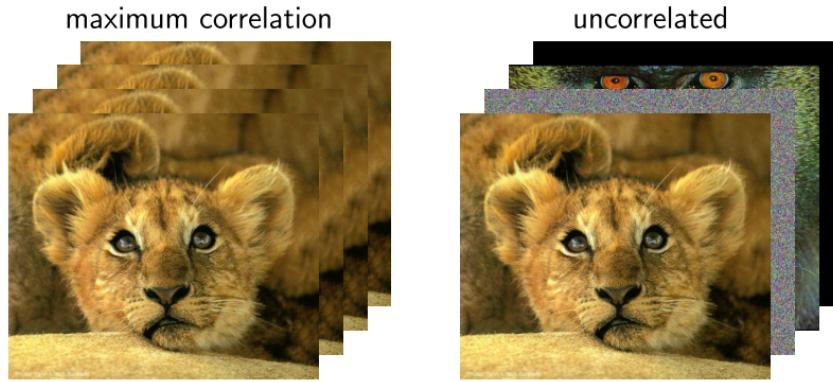
**Figure 3.4** – Digital video produced through compression techniques

When videos are captured by a camera, they are usually stored in an uncompressed format where each frame contains all the original data recorded by the capture device. This process is shown in Fig. 3.4. As evident in the figure, there are two sampling subprocesses, namely Spatial Sampling and Temporal Sampling, that are executed serially to produce a compressed video.

The data present in video frames can have 4 kinds of redundancies [6].

## || TEMPORAL REDUNDANCY

**Temporal Redundancy:** Since the elapsed time between two consecutive frames is generally very short, consecutive frames tend to be very similar in content and thus, contain a lot of data redundancy. The differences between consecutive frames may be expressed by considering the displacements of objects in the frames and encoding this motion's vectors and differences.

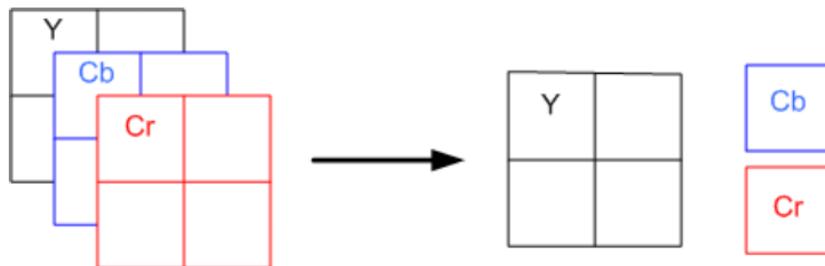


**Figure 3.5** – Demonstration of temporal redundancy[1]

To simplify this procedure, a frame is divided into small fixed (H.261, MPEG-1 encoders) or variable sized blocks(H.263, MPEG-4, H.264 encoders). Motion detection may then be performed by using a statistical measure to determine the best match for a block in a window centered at the block's position in the second frame [2].

## || PSYCHO VISUAL REDUNDANCY

Psycho Visual Redundancy: Since the target audience for 99.99% of all videos is a human recipient, the capabilities of the human visual system (HVS) need to be taken into account before encoding. The HVS is very sensitive to changes in luminance aka intensity compared to changes in chromaticity (color). In fact, the HVS is extremely good at discerning color details based on the intensity levels in an image. This knowledge may be used to selectively subsample the color data in a frame while keeping the intensity data unchanged.



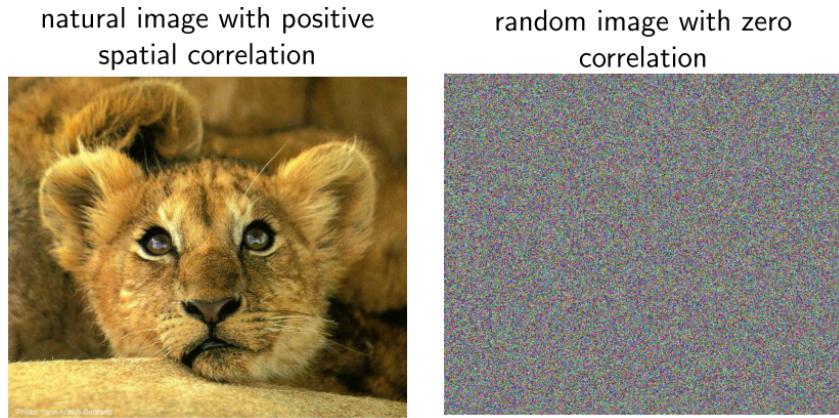
**Figure 3.6** – Subsampling chromaticity to achieve data reduction[2]

A frame maybe divided into macro blocks which are further divided into

three layers with each layer holding one of the color components of the macroblock pixels. YCbCr is the color space used in almost all video coding standards because of its compatibility with the YUV color space used in most displays and televisions and its ability to separate chromaticity from intensity. If the macro block layer for each of Y, Cb and Cr components has 4 8x8 blocks, the Cb and Cr components can each be subsampled to 1 8x8 blocks. The format obtained in this way is referred to as 4:2:0 and the subsampling creates a data reduction of about 50%. A macro block converted to 6 blocks of 8x8 each using 4:2:0 mode is shown in Fig. 3.6

