



# Office of Graduate Studies

## Dissertation / Thesis Approval Form

This form is for use by all doctoral and master's students with a dissertation/thesis requirement. Please print clearly as the library will bind a copy of this form with each copy of the dissertation/thesis. All doctoral dissertations must conform to university format requirements, which is the responsibility of the student and supervising professor. Students should obtain a copy of the Thesis Manual located on the library website.

**Dissertation/Thesis Title:** New Methods for Detecting Frame Deletion in Modern Video

**Author:** Hunter Kippen

**This dissertation/thesis is hereby accepted and approved.**

**Signatures:**

Examining Committee

Chair \_\_\_\_\_

Members \_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_

Academic Advisor \_\_\_\_\_

Department Head \_\_\_\_\_

**New Methods for Detecting Frame Deletion in Modern Video**

A Thesis

Submitted to the Faculty

of

Drexel University

by

Hunter Kippen

in partial fulfillment of the

requirements for the degree

of

Master of Science

May 2019



© Copyright 2019  
Hunter Kippen. All Rights Reserved.

This work is licensed under the terms of the Creative Commons Attribution-ShareAlike 4.0 International license. The license is available at  
<http://creativecommons.org/licenses/by-sa/4.0/>.

## Table of Contents

LIST OF TABLES . . . . .	iii
LIST OF FIGURES . . . . .	iv
1. BACKGROUND . . . . .	1
1.1    Video Encoding . . . . .	1
1.2    Frame Deletion Detection . . . . .	5
1.2.1    Frame Deletion Detection on H.264 . . . . .	7
1.3    Autoregressive Models . . . . .	8
1.4    Support Vector Machines . . . . .	8
2. PROBLEM FORMULATION . . . . .	10
2.1    Video Frame Deletion Detection . . . . .	11
3. PROPOSED APPROACH . . . . .	13
3.1    Prediction Error Sequence Extraction . . . . .	13
3.1.1    Proposed Prediction Error Sequence Extractor for H.264 . . . . .	16
3.2    Proposed Detection Algorithm . . . . .	18
4. EXPERIMENTAL RESULTS . . . . .	22
4.1    Single Camera Model . . . . .	23
4.2    Multiple Camera Models . . . . .	25
4.3    Additional Work . . . . .	27

## List of Tables

4.1 Composition of Frame Deletion Dataset . . . . .	22
---	----

## List of Figures

1.1 Visualization of a Video Prediction Error Residual . . . . .	2
1.2 Example of a GOP Sequence . . . . .	3
1.3 Visualization of Using Multiple Anchor Frames for Prediction in H.264 . . . . .	4
1.4 Visualization of the Effects of Frame Deletion on a Sample MPEG-2 Encoded Video. . . . .	6
1.5 Visualization of the Optimal Hyperplane Learned by an SVM . . . . .	9
2.1 Generalized Approach to Frame Deletion Detection . . . . .	10
3.1 Comparison Between MPEG-2 Detection Methods used on (a) MPEG-2 and (b) H.264 .	15
3.2 Identification of the source macroblock in 3 previous anchor frames . . . . .	16
3.3 Methodology for forming $P[n]$ . . . . .	18
3.4 <i>Top Left</i> - The extracted prediction error sequence for a high motion video. <i>Top Right</i> - The extracted prediction error sequence for a low motion video. <i>Bottom Left</i> - Estimated fingerprint signal for a high motion video. <i>Bottom Right</i> - Estimated fingerprint signal for a low motion video. . . . .	19
3.5 <i>Top</i> - The extracted prediction error sequence for a low motion video with frame deletion. <i>Bottom</i> - The Estimated fingerprint signal for a low motion video with frame deletion. . . . .	20
4.1 ROC Curves for Fingerprint Energy Detector Obtained by Using Different Prediction Error Extraction Methods . . . . .	23
4.2 ROC Curves for Inferred Prediction Error Energy Detector and SVM on One Camera Model . . . . .	24
4.3 ROC Curves for Fingerprint Energy Detector Obtained by Using Different Prediction Error Extraction Methods on Multiple Camera Models . . . . .	25
4.4 Search for Optimal SVM Hyperparameters on (a) Log Scale and (b) Linear Scale . . . . .	26
4.5 ROC Curves for Various Detection Methods on Data from Multiple Camera Models . . . . .	27
4.6 Confusion Matrix for Three Class Classification . . . . .	28
4.7 Confusion Matrix for Three Class Classification using Scaled Data . . . . .	29



## Chapter 1: Background

### 1.1 Video Encoding

Computer storage and processing power is often obtained at a premium. While this is less true today, as reliable high-capacity Hard Disk and Solid-state Drives have become widely available, high definition raw video footage still places a large constraint on these high-capacity drives. They can quickly become filled with only a few minutes of content. For example, a 1080p video has a frame resolution of 1920x1080, meaning a single frame consists of 2,073,600 pixels. Each pixel is represented by three values, one for red, one for blue, and one for green. For standard 8-bit color, a single pixel requires 24 bits of storage. Thus, a single frame requires 49,766,400 bits or about 6 Megabytes of storage. Videos are often played at 24, 30, or even 60 frames per second (fps). A raw 1080p video at 24 fps requires 142.4 Megabytes of storage per second of content. A 1 Terabyte hard drive could store only 122 minutes of content, or a little over 2 hours. If the video was 30 or 60 fps, then the maximum capacity of the hard drive would be even less.

Thus there is a need to be able to compress video files such that they retain the same or similar visual quality while being tens or hundreds of times smaller, depending on the use case. The information stored in video files has a large amount redundancy. This redundancy is primarily along two axes, spatial and temporal. Video compression standards (or codecs) smartly exploit these redundancies to produce highly compressed video files. Two popular video codecs are MPEG-2 [1], used in DVDs and H.264 [2], used for encoding high definition video.

Spatial redundancy in video is similar to that of still images. Most of the information in the visual domain in an image is encoded in low frequency components. One can reduce the amount of information in the high frequency regions of an image and incur very little perceptual loss while reducing the overall size of the image. Joint Photography Experts Group (JPEG) encoding is a widely used compression standard for images that makes use of this particular phenomenon. A simple way to encode video is to treat each video frame as a still image, and encode it using a

JPEG-like process. One of the first video encoding standards, Motion JPEG, did this. However, this method of video compression does not reach the levels of compression when exploiting both spacial and temporal redundancy.

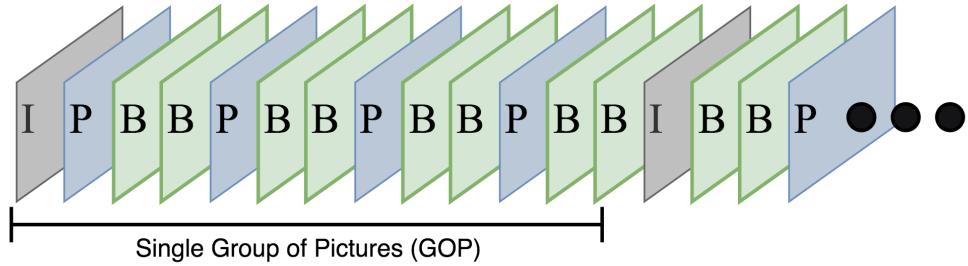
Temporal redundancy is exploiting the fact that in most cases, one video frame does not change too much from the previous frame or frames. Often portions of a video will remain completely static, while certain sections move (think a traffic camera). As such, for a sufficiently small time interval, the current frame can be predicted from the previous frame. Obviously, the current frame is not an exact match for the previous frame. Thus, the difference between the two frames, referred to as the prediction error residual, is stored alongside the previous frame. As shown in figure 1.1, an error residual for a small amount of motion has an extremely low dynamic range. Most values are black, or close to it. It is possible then, to also highly compress the residual frame using a JPEG-like process to get even more space savings without losing much in the way of prediction accuracy.



**Figure 1.1:** Visualization of a Video Prediction Error Residual

MPEG-2 takes both of these concepts a step further. Both Temporal and Spatial redundancy are exploited in the frame prediction process. Successive video frames are first grouped together to minimize prediction error. The first frame in the group is encoded entirely using a JPEG-like process. This frame is called an intra-coded frame or I-frame. Following this I-frame is a sequence of predicted frames. Predicted frames come in two types. Regular predicted frames or P-frames derive their predictions from past frames. Bi-directional frames, or B-frames can mix predictions from both past and future frames. These groupings of frames are called groups of pictures (GOPs).

An example GOP sequence can be seen in Fig. 1.2. In most modern codecs, including MPEG-2,



**Figure 1.2:** Example of a GOP Sequence

these GOP structures do not have to be fixed in size. The GOP length is determined by the distance between subsequent I-frames. Often, it is more efficient to make the GOP smaller during segments of high motion in a video. This is known as variable GOP encoding, versus fixed GOP encoding where the GOP sequence is always the same length.

