

# Video Fingerprinting Using Robust Hashing of Scale Invariant Features of Frames

R. Moolman and W.C. Venter

School of Electrical, Electronic and Computer Engineering

North-West University, Potchefstroom Campus

Tel: +27 18 299 1961, Fax: +27 18 299 1977

Email: {20673892, willie.venter}@nwu.ac.za

**Abstract**—Video fingerprinting is a technique used to identify an unknown video by matching the video’s fingerprint to a database of video fingerprints. The fingerprints are derived through algorithms and characterizes the video’s content. Video fingerprinting automates video identification and it has many uses these days with the exponential increase in video usage. A high quality video fingerprinting system should be fast, robust to distortions and use storage space efficiently. All these aspects are dependant on the technique used to fingerprint the videos. In this paper a novel video fingerprinting technique is proposed. This technique was developed with the intention to detect videos in real time, with its primary use being advertisement tracking. The technique makes use of the Scale Invariant Feature Transform (SIFT) algorithm, combined with the idea of hashing introduced by Shazam, an audio fingerprinting technique, to fingerprint and detect single frames. The aspects of the frame detector are discussed and tested, after which the optimal variable values are chosen for frame detection.

**Keywords:** video fingerprinting; copy detection; automatic video recognition; perceptual frame hashing; content-based video identification; robust matching

## I. INTRODUCTION

In recent years the use of video has increased dramatically. Video fingerprinting is a technique that can be implemented to automate the management of this large amount of video data. For this reason it has received a lot of attention in research over the last few years. Video fingerprinting can be used for copyright management [1], advertisement tracking [2][3] and video management on personal systems.

In the video detection system a known video is run through an algorithm to extract its fingerprint and then it is saved to the database. When detecting an unknown video, the same algorithm is used to fingerprint the video, except this time the fingerprint is compared to the database to find the best match. If the video’s fingerprint was added to the database, the unknown video should be detected. As all video fingerprinting systems aim for the same goal, there are certain traits that are expected of the system, namely:

- **Robustness:** The ability to detect videos while distortion is present.
- **Accuracy:** The ability to give unique fingerprints to different videos.
- **Speed:** The ability to fingerprint and detect videos quickly.
- **Efficiency:** The ability to effectively store and match fingerprint data in a database.

A lot of previous work has been done in this field, including [4], which uses Radial Projection of key frames, Gradient of Centroids [5][6], Visual Attention Areas [3], Average Luminance over Time [2], Colour Histograms [7] and Visual Digest of Local Fingerprints [8]. Most of the existing techniques focus on the robustness and accuracy of the video fingerprinting, but seem to neglect the speed of the video fingerprint matches. This research aims to expand what has been done in the field by developing a new technique that focuses on speed, while maintaining robustness, accuracy and efficiency.

In Section II the technique is explained in detail, and in Section III video detection is discussed. Section III is separate because it makes use of a basic key frame detector and the frame fingerprinting to detect videos. Sections II and III both include subsections where tests are discussed and results are shown. Lastly, Section IV contains a conclusion along with proposed future work.

## II. FRAME FINGERPRINTING

While creating a robust frame detection algorithm to be used in a video fingerprinting system, the main focus was on developing an algorithm that can process and detect frame matches as fast as possible while still maintaining robustness and database efficiency.

Shazam is an audio fingerprinting technique that makes use of hashes, created from pairs of key points found on the song’s spectrogram, to quickly match an unknown audio file to the database [9]. To utilize this technique to detect frames as fast as possible, an algorithm called SIFT was used. SIFT detects robust, scale and rotation invariant key points in an image and was originally designed for object recognition [10]. Each SIFT key point consists of a x and y coordinate, a magnitude and a direction. By using the key points detected by the SIFT algorithm, Shazam-like hashes can be created to fingerprint the frames. The proposed algorithm, database and searching method will be discussed in the section below.