## || SPATIAL REDUNDANCY

In any given frame, a pixel's value is correlated to its neighboring pixel values most of the time. Thus, this value maybe predicted to a certain extent given the values of its neighboring pixels. An example is shown in Fig. 3.7. High correlation means that pixels within a neighborhood have similar colors and zero correlation can mean that pixels in a neighborhood are unrelated in color. Spatial redundancy may be reduced by using transforms such as Discrete Co-



**Figure 3.7** – Demonstration of spatial redundancy[1]

sine Transform (DCT) or Discrete Wavelet Transform (DWT). These values are then quantized and converted to 1-D vectors by reading their values in zig-zag order. These transforms eliminate high frequency pixel values with low energy content.

The quality of the image/frame is directly related to this elimination, which means a trade-off between quality and compression ratio can be achieved based

on constraints imposed by the transmission/playback medium of the video. A video meant to be consumed from disk based file systems can be allowed to retain more data during encoding while videos created for a streaming medium would need to take network bandwidth into consideration to determine an optimal compression ratio.

In case of our application, the video is meant to be recorded, encoded and sent over the network like a streamable video. So, higher compression ratios are desired while still being able to maintain video resolution and clarity rivaling HDMI (720p HD or 1280x720 frame dimensions).

The process followed to achieve an optimum compression ratio is explained in later sections.

## || STATISTICAL REDUNDANCY

The process of reducing statistical redundancy is known as entropy coding. Entropy coding needs to occur after spatial redundancy reduction for optimum compression. The quantized frame data obtained after spatial redundancy reduction is then compressed by Run Length Encoding (RLE) and the resulting values are coded (each unique value gets a unique binary representation) using Huffman encoding.

## || LIVE VIDEO CONSIDERATIONS

Since temporal redundancy reduction makes use of the differences between consecutive frames, this may result in an accumulation of errors even if a frame experiences corruption during network transport.

This scenario can happen if the network transport method used is UDP (User Datagram Protocol), as UDP does not guarantee packet delivery and correct order of packet delivery as is the case with the TCP/IP protocol (used by HTTP). However, UDP is extremely useful for low latency data transmission due to its lack of error correction mechanisms that can guarantee packet delivery without corruption. The operational nature of UDP is commonly referred to as "send and forget".

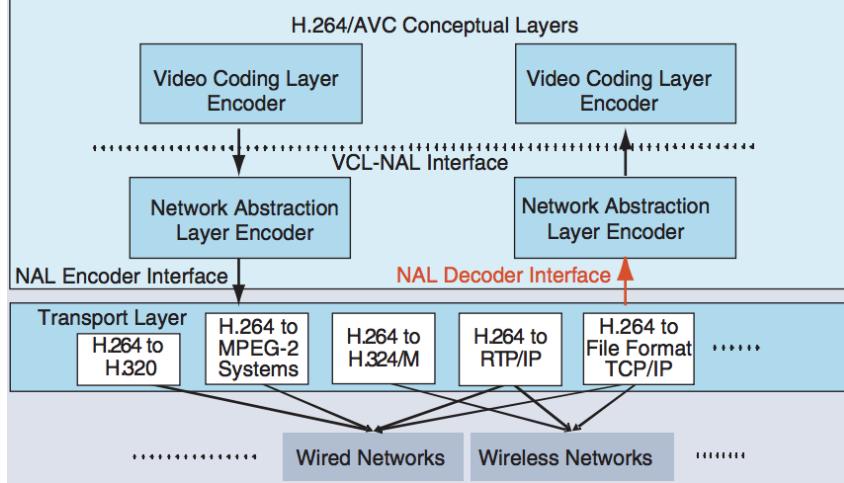
Encoding methods that depend only on the previous frame create a serially accessible frame sequence that requires that the end user download and decode all frames before a particular frame to be able to correctly view the frame. This drawback has been overcome in various ways by different codecs. For example, the MPEG-2 codec uses multiple types of frames, with each type differentiated by the amount of data they hold and dependence on previous or future frames.

Further improvements to video codecs resulted in the development of the H.264/AVC standard described in the following section. The performance and interoperability offered by the H.264 encoder absolutely blew every other encoder out of the water, especially for applications dependent on network transport.

### 3.2.2 | H.264/AVC ENCODER

The H.264 encoder includes a set of improvements to the video coding process that provides enhanced compression performance relative to other encoders like MPEG-2 and VP6. The enhancements provided by H.264 specifically target broadcast, streaming, video telephony and other network friendly video representations. It provides significant improvement in rate distortion efficiency relative to existing standards [7]. All of these features have enabled the H.264 codec to sort of become the de facto standard for video compression for network streaming applications.