Predictions can only be derived from I or P-frames, sometimes referred to as anchor frames. In MPEG-2, P-frames can only derive their predictions from a single previous anchor frame, and B-frames can only derive their predictions from a single previous and a single future anchor frame [1]. To decrease the prediction error, predicted frames are partitioned into 16x16 pixel regions called macroblocks. These macroblocks are compared with blocks in the previous anchor frame and/or the next anchor frame in the case of B-frames. A fast search algorithm is performed to find a macroblock sized region with the smallest error when subtracted from the macroblock of interest.

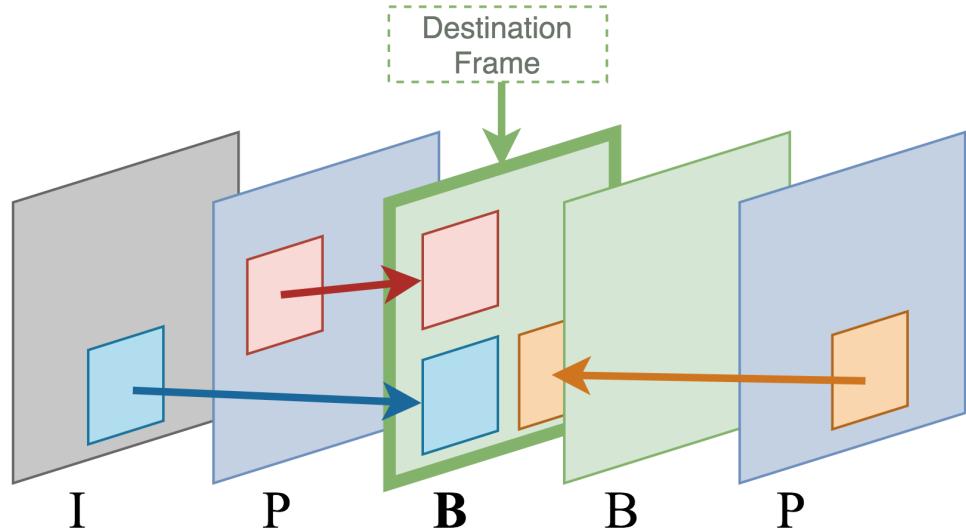
The pixel displacement in both the x and y directions from the centroids of the macroblocks is stored in what is known as a motion vector. Then the error residual between the two macroblocks mapped by the motion vector is also stored. In this way, predicted frames are transmitted as a set of motion vectors and macroblock sized prediction error residuals. Decoding a predicted frame requires using the motion vectors to obtain the pixel values of the source macroblock for each macroblock in the frame. Then the error residual for each macroblock is added back, thus reproducing the original picture. The entire prediction error and motion vector generation process is known as motion compensation and estimation.

The H.264 specification contains many advancements in the motion compensation and estimation

process [2] [3]. Macroblocks no longer have to be 16x16 pixel blocks, and instead can be a variety of shapes including 8x16, 16x8, and 8x8. Single macroblocks can be encoded like an I-frame, instead of being predicted. These I-blocks are useful in scenes with high motion. Also, macroblocks can derive predictions from other portions of their frame. Intra-frame predictions are useful for large, mostly single colored regions like the sky. In this way, only one inter-frame motion vector has to be produced for an entire section of the image.

It is also possible for the encoder to entirely skip a prediction for a given macroblock. Instead, the macroblock is directly copied from a previously decoded frame. These are known as skip macroblocks, or S-blocks.

The most notable improvement from H.264 over MPEG-2 is the ability for the encoder to have a substantially frame buffer. This allows for macroblocks in predicted frames to derive the motion vectors from across multiple different anchor frames. Effectively, a B-frame could be composed of macroblocks sourced from all anchor frames in the encoder's frame buffer. This substantially reduces the prediction error of any given macroblock. Fig. 1.3 shows an example of this. The first B-frame can source macroblocks from any I and P-frame within the GOP that is resident in the encoder's frame buffer.



**Figure 1.3:** Visualization of Using Multiple Anchor Frames for Prediction in H.264

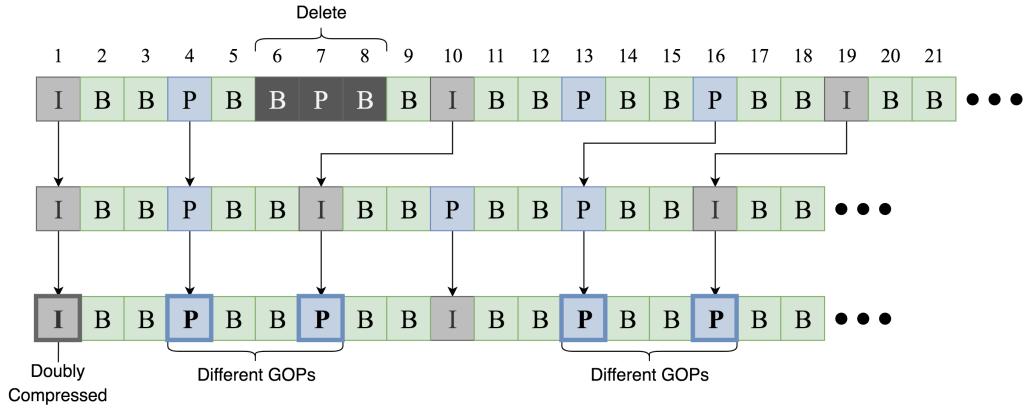
## 1.2 Frame Deletion Detection

Detecting the tampering of a video uses similar general techniques that were originally used for detecting tampering in images. In most cases, a statistical fingerprint is extracted from the media source which contains information about the specific method of tampering an investigator wishes to detect. Frame deletion detection in video was originally found to create two distinct fingerprints in the underlying video [4].

One of these fingerprints was spatial in nature. When a video has had frames removed from it using video editing software, the video had to first be decoded. One could simply store the decoded video without recompression, but then they would run into storage limitations as shown previously. Thus, it is common to recompress a video after removing frames. In this manner, detecting frame deletion in video is related to detecting double compression of video. Due to the JPEG-like process of encoding I-frames, Wang and Farid showed that a similar fingerprint to the one used to detect double JPEG compression [5] [6] is found in the I-frames of doubly compressed MPEG video. Note that this fingerprint occurs even if no frames are removed from the video. Thus, more information is needed to determine whether or not a video has had frames removed from it.

The temporal fingerprint observed by Wang and Farid can make this distinction. The temporal fingerprint occurs due to the structure of the motion compensation and estimation process in MPEG-2 encoding. In each GOP, the statically compressed I-frame is there to prevent excessive buildup of motion estimation errors. The first P-frame in a GOP is thus directly encoded with respect to the I-frame. As each subsequent P-frame in the GOP is encoded based on the previous P-frame, it can be said the these P-frames are indirectly encoded with respect to the initial I-frame. As well, the content of each P-frame was found to be correlated with the I-frame. Wang and Farid explain this phenomenon as the result of the JPEG-like compression process of the I-frame [4]. The compression artifacts propagate through the P-frames resulting in each P-frame being correlated to its previous anchor frame. In addition, the correlation between P-frames in different GOPs is reduced compared to the correlation within a single GOP.

When frames are removed from a video, the new sequence of frames has GOPs that contain



**Figure 1.4:** Visualization of the effects of frame deletion on a sample MPEG-2 encoded video. The top row shows the original encoded sequence, while the bottom rows show the effects of deleting the three highlighted frames.

frames from across multiple GOPs in the original unaltered video. As a result, some frames can change type, and predicted frames can derive their predictions from frames originally in another GOP. This effect can be seen in Fig. 1.4. Frames 6-8 were removed from the original video sequence. Note that frame 13 has become an I-frame when it was originally a P-frame, and frames 10 and 19 have become P-frames when they were originally I-frames. When frames derive their predictions across GOP boundaries, the correlation between the frames is lower, and thus there is an increase in prediction error. Wang and Farid show that for fixed GOP videos, this increase is periodic with respect to the length of a GOP.

As such, Wang and Farid proposed observing the average inter-frame prediction error for each P-frame in a given video. If there are periodic spikes in this P-frame prediction error sequence, then the video has had frames removed from it. While Wang and Farid showed that this fingerprint exists, they did not formulate a detection algorithm for this fingerprint, nor did they test the effects of frame deletion on variable GOP video. They relied on visual inspection of the Discrete Fourier Transform (DFT) of the P-frame prediction error sequence for detecting frame deletion.

Further work was done by Stamm et al. [7] on using the temporal fingerprint proposed by Wang and Farid to detect frame deletion in MPEG-2 video. Stamm et al. proposed an initial model for the

fingerprint itself, and how to separate an estimate of the fingerprint from a given P-frame prediction error sequence. They then formulated two decision rules based on the estimated fingerprint signal. One was for fixed GOP video, and the other was for variable GOP video. The variable GOP decision rule was found to also work for fixed GOP video as well, since it was based on the total energy in the estimated fingerprint signal.