### A. Algorithm

First, the video is sampled and each frame normalised by re-sampling it to a set resolution of  $320 \times 240$ , after which the frames are converted to gray-scale. Key points for the frame are then calculated using the SIFT algorithm. The key points with the largest magnitudes are used to create the hashes, because they are more robust and repeatable than smaller key points and the processing time is greatly reduced if only the largest key points have to be detected.

Each key point in the frame is matched with other key points found in the frame that are within the specified region

around the key point. The values of the magnitude ratio, the direction difference, relative distance and direction to the secondary key point are then calculated and normalised to a set number of bits per value, after which they are then combined to create a hash value. Key points are only matched with smaller key points to ensure that the magnitude ratio is smaller than 1 and to remove any hash redundancy by not repeating key point pairs. Once a certain number of hashes per frame is calculated, the hashing ceases. This hash limit is discussed in Section II-D2.

### B. Database

A Video Frame Identification (VFID) is a number each fingerprinted frame of a video is given, to identify it with. It is these VFIDs that are saved in a database in such a way that it can be retrieved very quickly. For each hash value a specific number of spots are allocated for VFIDs within the database. If a hash occurs in a frame, the frames' VFID is saved in one of that hash's spots. The hashes are very unique, so the spots will fill up slowly. The database may be sized according to the requirement of the system by changing the amount of bits used to represent the VFIDs or by changing the number of available VFID spots per hash. The database will be optimised to use as little storage space as possible per frame fingerprint to allow it to contain as many video's fingerprints as possible.

Video and frame data is also saved with the VFIDs in a supplementary database, so the VFIDs carry meaning. This way, when a certain VFID is detected when querying a frame, the information about the originating video and frame number is available to use.

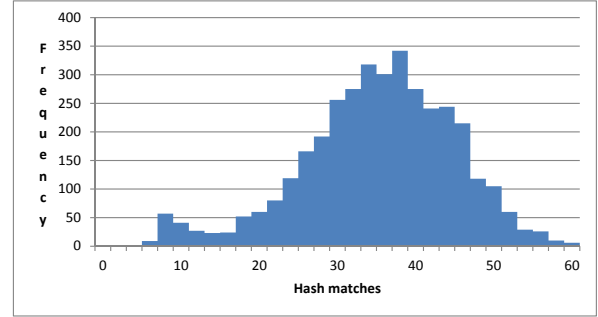
### C. Searching

To match a frame's fingerprint to the database, the video fingerprinting algorithm is applied to the frame to generate a set of hashes. All the VFIDs for every generated hash value are extracted from the database and every time a VFID occurs, it gets a point. After all the hashes' corresponding VFIDs have been accounted for, the VFID with the most accumulated points above a specified threshold is the most credible match.

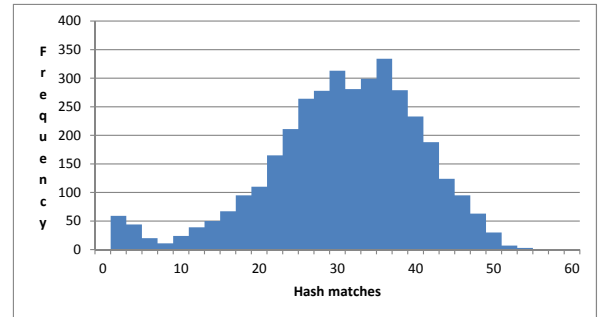
This searching technique is based on the Shazam algorithm's searching technique. The search is very quick and the speed stays more or less constant even when the database becomes really large [9]. The reason for this is that it doesn't have to do an exhaustive search by scanning through the entire database for a match, it only has to jump to the hashes that match exactly and extract its information. The maximum search time is a combination of the time it takes to read the VFIDs (maximum of  $hashes \times number\ of\ VFID\ slots$ ) and the time it takes to determine which VFID has the most hash matches.