A general architecture of the H.264/AVC codec is provided in Fig. 3.8. For efficient transmission in different environments not only coding efficiency is relevant, but also the seamless and easy integration of the coded video into all current and future protocol and network architectures. This includes the public Internet with best effort delivery, as well as wireless networks expected to be a major application for the new video coding standard. The adaptation of the coded video representation or bitstream to different transport networks was typically defined in the systems specification in previous MPEG standards or separate standards like H.320 or H.324. However, only the close integration of network adaptation and video coding can bring the best possible performance of a video communication system.



**Figure 3.8 – Structure of H.264/AVC encoder[3]**

Therefore H.264/AVC consists of two conceptual layers. The video coding layer (VCL) defines the efficient representation of the video, and the network adaptation layer (NAL) converts the VCL representation into a suitable format for specific transport layers or storage media. For circuit-switched transport like H.320, H.324M or MPEG-2, the NAL delivers the coded video as an ordered stream of bytes containing start codes such that these transport layers and the decoder can robustly and simply identify the structure of the bitstream. For packet switched networks like RTP/IP or TCP/IP, the NAL delivers the coded video in packets without these start codes [3].

The following features describe the important features of the H.264 codec that make it a better choice over previous coding standards:

## || INTRA PREDICTION

Intra prediction means that the samples of a macroblock in a frame (slice, in case of H.264) are predicted by using information of already transmitted macroblocks of the same frame. It is to be noted that each image is divided up into smaller packets (NALs) which can be read into macroblocks. H.264 uses varying modes for Intra Frame Prediction depending upon the rates of change of luminance and chromaticity in the image.

## || MOTION COMPENSATED PREDICTION

This is a form of inter-frame (image) prediction. In this case, the macroblocks of an image can be predicted from already transmitted macroblocks of previous reference images. H.264 differs from previous standards (specifically, MPEG) in that it can use several preceding reference images for motion compensation prediction. For this purpose, an additional picture reference parameter has to be transmitted along with the standard motion displacement vectors usually needed for motion compensation prediction as described in Section [3.2.1.1](#).

## || BLOCK TRANSFORM CODING

Former standards such as MPEG-1 and MPEG-2 used a Discrete Cosine Transform (DCT) with block size 8x8 for the purpose of transform coding. H.264 mainly uses 4x4 block sizes while switching to 2x2 blocks in special cases. It also uses 3 different kinds of applied integer transforms instead of a DCT. The first transform type of size 4x4 is applied to all samples of luminance and chromaticity components regardless of whether motion compensation prediction or intra prediction was applied. The other two types of transforms are Haddard transforms of sizes 4x4 and 2x2 respectively.

Compared to the DCT, the applied integer transforms used in H.264 have only integers between -2 and 2 in their transform matrix. This allows computing the transform and inverse transform in 16-bit arithmetic using only low complexity shift, add and subtract operations [3].

## || ENTROPY CODING SCHEMES

Entropy coding is used to reduce statistical redundancy, i.e., use lower number of bits to represent values that occur with high frequencies and a high number of bits to represent values that occur with low frequencies. This reduces the amount of data needed to represent the overall data required to make up the data.

### 3.2.3 | TRANSCODING

Video transcoding refers to the process of data exchange between heterogeneous multimedia networks to reduce the complexity and transmission time by

avoiding total decoding and re-encoding of a video bitstream. Despite the fact that a video stream is generated by eliminating all redundancies, many network channels may not have the necessary capabilities to handle these streams. This restriction may be overcome by reducing the video data size through a change in video format. In terms of video properties, this change can be affected by changing bits per pixel, pixels per frame (pixel density reduction), frames per second, video content or coding standard [2].

Video transcoding for real-time applications on raw video data is extremely time consuming because of the motion estimation and data transformation operations. Acceptable transcoding performance for real time operations can be achieved however, if the conversion of video formats is performed on compressed data rather than raw data. A few effective compressed data video transcoding techniques include:

- Bitrate transcoding
- Spatial transcoding
- Temporal transcoding
- Standard transcoding

A description of these techniques was deemed to be beyond the scope of this report, but more information may be found in the paper cited here [2].

In summary, it can be said that transcoding is something of an art form whereby one must balance dozens of requirements, formats, parameters and more. General video transcoding best practices are presented as follows:

- Always encode for a specific quality rather than relying on bitrates. With bandwidth availability increasing across the board there is no need for using a target bitrate unless a specific limited device is being targeted (applicable to StreamBox) or the quality required is unrealistic within bitrate constraints (in which case quality expectations have to be lowered)