This work still had limitations. First, it is not clear whether this detection method is robust to advancements in video encoding, particularly the ability for H.264 to derive motion vectors across multiple anchor frames. As well, both Stamm et al. and Wang and Farid were limited in the type of data they were able to generate for their experiments. They both used video with carefully controlled processing histories that are not necessarily representative of video that can be found in real world scenarios.

### 1.2.1 Frame Deletion Detection on H.264

Work on frame deletion detection in more modern video codecs, such as H.264 or MPEG-4 so far has been limited. Techniques have been proposed based on Wang and Farid and Stamm et al.'s original solutions for MPEG-2. However, they only cover video without B-Frames and the scenario where the altered video undergoes recompression with a different GOP size after frames have been removed [8]. Another technique based on detecting mismatches in image sensor based traces within a video [9] is able to detect video splicing, where content is stitched together from multiple camera sources. However, this technique may not work on video take only from one camera, as well as video that has undergone motion stabilization, a post-processing technique found in many smartphone cameras.

Of particular note is a technique proposed Vázquez-Padín et al. They found that when H.264 video is reencoded, particularly when GOP sizes do not match between encodings, the number of I-blocks and S-blocks vary from expected values when P-frames are encoded as I-frames and vice-versa. Vázquez-Padín et al. showed that it is possible to use this mismatch to determine whether an H.264 video has been recompressed [10]. This technique was expanded upon by Gironi et al. [11]. As videos that have undergone frame deletion will naturally have frames change type upon

reencoding, the same mismatch in the number of expected I and S-blocks appears in videos with frame deletion as well. Gironi et al. reformulated the proposed fingerprint by Vázquez-Padín et al. to work in this context. However, the method proposed by Gironi et al. is unable to work on video that has a variable GOP size.

### 1.3 Autoregressive Models

An autoregressive (AR) model is a statistical model of a stochastic process. The model seeks to predict future values of the process based on past observations. They are used for forecasting when there is some correlation between the values within the process. The future values are modeled as a linear combination of past observations (called lag variables) and white noise. You only use data from the same process to model the future values, thus the name "autoregressive". Effectively an AR model is a linear regression of the past data in a stochastic process where the target value is a future value of the process.

The model order of an AR model depends on the number of lag variables used in calculating the future values. Thus for a  $K^{th}$  order AR model, the current time series  $X_t$  is defined as

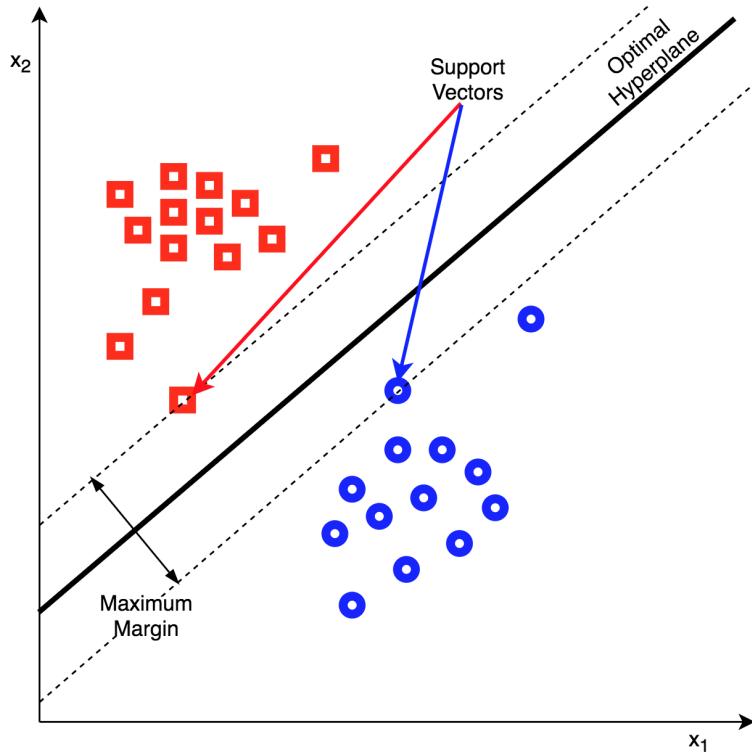
$$X_t = c + \sum_{i=1}^K \varphi_i X_{t-i} + \varepsilon_t \quad (1.1)$$

where  $\varphi_1, \dots, \varphi_K$  are the AR model parameters,  $c$  is a constant and  $\varepsilon_t$  is white noise. The variance of the white noise component determines the degree to which the model is able to fit the stochastic process. AR models are solved using the Yule-Walker equations as found in [12].

### 1.4 Support Vector Machines

Support Vector Machines (SVMs) are a widely popular machine learning technique for classification [13]. The goal of an SVM is to learn an optimal separating hyperplane between classes of data by maximizing the margin between the data points in each class that are closest to data points in the other classes. These closest points are known as support vectors. The support vector machine learns a discriminant function based on these support vectors. The SVM does not make a probabilistic model of the data to determine classification. However, there are techniques [14] to scale the data

in such a fashion that the SVM can return an estimate of the probability that a data point is from a particular class. Fig. 1.5 shows an example of creating a binary decision function using an SVM.

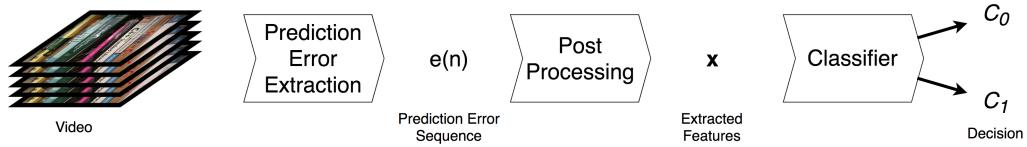


**Figure 1.5:** Visualization of the Optimal Hyperplane Learned by an SVM

An SVM uses a kernel function to determine the similarity or distance between data points. A kernel function is simply a function that can compare the similarity of two input data points. This means that an SVM can learn nonlinear decision surfaces, given the right kernel. One of the most popular kernels to use is the radial basis kernel. The radial basis kernel projects the data onto a Taylor series expansion of a gaussian distribution.

## Chapter 2: Problem Formulation

Detecting frame deletion in a video requires detecting the structural changes in a video due to the deletion process. In particular, Wang and Farid's work on temporal traces for detecting frame deletion shows that for MPEG-2 video, the P-frame prediction error can be formulated into a sequence. This sequence can then be monitored to detect frame deletion. Both Wang and Farid, and Stamm et al. use a system like in Fig. 2.1 to detect frame deletion. The prediction error sequence  $e(n)$  is extracted from the decoded video file and processed to produce detection features. Wang and Farid's work did not propose features for automatic detection, and instead relied on visual inspection of the DFT of the prediction error sequence [4] [7].



**Figure 2.1:** Generalized Approach to Frame Deletion Detection

This broad approach can also be applied to work with H.264 encoded video as well. While video encoding has advanced significantly, the fundamental structures of a compressed digital video have remained unchanged. Regardless of codec, the hallmark of video compression is the motion compensation and estimation process. Video frames are organized into GOPs that begin with an I-frame, and have varying structures of P and B-frames. In H.264, GOP structures are more dynamic due to the ability to derive motion-vector predictions from across multiple anchor frames. To remain robust to these advances in video compression, the prediction error extraction and post processing steps must be altered or augmented.

This work is concerned particularly with the detection of frame deletion in H.264 and similar modern video codecs. Frame addition has also been observed to introduce similar traces in the P-

frame prediction error sequence as frame deletion. Our proposed system can be applied to detecting frame addition and for simplicity we will not discuss the detection of frame addition for the remainder of this thesis. We have made the following assumptions regarding our proposed system. First, we assume that all altered video has undergone re-compression. In fact, since most consumer video recording devices do not have the storage capability or processing power to record high-definition raw video, it is assumed that all video sources have been compressed by either MPEG-4 or H.264, and that all frame deleted video will be re-compressed using H.264 or a similar codec, where the reencoding is set to match the GOP structure of the source video.

In addition, it is assumed that all videos that are passed to the detector are of sufficient length to make a classification. Without multiple full GOPs, the presence of a deletion fingerprint is negligible. Lastly, we make the assumption that if indeed frames have been removed from a video, they have not been removed from the end of the video. The detection features are dependent on differences between the structure of the prediction error sequences in natural videos versus videos with frame deletion. When frames are removed from the end of the video sequence, this difference is not observable.