### D. Tests

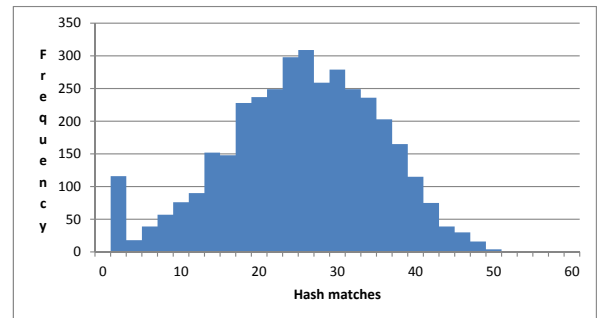
The frame fingerprint consists of a set of hashes, that is created from key points, as mentioned above. While comparing an unknown frame to the database, the hashes created for the frame should match exactly to the hashes created from the original frame for it to be detected as a match, even if there are distortions present. The variables that affect detection are the number of bits used for a hash, the number of hashes used to fingerprint a frame and the



(a) Frequency histogram of hash matches for AVI XviD using 12 bit hash



(b) Frequency histogram of hash matches for AVI XviD using 16 bit hash

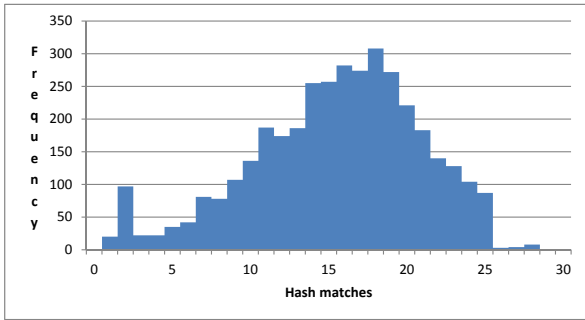


(c) Frequency histogram of hash matches for AVI XviD using 20 bit hash

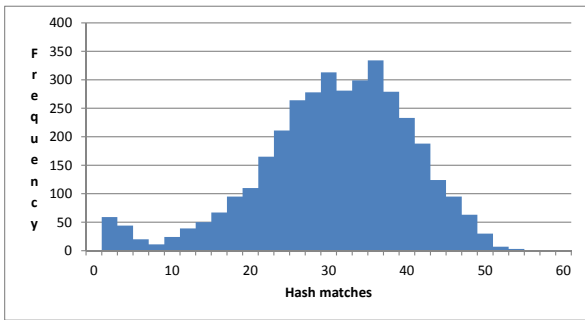
**Figure 1:** Results for *Bits per Hash* test

detection threshold. Each of these variables will be discussed and the tests done, to determine the best value for each variable, will be covered in this subsection.

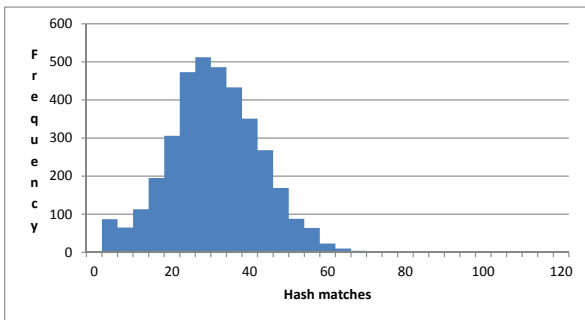
1) *Bits per Hash:* Each hash value consists of four variables normalised to use only a certain amount of bits. The number



(a) Frequency histogram of hash matches for AVI XviD using 25 hashes



(b) Frequency histogram of hash matches for AVI XviD using 50 hashes



(c) Frequency histogram of hash matches for AVI XviD using 100 hashes

**Figure 2:** Results for *Detection Threshold and Hashes per Frame* test

of bits used for each hash is important, because hash values have to be exactly the same to get a match with the database. The higher the number of bits used, the higher the resolution, thus the hashes are more specific and its harder to get matches. On the other hand, using less bits to create hashes, means that its easier to get matches, but it increases the chances of false matches. It also increases the change of getting two or three of the same hash values in one frame.