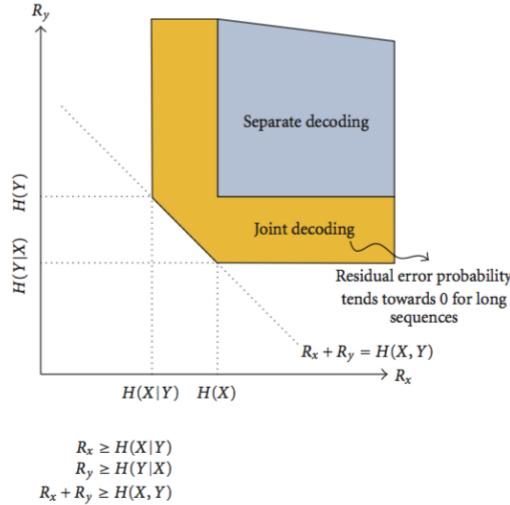
- Avoid upscaling video dimensions from the original dimensions as this only blurs the video. In general, video players do automatic upscaling to make the video fit device screen, so it isn't necessary to do this during transcoding.
- Using Free and Open Source software like FFmpeg and libx264[8] is highly recommended. These libraries have a community of video experts offering help and are very friendly. A lot of money may be saved by not using licensed and proprietary tools while still being able to provide insanely good video quality.[9]

## 4 | DISTRIBUTED VIDEO CODING

### 4.1 | OVERVIEW

Distributed Video Coding (DVC) is ideally suitable to fulfill the above demand. DVC proposed a dramatic structural change to video coding by shifting the majority of complexity conventionally residing in the encoder towards the decoder by implementing distributed source coding concepts. The major task of exploiting the source redundancies to achieve the video compression is accordingly placed in the decoder. The DVC encoder thus performs a computationally very inexpensive operation enabling a significantly low cost implementation of the signal processor in DVC based video cameras[10].

#### 4.1.1 | SLEPIAN-WOLF (SW) THEOREM FOR LOSSLESS CODING



**Figure 4.1** – Achievable rates of lossless coding by distributed coding of two statistically dependent random signals[4]

The SW theorem establishes some lower bounds on the achievable rates for the lossless coding of two or more correlated sources. More specifically, let us

consider two statistically dependent random signals  $X$  and  $Y$ . In conventional coding, the two signals are jointly encoded and it is well known that the lower bound for the rate is given by the joint entropy  $H(X, Y)$ . Conversely, with distributed coding, these two signals are independently encoded but jointly decoded. In this case, the SW theorem proves that the minimum rate is still  $H(X, Y)$  with a residual error probability which tends towards 0 for long sequences. Figure 4.1 illustrates the achievable rate region. In other words, SW coding allows the same coding efficiency to be asymptotically attained. However, in practice, finite block lengths have to be used. In this case, SW coding entails a coding efficiency loss compared to lossless source coding, and the loss can be sizeable depending on the block length and the source statistics.

#### 4.1.2 | WYNER-ZIV(WZ) THEORY

Wyner and Ziv extended the Slepian-Wolf theorem by characterizing the achievable rate- distortion region for lossy coding with Side Information (SI). More specifically, WZ showed that there is no rate loss with respect to joint encoding and decoding of the two sources, under the assumptions that the sources are jointly Gaussian and an MSE distortion measure is used. This result has been shown to remain valid as long as the innovation between  $X$  and  $Y$  is Gaussian.

## 4.2 | PRISM ENCODER

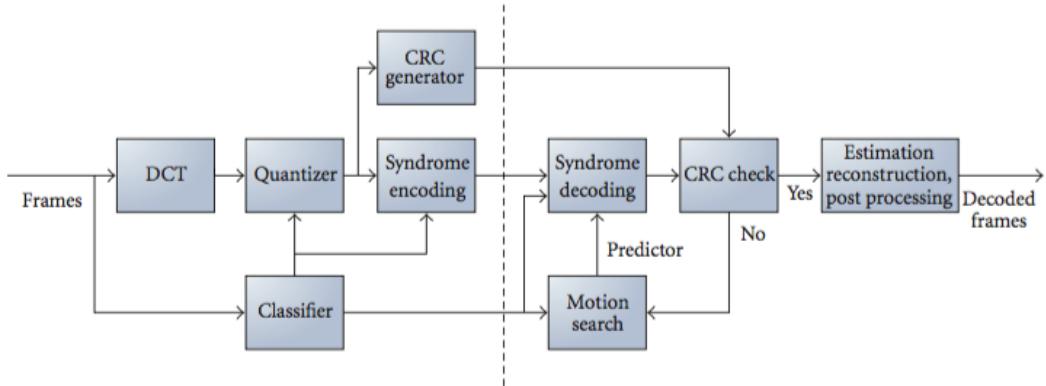


FIGURE 2: PRISM architecture.

**Figure 4.2** – Architecture of PRISM coding[4]

PRISM, which stands for Power-efficient, Robust, hIgh compression Syndrome-based Multimedia coding, is one of the early practical implementations of DVC. This architecture is shown in Figure 4.2. More specifically, each frame is split into  $8 \times 8$  blocks which are DCT transformed. Concurrently, a zero-motion block difference is used to estimate their temporal correlation level. This information is used to classify blocks into 16 encoding classes. One class corresponds to blocks with very low correlation which are encoded using conventional Intra- coding. Another class is made of blocks which have very high correlation and are merely signaled as skipped. Finally, the remaining blocks are encoded based on distributed coding principles. More precisely, syndrome bits are computed from the least significant bits of the transform coefficients, where the number of least significant bits depends on the estimated correlation level.