A user of our proposed system will not need physical access to a specific device to analyze a video captured by the device. The system should accept videos of an arbitrary length, and will not require metadata unrelated to video playback to be intact. It will work with videos of any resolution, frame rate, or GOP structure. Also, as our approach will be data driven, it is imperative that a user have access to a sufficient database of videos with known labels.

## 2.1 Video Frame Deletion Detection

Detecting frame deletion is a binary classification problem. Given a Video  $V$ , there are two possible classes:

$$C_0 : \text{The video is genuine, and has not had frames removed from it.} \quad (2.1)$$

$$C_1 : \text{The video is altered, and has had frames removed from it.}$$

Note that in this case, *genuine* refers to the fact that the video has not undergone any additional

post processing after its original capture. We will simply be considering the limited scenario whereby video has either come directly from the camera that captured it, or frames have been removed from the video and it has been recompressed. Any mention of a genuine video for the rest of this thesis refers to a video that has not been modified since it's original capture.

In general, it is difficult to classify whether or not a video has had frames removed based on the entirety of a video directly. Thus, the problem must be reworked. As shown above, a feature extraction system will be used to produce the P-frame prediction error sequence  $e(n)$ , and a feature vector  $\mathbf{x}$ . The feature vector ideally contains information about the prediction error sequence that can perfectly separate the two classes. As such, the classification problem is as follows. Given a feature vector  $\mathbf{x}'$ , it belongs to one of two classes:

$$C_0 : \mathbf{x}' \text{ resulted from a genuine video that has not had frames removed from it.} \quad (2.2)$$

$$C_1 : \mathbf{x}' \text{ resulted from an altered video which has had frames removed from it.}$$

In the following chapter, we will propose both a new method for extracting  $e(n)$ , and additional augmentations to  $\mathbf{x}$  that allow for improved separation of data and increased robustness of the overall system.

## Chapter 3: Proposed Approach

### 3.1 Prediction Error Sequence Extraction

In previous work on frame deletion detection in MPEG-2, the prediction error sequence was extracted directly from the video decoder using the DCT coefficients of the prediction error residuals located in the compressed video file. The prediction error was averaged over all macroblocks in a frame. This prediction error was then stored as a sequence. Due to the nature of the correlation between P-frame prediction errors across a single GOP, any prediction made across GOP boundaries would result in increased prediction error [4]. Wang and Farid showed that for fixed GOP video, the increase in average prediction error is periodic with respect to the number of frames deleted from the video. Stamm's work expands the idea of the prediction error trace by introducing the formulation of the problem as detecting the presence of a fingerprint signal  $s(n)$ . As H.264 uses variable GOP structures we will only be concerned with the model defined for variable GOP video. Stamm et al. defines the model of  $s(n)$  as

$$s(n) = \beta \mathbf{1} (\Theta(n) = 0). \quad (3.1)$$

where  $\beta > 0$  is a constant and  $\Theta(n)$  is a random variable distributed over the set  $\{0, 1\}$  [7]. This model corresponds to modeling the fingerprint signal as randomly occurring sequence of discrete impulses with a magnitude of  $\beta$ . From this model they pose the detection of frame deletion as distinguishing between two hypotheses:

$$\begin{aligned} H_0 : e(n) &= e_1(n). \\ H_1 : e(n) &= e_2(n) = e_1(n) + s(n)e_1(n). \end{aligned} \quad (3.2)$$

Thus, detection of frame deletion is detection of the presence of the modulated fingerprint signal  $s(n)e_1(n)$ . Given an unknown video, Stamm et al. first makes an approximation of the unaltered

P-frame prediction error sequence. To do this they use a median filter with a filter width of 3.

$$\hat{e}(n) = \text{median}\{e(n-1), e(n), e(n+1)\}. \quad (3.3)$$

Thus, the relationship between the estimate and  $e_1(n)$  is

$$e_1(n) = \hat{e}(n) + \epsilon(n). \quad (3.4)$$

where  $\epsilon(n)$  is a zero mean random variable representing estimation error.

Using the estimate of the unaltered P-frame prediction error sequence, Stamm et al. calculates  $\hat{s}(n)$ , which is an estimate of the fingerprint signal modulated by the prediction error sequence as defined by

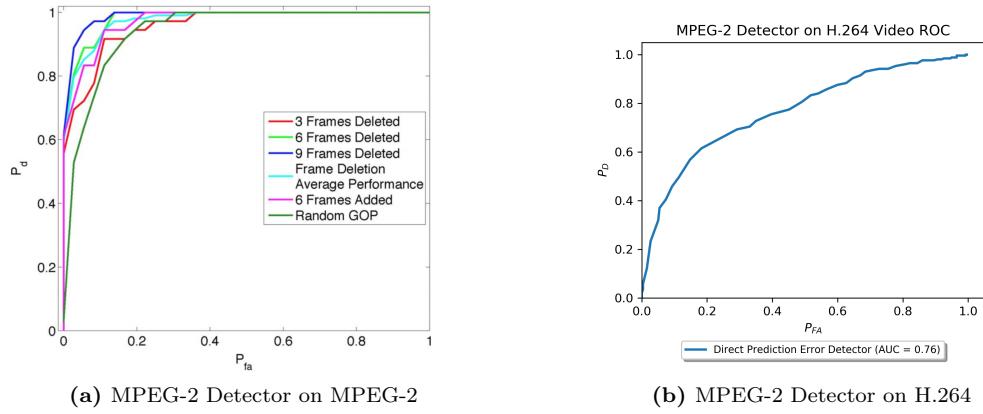
$$\hat{s}(n) = \max(e(n) - \hat{e}(n), 0). \quad (3.5)$$

The estimate of the fingerprint signal is floored at 0, as the model of the  $s(n)$  dictates that it must be greater than or equal to 0. This estimate of the fingerprint signal can be used to build a detector. The decision function found by Stamm et al. for variable GOP video is

$$\delta_{var} = \begin{cases} H_0 & \text{if } \frac{1}{N} \sum_{n=1}^N |\hat{s}(n)| < \tau_{var} \\ H_1 & \text{if } \frac{1}{N} \sum_{n=1}^N |\hat{s}(n)| \geq \tau_{var} \end{cases} \quad (3.6)$$

where the decision is made on the basis of the energy in  $\hat{s}$ .

In MPEG-2, a P-frame is encoded by searching the previous anchor frame for the macroblock which incurs the least error [1]. This means that the average prediction error for a single P-frame is only associated with the previous I or P-frame. H.264 expands the capabilities of its motion compensation and estimation system by allowing prediction from multiple previous frames (and subsequent frames in the case of B-frames) [2]. If the prediction error trace is extracted via the codec for H.264, the average prediction error associated with one frame is comprised of a linear



**Figure 3.1:** Comparison Between MPEG-2 Detection Methods used on (a) MPEG-2<sup>a</sup> and (b) H.264

<sup>a</sup>Figure reprinted with permission from Stamm et al.

combination of the average prediction error associated with motion vectors that map to the different anchor frames used in the motion estimation and compensation process. Thus, cross GOP predictions are smoothed out in such a way that it makes the fingerprint energy detector in Stamm et al.'s paper perform inadequately.

To test this, we collected 257 videos from a cell phone camera (the ASUS ZenFone 3 Laser), and generated an altered video with 15 frames removed from the beginning corresponding to each collected video. The encoding parameters were kept constant, and we fixed the GOP of the altered videos to that of the unaltered videos. In this particular case, the GOP structure was 30 frames in length, with 1 I-frame followed by 29 P-frames. We extracted the prediction error sequence directly from the codec, and measured the estimated fingerprint energy as described above.

As shown in Fig. 3.1 the performance of the detector using the methodology derived for MPEG-2 videos suffers a significant decrease when used on H.264 video, particularly at low false alarm rates. This is due to a limitation in the model used by Stamm et al. above in Equation 3.1. As it is possible to predict across multiple previous anchor frames in H.264, the contribution of the fingerprint signal is variable over time. This variation is also not regular, as scene content and motion determine how many cross GOP predictions are present in each P-frame. Thus we propose the updated model for

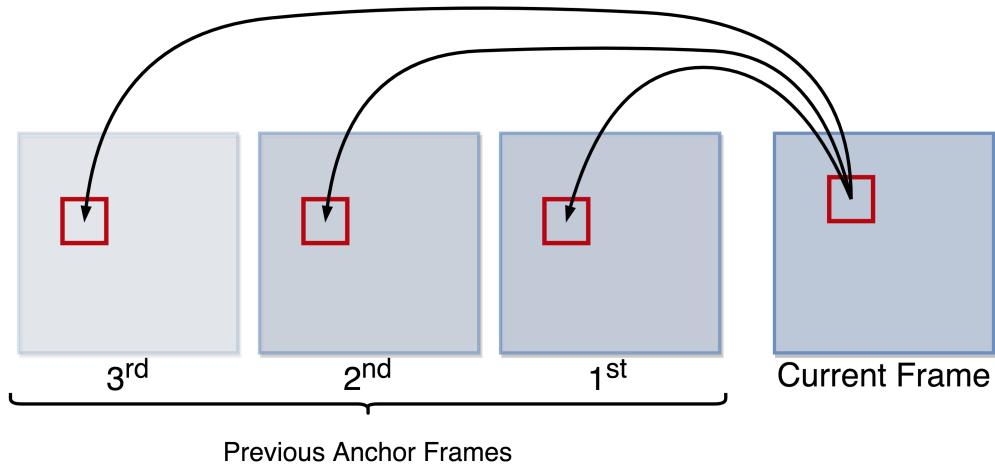
the fingerprint signal as

$$s(n) = \beta(n) \mathbb{1} (\Theta(n) = 0). \quad (3.7)$$

where  $\beta(n)$  is now a random variable that takes values in  $\mathbb{R}_{\geq 0}$ . Thus, we propose the following methodology for extracting the prediction error sequence in H.264.