The lower part of the least significant bit planes is entropy coded with a (run, depth, path, last) 4-tuple alphabet. The upper part of the least significant bit planes is coded using a coset channel code. For this purpose, a BCH code is used, as it performs well even with small block-lengths. Conversely, the most significant bits are assumed to be inferred from the block predictor or Side Information (SI). In parallel, a 16-bit Cyclic Redundancy Check (CRC) is also computed. At the decoder, the syndrome bits are then used to correct predictors, which are generated using different motion vectors. The CRC is used to confirm whether the decoding is successful.

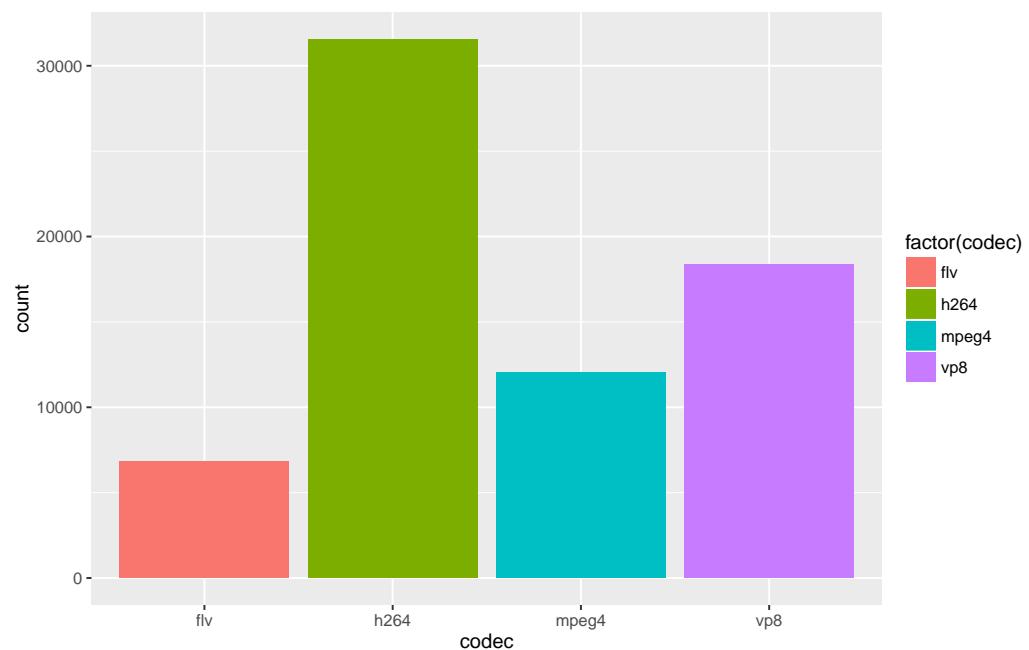
#### 4.2.1 | PRISM ENCODER PERFORMANCE

### 4.3 | DISCOVER ENCODER

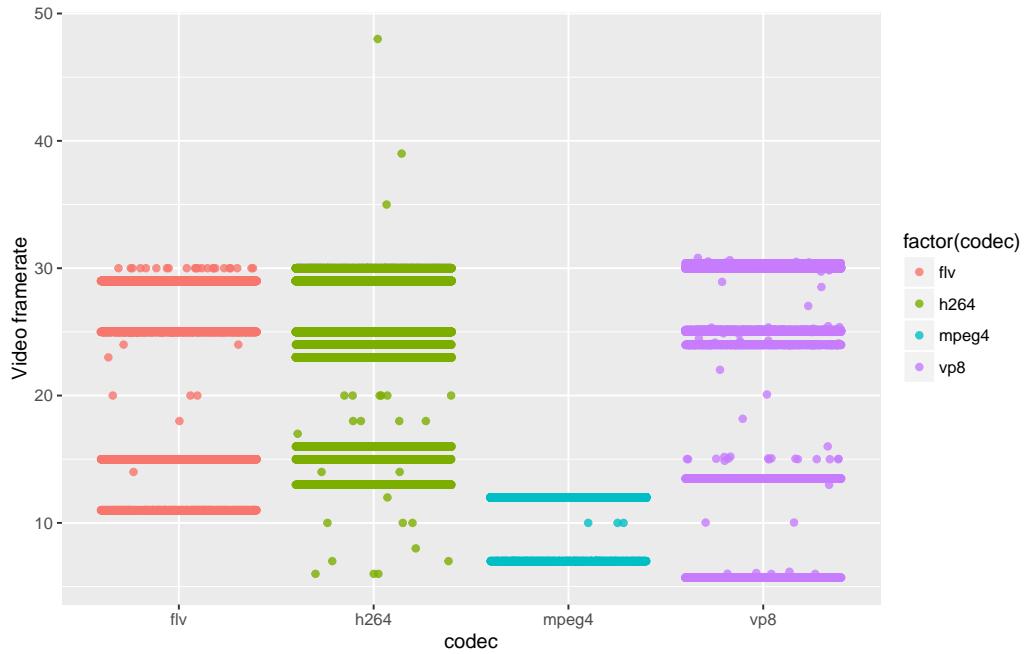
#### 4.3.1 | DISCOVER ENCODER PERFORMANCE

# 5 | REGRESSION ANALYSIS OF YOUTUBE VIDEOS

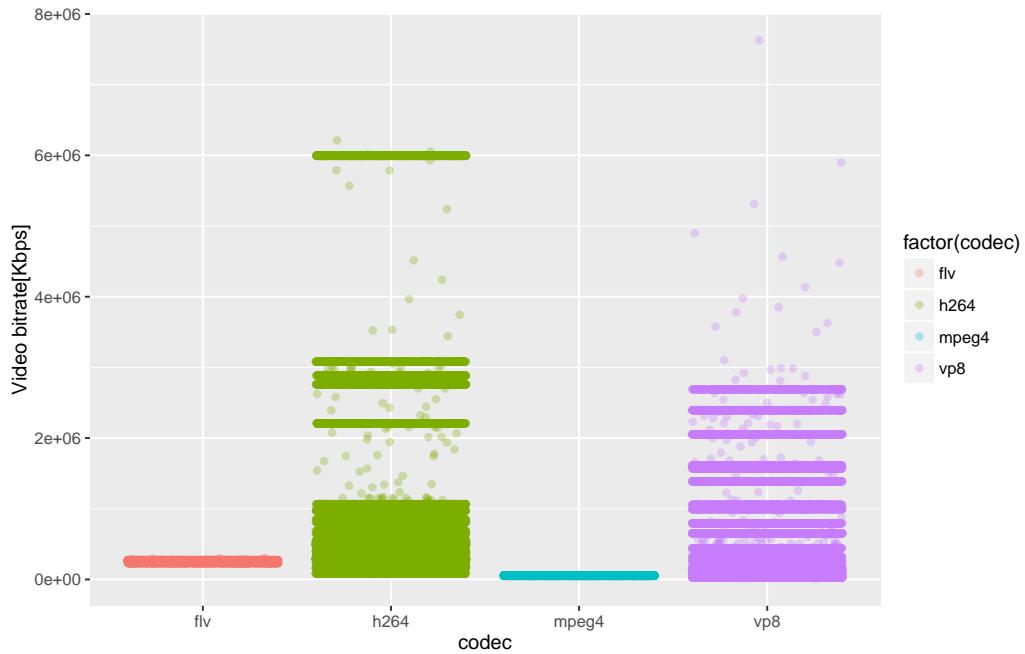
## 5.1 | DATASET CHARACTERISTICS



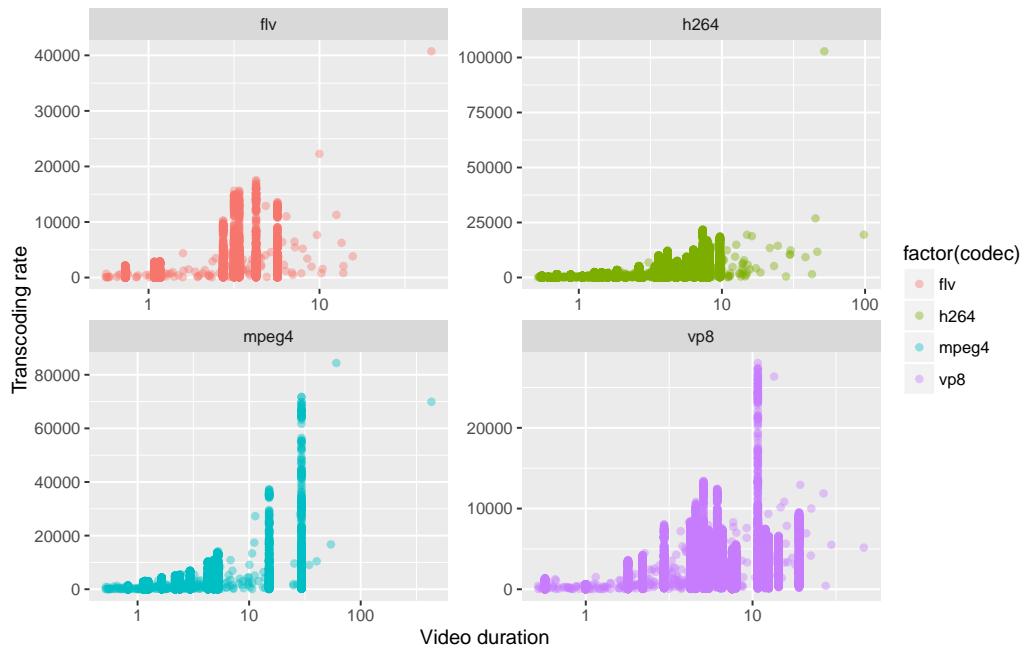
**Figure 5.1** – Histogram of videos by codec type in the youtube dataset