### 3.1.1 Proposed Prediction Error Sequence Extractor for H.264

The goal of the proposed extraction algorithm is to maximize the probability that should frame deletion exist, a given measurement of the prediction error comes from a cross-GOP prediction. To this end, instead of directly measuring the prediction error from the DCT coefficients from the decoder, we decode the frame of interest and store the motion vectors associated with said frame. For each motion vector in the current P-frame, we find the x and y coordinates defining the source macroblock which provides the least error mapping from a particular previous anchor frame. Then for that previous anchor frame, we subtract the pixels in the source macroblock from the destination macroblock in the current frame. This leaves us with a prediction error residual associated with the motion vector. We then calculate the average absolute value of this residual.



**Figure 3.2:** Identification of the source macroblock in 3 previous anchor frames

We repeat this process for each of D previous anchor frames. Then we store these prediction error values in a matrix  $M[n]$ , where  $n$  denotes the P-frame index of the current frame. Each row

of the matrix corresponds to the errors associated with a single motion vector, and the columns are the errors associated with each of the previous anchor frames.

After obtaining the  $M[n]$  matrix for the current frame, we create the matrix  $\tilde{M}[n]$  defined like so:

$$\tilde{M}_{i,j}[n] = \mathbb{1} \left( j = \operatorname{argmin}_l (M_{i,l}[n]) \right) * M_{i,j}[n] \quad (3.8)$$

Thus  $\tilde{M}[n]$  is a copy of  $M[n]$  where the non-zero entries in each row correspond to the minimum average error associated with the macroblock and all other elements are zero. Effectively, the column index of the non-zero entry is an estimation of which previous frame the motion vector associated with the macroblock maps to in the decoding process.

Further processing is done on  $\tilde{M}[n]$  to output only a single prediction error value. First,  $\tilde{M}[n]$  is reduced into a vector  $P[n]$ , such that:

$$P_j[n] = \frac{1}{N_j} \sum_i \tilde{M}_{i,j}[n] \quad (3.9)$$

Where  $N_j$  is the number of non-zero elements in the  $j^{th}$  column of  $\tilde{M}[n]$ . Then, the reported prediction error for the current frame  $e^*[n]$  is calculated as

$$e^*[n] = \max_j P_j[n] \quad (3.10)$$

This entire process is repeated for every P frame in the video. This method of error extraction estimates which previous anchor frame contributes the maximum error per macroblock to the overall prediction error residual obtained by the codec for a given P frame. Since prediction across GOP boundaries results in spikes in the prediction error, the anchor frame that contributes the most error is most likely to be from a different original GOP. In this manner, we obtain a trace that is resilient to advances made in the motion compensation and estimation process in modern codecs as well as robust to variable frame rates and dynamic GOP structures.

$$M[n] = \begin{bmatrix} 30.4 & 20.5 & 40.2 \\ 16.3 & 22.1 & 30.4 \\ 25.5 & 23.4 & 19.8 \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$\tilde{M}[n] = \begin{bmatrix} 0.00 & \mathbf{20.5} & 0.00 \\ \mathbf{16.3} & 0.00 & 0.00 \\ 0.00 & 0.00 & \mathbf{19.8} \\ \vdots & \vdots & \vdots \end{bmatrix}$$

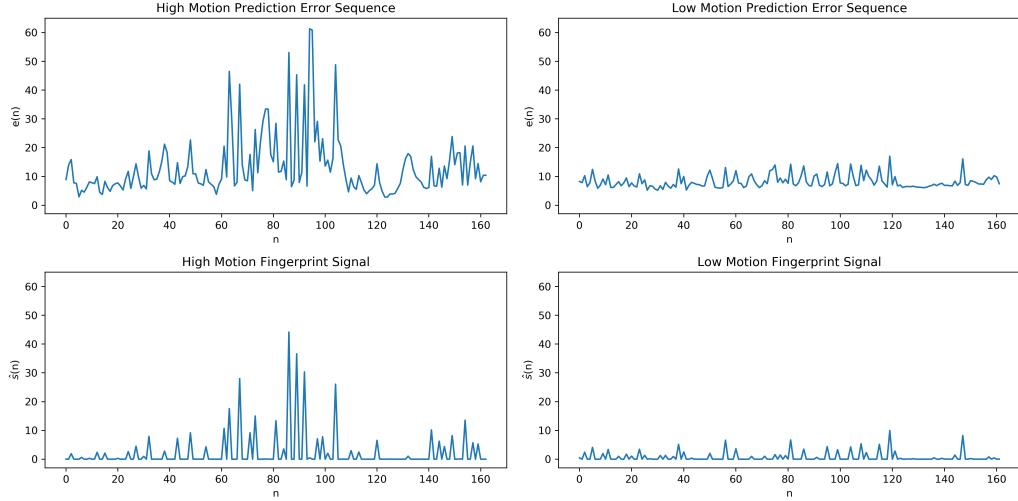
$$P[n] = [P_1[n] \quad P_2[n] \quad P_3[n]]$$

**Figure 3.3:** Methodology for forming  $P[n]$

### 3.2 Proposed Detection Algorithm

In addition to the new methods for prediction error extraction, we propose an expanded detection algorithm to better capture the statistical differences between videos. In fact, depending on scene content, video capture settings, and the amount of motion captured in a single recording, the prediction error sequence and fingerprint signal exhibit different structural behavior. This is true even for videos captured from a single camera model. Figure 3.4 shows this clearly. The two videos were captured from an LG Nexus 5X using the high quality 1080p capture mode. Both videos were shot using similar scene content, but the amount of motion in each video is different. The first video was shot with high motion, while the second video was comparatively low motion. The top row shows the different prediction error sequences, while the bottom row shows the different estimated fingerprint signals.

Notice the large discrepancy in magnitude between both types of signals depending on the amount of motion in the video. The original detection criteria defined by Stamm et al. in Equation 3.6 is based on the energy of the estimated fingerprint signal. This will lead to undesirable misclassifica-

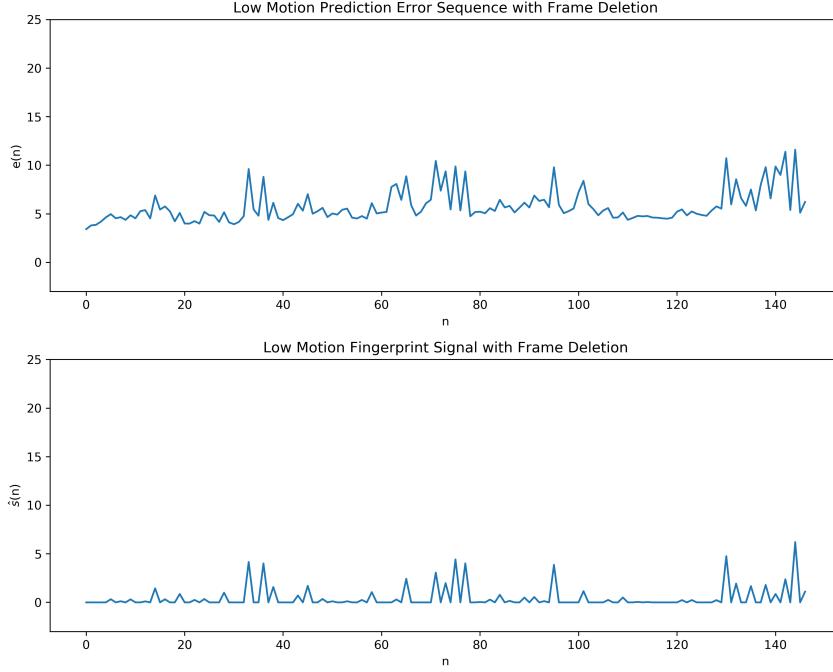


**Figure 3.4:** *Top Left* - The extracted prediction error sequence for a high motion video. *Top Right* - The extracted prediction error sequence for a low motion video. *Bottom Left* - Estimated fingerprint signal for a high motion video. *Bottom Right* - Estimated fingerprint signal for a low motion video.

tions when only using signal energy as the detection feature. Thus, we need a set of new features that can help account for this difference in fingerprint signal energy between videos.