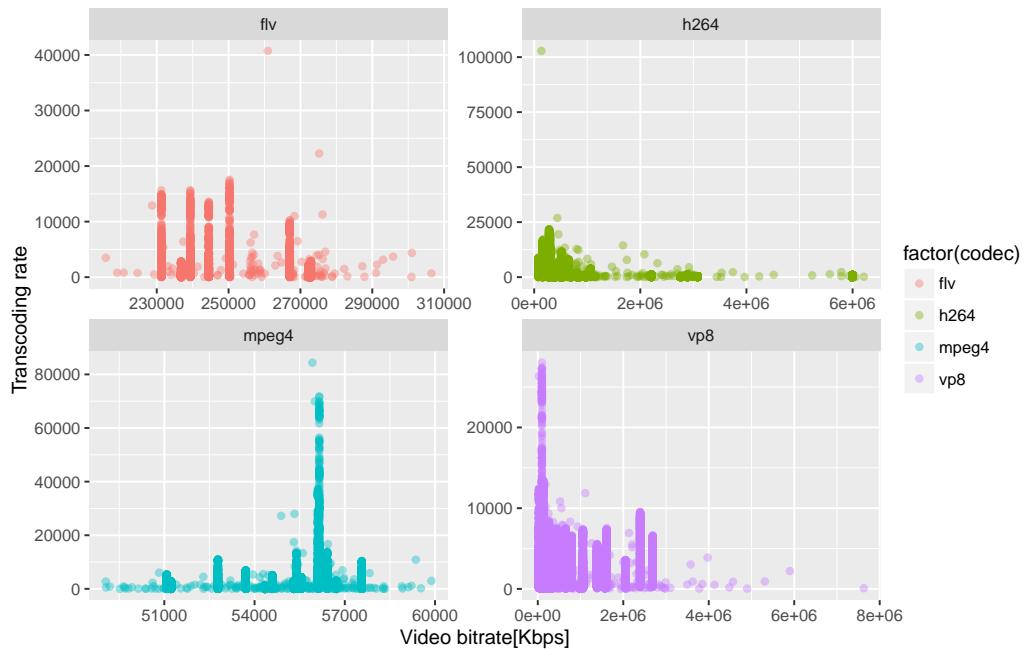
**Figure 5.2** – Jitter plots of framerates separated by codec type



**Figure 5.3** – Jitter plots of bitrates separated by codec type



**Figure 5.4 – Transcoding rate [fps] v/s video duration [min]**



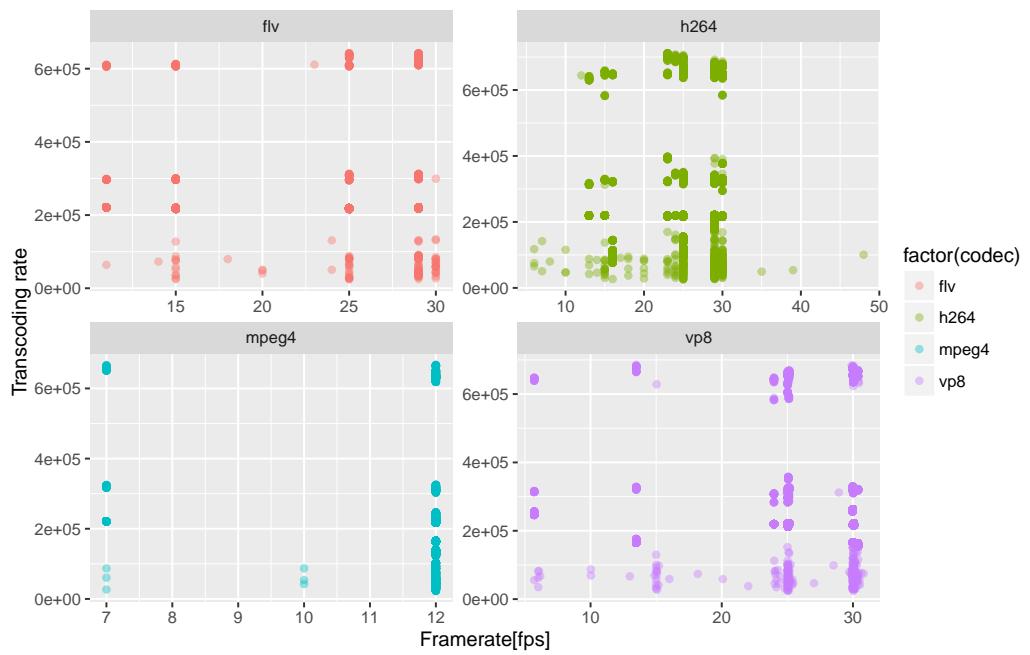
**Figure 5.5 – Transcoding rate [fps] v/s video bitrate [Kbps]**

## 5.2 | METHODOLOGY FOR REGRESSION BASED TRAIN-

ING

## 5.3 | EVALUATION OF SOFTWARE PACKAGES

## 5.4 | REGRESSION TRAINING USING



**Figure 5.6 – Transcoding rate [fps] v/s video framerate [fps]**

## 6 | CONCLUSIONS AND RECOMMENDATIONS

### 6.1 | CONCLUSIONS

### 6.2 | RECOMMENDATIONS

## REFERENCES

- [Kundur] [1] D. Kundur, “Introduction to Video Processing.” [Online]. Available: [http://www.comm.utoronto.ca/~dkundur/course\\_info/real-time-DSP/notes/13\\_Kundur\\_Intro\\_Video\\_Edge\\_Detection.pdf](http://www.comm.utoronto.ca/~dkundur/course_info/real-time-DSP/notes/13_Kundur_Intro_Video_Edge_Detection.pdf)
- [Choupani] [2] R. Choupani, S. Wong, and M. Tolun, “Video coding and transcoding: A review,” 2007.
- [Ostermann2004] [3] “Video coding with H.264/AVC: Tools, performance, and complexity,” *IEEE Circuits and Systems Magazine*, vol. 4, no. 1, pp. 7–28, 2004.
- [Dufaux2009] [4] F. Dufaux, W. Gao, S. Tubaro, and A. Vetro, “Distributed Video Coding: Trends and Perspectives,” *EURASIP Journal on Image and Video Processing*, vol. 2009, no. 1, pp. 1–13, apr 2009. [Online]. Available: <http://jivp.eurasipjournals.com/content/2009/1/508167>
- [Lefevre2003] [5] S. Lefèvre, J. Holler, and N. Vincent, “A review of real-time segmentation of uncompressed video sequences for content-based search and retrieval,” *Real-Time Imaging*, vol. 9, no. 1, pp. 73–98, 2003.
- [Saggi2010] [6] A. Saggi, “A framework for multimedia playback and analysis of MPEG-2 videos with FFmpeg,” Ph.D. dissertation, 2010.
- [Wiegand2003] [7] T. Wiegand, “Overview of the H. 264/AVC video coding standard,” *... and Systems for Video ...*, vol. 13, no. 7, pp. 560 –576, 2003. [Online]. Available: [http://ieeexplore.ieee.org/ielx5/76/27384/01218189.pdf?tp=&arnumber=1218189&isnumber=27384\\$delimiter\\$026E30F\\$nhttp://ieeexplore.ieee.org/xpls/abs{\\_}all.jsp?arnumber=1218189](http://ieeexplore.ieee.org/ielx5/76/27384/01218189.pdf?tp=&arnumber=1218189&isnumber=27384$delimiter$026E30F$nhttp://ieeexplore.ieee.org/xpls/abs{_}all.jsp?arnumber=1218189)
- [Video68:online] [8] “Videolan - x264, the best h.264/avc encoder,” <http://www.videolan.org/developers/x264.html>, (Visited on 01/07/2016).
- [Trans44:online] [9] “Transcoding best practices - jwplayer,” <http://www.jwplayer.com/blog/transcoding-best-practices/>, (Visited on 01/07/2016).

- Weerakkody2007 [10] W. Weerakkody, W. A. C. Fernando, and a. B. B. Adikari, “Unidirectional Distributed Video Coding for Low Cost Video Encoding,” *IEEE Transactions on Consumer Electronics*, vol. 53, no. 2, pp. 788–795, 2007. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4266974>
- FFmpeg90:online? [11] “Ffmpeg devices documentation,” <https://www.ffmpeg.org/ffmpeg-devices.html#Options-1>, (Visited on 01/04/2016).
- ffmpeg:online? [12] “Ffmpeg - open source audio video processing library,” <https://www.ffmpeg.org/>, (Visited on 01/04/2016).
- Chandra2012? [13] S. Chandra, J. T. Biehl, J. Boreczky, S. Carter, and L. a. Rowe, “Understanding screen contents for building a high performance, real time screen sharing system,” *Proceedings of the 20th ACM international conference on Multimedia - MM '12*, p. 389, 2012. [Online]. Available: <http://dl.acm.org/citation.cfm?doid=2393347.2393404>
- Maint44:online? [14] “Maintenance mode,” <http://westseattleblog.com/category/seattle-police-surveillance-cameras/>, (Visited on 01/09/2016).
- Dakot10:online? [15] “Dakota security — ultra-low power day/night camera,” <http://dakotasecurity.com/ultra-low-power-daynight-camera/>, (Visited on 01/09/2016).
- video11:online? [16] “videos.cctvcamerapros.com/pdf/zavio/f312a-wireless-ip-camera.pdf,” <http://videos.cctvcamerapros.com/pdf/Zavio/F312A-wireless-ip-camera.pdf>, (Visited on 01/09/2016).