Under the old model for the fingerprint signal,  $\beta$  was a constant, meaning the fingerprint signal was thought of as a randomly occurring sequence of discrete impulses with a magnitude of  $\beta$ . With the new model defined in Equation 3.7,  $\beta(n)$  is a random variable that takes nonnegative real values. The effect of  $\beta(n)$  can be seen in Fig. 3.5. The sample low motion video was reencoded with the first 15 frames removed. This accounts for half of a GOP. As the GOP structure of the video is fixed, it is expected that the fingerprint signal will be periodic [7]. The estimated fingerprint signal of the sample video with frame deletion shows some amount of periodicity but the amplitude of the signal varies with time. Modeling this variation of  $\beta(n)$  for a given video can help aid in the detection of frame deletion.

Over several seconds of video both the prediction error sequence and the fingerprint signal are wide sense stationary. Due to this, we propose modeling both the fingerprint signal and prediction error sequence as autoregressive (AR) processes. The model parameters capture some of the statistical information about  $\beta(n)$ , and thus are added to a feature vector along with the fingerprint



**Figure 3.5:** *Top* - The extracted prediction error sequence for a low motion video with frame deletion. *Bottom* - The Estimated fingerprint signal for a low motion video with frame deletion.

energy. In order to capture the degree to which the model fits a given sequence, the error variance of each AR model is also included. In addition, some basic statistical features are included to scale the overall decision surface. We propose including the mean and variance of both the prediction error sequence and fingerprint signal to the feature vector as well.

For a given video  $V$ , the feature vector used for classification  $\mathbf{x}_V$  is structured as

$$\mathbf{x}_V = \begin{bmatrix} \frac{1}{N} \sum_{n=1}^N |\hat{s}[n]| \\ \mu_{\hat{s}[n]} \\ \sigma_{\hat{s}[n]}^2 \\ \mu_{e^*[n]} \\ \sigma_{e^*[n]}^2 \\ \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix} \quad (3.11)$$

Where  $\hat{s}[n]$  is the estimated fingerprint sequence of  $e^*[n]$  obtained by using Equation 3.5 above,

but substituting  $e^*[n]$  for  $e[n]$ ,  $\mu$  is the sample mean,  $\sigma^2$  is the sample variance,  $\mathbf{a}_1$  are the  $Q^{th}$  order AR model parameters for  $e^*[n]$ , and  $\mathbf{a}_2$  are the  $Q^{th}$  order AR model parameters for  $\hat{s}[n]$ .

Note that after creating a feature vector for each video, there is no clear way to make a simple classification. It is inadvisable to create a probabilistic model of the feature vector for classification as there is no clear distribution for the AR model parameters. Instead, we propose using a discriminative function to map an incoming feature vector directly to the set of natural videos or the set of videos altered by frame deletion. As such, we propose using a Support Vector Machine (SVM) classifier with a Radial Basis Kernel function for classification [13].

## Chapter 4: Experimental Results

We conducted a series of experiments to evaluate both the performance of our proposed prediction error extraction algorithm and new detection features. In order to create a suitable dataset for our experiments, we gathered a set of at least 250 five second long video clips from a variety of camera models. As the video processing pipeline is different between camera models, we aimed to include a diverse set of GOP structures. To create our frame deleted videos, we then removed 15 frames from the beginning of each video in our dataset. The reencoding was handled using the x264 library, and we matched the original GOP structure of each video for the frame deleted copy. We then filtered out all videos with less than 60 P-frames, as having fewer than 60 P-frames is not sufficient for making a reasonable classification of a video sequence. The resulting dataset can be found in Table 4.1.

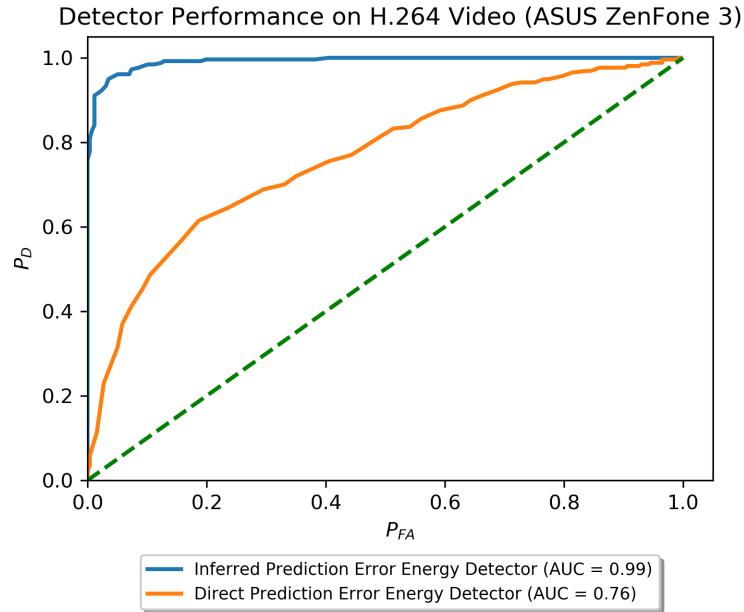
**Table 4.1:** Composition of Frame Deletion Dataset

Camera Model	GOP Type	GOP N	GOP M	Genuine Videos	Frame Deleted Videos
ASUS ZenFone 3 Laser	Fixed	30	1	257	257
Apple iPhone 8 plus	Variable	30	2	248	245
Google Pixel 1	Variable	29	2	776	777
Google Pixel 2	Variable	29	2	212	212
Kodak Ektra	Fixed	30	1	503	502
Nokia 6.1	Fixed	30	1	501	500
Samsung Galaxy S7	Fixed	30	1	231	231
<b>Total Count</b>				<b>2728</b>	<b>2724</b>

Note that GOP N is the maximum distance between I-frames in a video, which is the GOP length. GOP M is the maximum distance between P-frames in a video, thus a video with an M value of 1 has no B-frames at all. Note that some deleted videos had just enough frames removed to put them below the 60 P-frame threshold. Thus, there is a slight imbalance between the total number of videos present in each class. As well, data gathered from camera models with more than 275 videos was captured using multiple individual devices. The videos from the Google Pixel 1 contain footage captured on three devices, while the Kodak Ektra and Nokia 6.1 contain video from two devices.

## 4.1 Single Camera Model

To gauge relative performance of our proposed detection methods, we first tested on a single camera model. We chose to use videos from the ASUS ZenFone 3 Laser. When testing for the effects of our proposed prediction error extraction method, we used the decision metric proposed by Stamm et al. for MPEG-2 video, as described in Equation 3.6. We varied the decision threshold used in each detector over a range of possible values. The probabilities of detection ( $P_D$ ) and false alarm ( $P_{FA}$ ) were measured for each value of the threshold. We then calculated the Receiver Operator Characteristic (ROC) curves for both prediction error extraction methodologies.

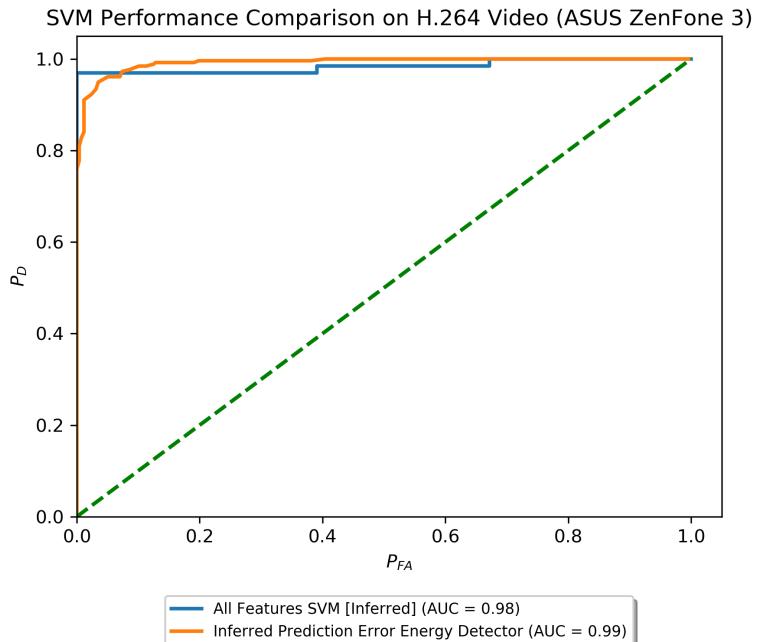


**Figure 4.1:** ROC Curves for Fingerprint Energy Detector Obtained by Using Different Prediction Error Extraction Methods

As shown in Fig. 4.1, extracting the prediction error sequence directly from the codec (direct prediction error) results in poor detector performance. The detector build from the direct prediction error sequence is only able to reach 60% detection rates at 20% false alarm rates. Compare to our proposed prediction error extraction methodology (inferred prediction error). For the ASUS ZenFone3, we achieve over 90% detection rates at less than 10% false alarm rates. These results indicate that indeed our new inferred prediction error extraction methodology is able to yield separation be-

tween the two classes of data where the original direct prediction error extraction methodology does not, particularly at low false alarm rates. There is a delta upwards of 70% in the probability of detection for a given false alarm rate.

Next, we examined the effects of using our proposed feature vector and an SVM on classification accuracy. We set the AR model orders for both the prediction error sequence and the fingerprint signal to 31. This is because of the periodic nature of the fingerprint on fixed GOP video. We want to predict when the signal will spike based on the previous GOP. We used the SVM with a radial basis function as the kernel and set the hyper-parameters of the SVM and the radial basis kernel to their defaults. To train the SVM, we partitioned our data into training and testing sets. The testing set was sized to be a quarter of the total dataset. Thus, there were 128 videos used for testing and 386 videos used for training. We set the SVM to use Platt scaling to return a probability instead of a class label [14]. By doing so, we were able to vary the detection threshold and construct an ROC curve from the testing set for comparison with the fingerprint energy detector. The results can be seen in Fig. 4.2.

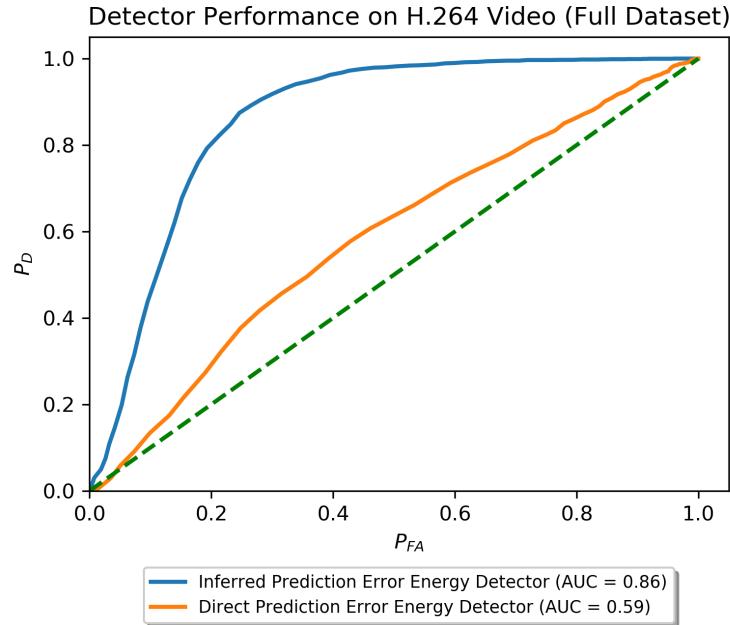


**Figure 4.2:** ROC Curves for Inferred Prediction Error Energy Detector and SVM on One Camera Model

As the detection rates for the fingerprint energy detector are already above 90% for low false alarm rates, the relative magnitude of improvement when adding the SVM is less than when changing the prediction error extraction methodology. However, there is a significant improvement at false alarm rates under 10%. The SVM detector is able to hold above 95% accuracy at a 0% false alarm rate. This suggests that some of the frame deleted videos do not have enough statistical differences from genuine videos to be classified as having frames removed from them. The videos in question may have low motion or content such that the fingerprint signal does not express itself in a detectable manner.

## 4.2 Multiple Camera Models

In order to obtain a better estimate of real world detector performance, we tested the effects of our proposed inferred prediction error extraction method using the entire dataset. We constructed ROC curves for detectors derived from both prediction error extraction methods. The resulting ROC curves can be seen in Fig. 4.3

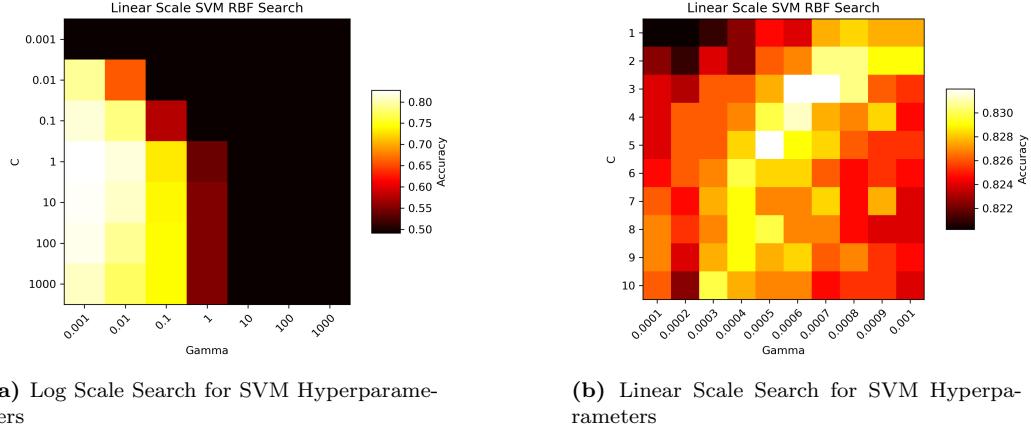


**Figure 4.3:** ROC Curves for Fingerprint Energy Detector Obtained by Using Different Prediction Error Extraction Methods on Multiple Camera Models

It is obvious that the results obtained when testing solely on the ASUS ZenFone 3 were the

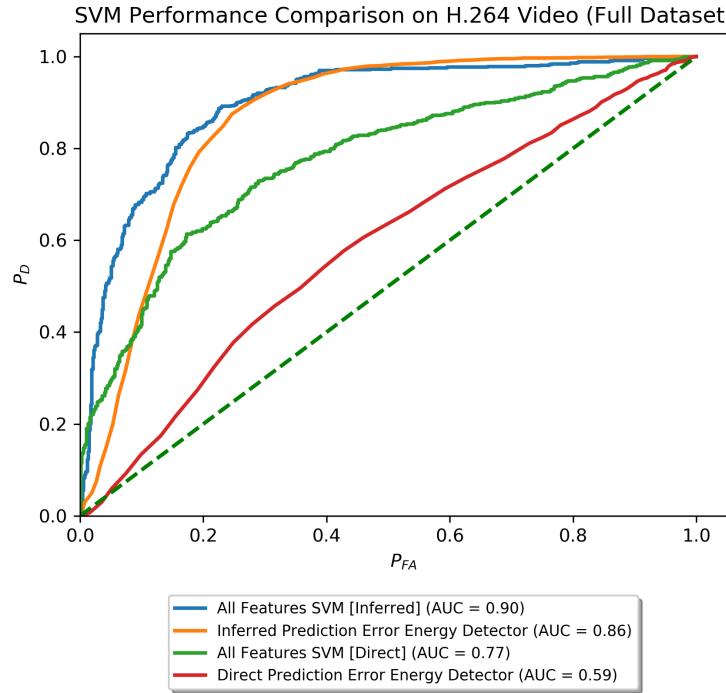
best case scenarios for both prediction error extraction methods. In fact, using the direct prediction error extraction method from prior research does not yield much better detector performance than random chance. While not optimal, our proposed prediction error extractor provides a much cleaner separability of the data, even when confronted with a wide variety of videos. Without using the SVM, we are able to raise the probability of detection from 15% to over 60% at a 10% false alarm rate.

Before we tested the effects of our expanded feature vector and SVM on the entire dataset, we tuned the SVM to have optimal hyperparameters. We tuned both the penalty of the error term ( $C$ ) and the size of gaussian used in the radial basis kernel ( $\text{Gamma}$ ). We performed a grid search first in the log-scale to determine their optimal order of magnitude.



**Figure 4.4:** Search for Optimal SVM Hyperparameters on (a) Log Scale and (b) Linear Scale

After determining the optimal order of magnitude for both of the hyperparameters, we performed a second grid search, but on a linear scale to find the set of optimal parameters. Fig 4.4 shows the results of both the log and linear scale grid searches. We chose to use  $C = 5$  and  $\text{Gamma} = 0.0005$  for our SVM. As we did previously, we split our entire dataset into training and testing sets to evaluate the SVM, where the testing set consisted of a quarter of the data. The dataset was shuffled prior to splitting so that a sampling of all camera models would be included in the testing set. Fig. 4.5 shows the ROC curves for the fingerprint energy detectors for both the direct and inferred prediction error extraction methods as well as SVMs.



**Figure 4.5:** ROC Curves for Various Detection Methods on Data from Multiple Camera Models

Using our proposed feature vector and SVM with the direct prediction error extraction method allows for some discrimination between the classes. However, the overall detector performance is still less than that of the fingerprint energy detector using our proposed inferred prediction error extraction methodology. Using our proposed feature vector and SVM with our proposed inferred prediction error extraction methodology yields the best detector performance in this scenario, particularly at low false alarm rates. At a 10% false alarm rate, our combined detector yields an almost 80% probability of detection.

### 4.3 Additional Work

As we have shown, our proposed detector is able to distinguish between genuine videos and videos with frame deletion with a reasonable level of accuracy. Indeed, it is also more effective than previously formulated motion vector based detectors, particularly those formulated for use on MPEG-2. However, this is only in a fairly limited scenario. Many videos, especially those shared on social

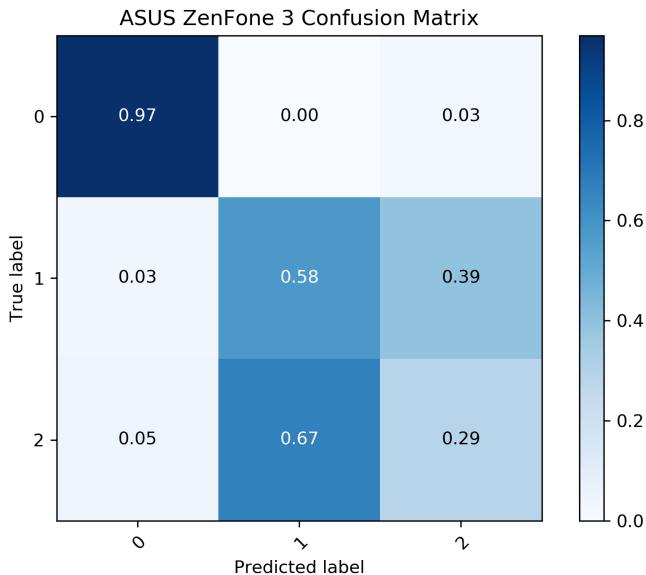
media, are often recompressed, even without having frames removed from them. We tested our proposed detector with a three class classification problem. Given a feature vector  $\mathbf{x}$ , it belongs to one of three classes:

$C_0 : \mathbf{x}$  resulted from a genuine video that has not had frames removed from it.

$C_1 : \mathbf{x}$  resulted from an altered video which has had frames removed from it. (4.1)

$C_2 : \mathbf{x}$  resulted from a video that has been recompressed without having frames removed.

To test this, we created a third set of videos for the ASUS ZenFone 3, where each unaltered video in the dataset was recompressed. Thus we had a dataset of 257 videos for each class. We trained our SVM using the same training and testing data ratio as above, and found optimal parameters for this problem. We tabulated our results into the confusion matrix in Fig. 4.6.

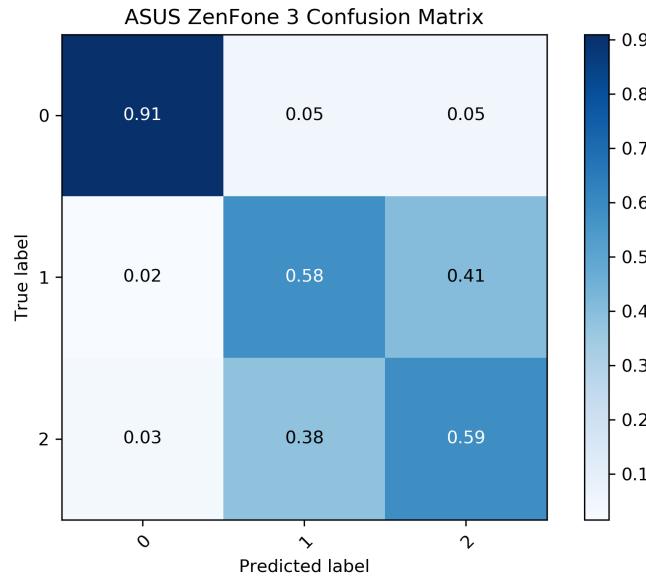


**Figure 4.6:** Confusion Matrix for Three Class Classification

From this confusion matrix, it becomes clear that more work is needed on our proposed detector. We are able to very easily distinguish whether videos have been recompressed or not, but it is a much more difficult problem to pinpoint the differences between videos that have had frames removed and videos that have only been recompressed. Additional work is needed to make an explicit detector for this scenario. Systems that can learn complex spacial relationships, like Convolutional Neural

Networks may be able to help.

It is possible to tune our proposed detector slightly to increase the overall accuracy of the system. If we scale our data to have zero mean and unit variance on a feature by feature basis, we can trade some accuracy in the genuine class for better distinction between videos with frame deletion and recompressed videos. Fig. 4.7 shows that we can now detect videos with frame deletion and recompressed videos with just under 60% accuracy. This puts our overall system accuracy at just under 70%. While this is a far cry from the over 95% accuracy on the same camera model with only two classes, it shows that finding differences between videos with frame deletion and recompressed videos is possible.



**Figure 4.7:** Confusion Matrix for Three Class Classification using Scaled Data

## Bibliography

- [1] ISO Central Secretary, “Information technology - Generic coding of moving video pictures and associated audio information - part 2: Video”, en, International Organization for Standardization, Geneva, Switzerland, Standard ISO/IEC IS 13818-2, 2008.
- [2] ——, “Information technology - coding of audio-visual objects - part 10: Advanced video coding (avc)”, en, International Organization for Standardization, Geneva, Switzerland, Standard ISO/IEC IS 14496-10, 2010.
- [3] T. Wiegand, G. J. Sullivan, G. Bjontegaard, and A. Luthra, “Overview of the h.264/avc video coding standard”, *IEEE Trans. Cir. and Sys. for Video Technol.*, vol. 13, no. 7, pp. 560–576, Jul. 2003, ISSN: 1051-8215. DOI: 10.1109/TCSVT.2003.815165. [Online]. Available: <https://doi.org/10.1109/TCSVT.2003.815165>.
- [4] W. Wang and H. Farid, “Exposing digital forgeries in video by detecting double mpeg compression”, in *Proceedings of the 8th Workshop on Multimedia and Security*, ser. MM&#38;Sec ’06, Geneva, Switzerland: ACM, 2006, pp. 37–47, ISBN: 1-59593-493-6. DOI: 10.1145/1161366.1161375. [Online]. Available: <http://doi.acm.org/10.1145/1161366.1161375>.
- [5] J. Lukáš and J. Fridrich, “Estimation of primary quantization matrix in double compressed jpeg images”, in *Proc. Digital forensic research workshop*, 2003, pp. 5–8.
- [6] A. C. Popescu and H. Farid, “Statistical tools for digital forensics”, in *International Workshop on Information Hiding*, Springer, 2004, pp. 128–147.
- [7] M. C. Stamm, W. S. Lin, and K. J. R. Liu, “Temporal forensics and anti-forensics for motion compensated video”, *IEEE Transactions on Information Forensics and Security*, vol. 7, no. 4, pp. 1315–1329, 2012, ISSN: 1556-6013. DOI: 10.1109/TIFS.2012.2205568.
- [8] J. Liu and X. Kang, “Exposing heterogeneous chain of video recompression”, in *2016 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA)*, 2016, pp. 1–6. DOI: 10.1109/APSIPA.2016.7820694.
- [9] S. Mandelli, P. Bestagini, S. Tubaro, D. Cozzolino, and L. Verdoliva, “Blind detection and localization of video temporal splicing exploiting sensor-based footprints”, in *2018 26th European Signal Processing Conference (EUSIPCO)*, 2018, pp. 1362–1366. DOI: 10.23919/EUSIPCO.2018.8553511.
- [10] D. Vazquez-Padin, M. Fontani, T. Bianchi, P. Comesana, A. Piva, and M. Barni, “Detection of video double encoding with gop size estimation”, in *2012 IEEE International Workshop on Information Forensics and Security (WIFS)*, 2012, pp. 151–156. DOI: 10.1109/WIFS.2012.6412641.
- [11] A. Gironi, M. Fontani, T. Bianchi, A. Piva, and M. Barni, “A video forensic technique for detecting frame deletion and insertion”, in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2014, pp. 6226–6230. DOI: 10.1109/ICASSP.2014.6854801.
- [12] “Adaptive filter theory”, in. Prentice Hall, 2002, ch. 1.
- [13] C. Cortes and V. Vapnik, “Support-vector networks”, *Machine Learning*, vol. 20, no. 3, pp. 273–297, 1995, ISSN: 1573-0565. DOI: 10.1007/BF00994018. [Online]. Available: <https://doi.org/10.1007/BF00994018>.
- [14] J. Platt, “Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods”, *Adv. Large Margin Classif.*, vol. 10, Jun. 2000.