The number of bits used per hash value is a trade-off between accuracy and robustness. Using more bits per hash will result in a higher hash resolution and thus increase the accuracy of a match, but the robustness will decrease as small changes may cause hash matches to be missed, resulting in false negative frame matches. On the other hand, if less bits are used the number of false positive matches will increase, but the robustness will also be greatly increased, because small errors will be allowed.

To determine the optimal number of bits to use for creating hashes, the following test was done: Firstly, the number of hashes per frame were set to 50. Then a set of 30 music videos were fingerprinted and their fingerprints saved to the database using a specific number of bits per hash. Music videos were used because of the wide array of images and scenes used in the videos. Once the database was created, the videos were converted to other formats (using FormatFactory) and run through the detection algorithm. In Figure 1 frequency histograms for the number of hashes matched are shown. The frequency histograms are made up of 3687 frame match's data. The test also included determining the second best match (first false positive) to frame, to determine the usual number of matches a frame gets if it's not in the database, but the data is only mentioned in Table I.

Figure	Best match average	Best match standard deviation	Second match average	Second match standard deviation
1(a)	34.83	10.02	8.51	2.6
1(b)	30.31	9.85	2.81	1.97
1(c)	24.93	9.9	1.3	1.38

Table I  
QUANTITATIVE RESULTS FOR BITS PER HASH TEST

If one looks at the graphs in Figure 1 one can see that the data distribution is almost normal, except for the outliers at the very low values. These outliers are caused by the frame detector algorithm, because frames are sometimes detected as key frames in the avi files, while they weren't detected in the mp4 files.

If one compares the data seen in Table I, one can see that the 12 bit hashes have a higher match average than the 16 bit hashes, but the lower resolution of the 12 hash matches cause a big increase in the average of the secondary matches. In this case the 16 bit hashes prove advantageous, as it keeps the hash uniqueness higher (good for database) with only a little decrease in hash matches. The 20 bit hashes have lower averages than the 16 bit hashes, but both options are acceptable for implementation. The deciding factor is the balance between uniqueness, robustness and efficiency, and thus, the 16 bit hashes are better.

2) *Detection Threshold and Hashes per Frame:* Every single fingerprinting system makes use of thresholds to determine if the query fingerprint has a legitimate match in the database. This frame fingerprinting technique, makes use of a minimum

hash matches threshold to determine if the match is successful. If the best VFID match's number of matches is not above the threshold, it is not considered a match. Another variable that also affects the frame detector is the number of hashes created for each frame fingerprint. Both the number of hashes per fingerprint and the minimum hash match threshold is determined by the same set of test as the minimum hash match threshold is dependant on the number of frames used to fingerprint a frame.

To determine these thresholds, the following test was done: As with the bits per hash test, a set of music video's fingerprints were added to the database, using 25 as the maximum hash threshold. The test was repeated again two times while using 50 and 100 as the maximum hash threshold. The original videos that were fingerprinted and added to the database, used the mp4 wrapper and DivX codec, while the videos used to match to the database, are in avi wrappers with XviD encoding. In Figure 2, the results for the tests are shown in the form of frequency histograms. The frequency histograms are made up of 3716, 3687 and 3650 frame match's data, respectively. The differences are due to the fact that frames are only added to the database if their number of hashes exceed the maximum hash threshold. The test also included determining the second best match (first false positive) to frame, as with the bits per hash test and are included in Table II.

Figure	Best match average	Best match standard deviation	Second match average	Second match standard deviation
2(a)	15.51	5.44	1.92	1.03
2(b)	30.31	9.85	2.81	1.97
2(c)	29.48	11.64	4.41	2.82

Table II  
QUANTITATIVE RESULTS FOR HASHES PER FRAME TEST

If the results in Table II are compared, it is clear to see that the average number of hash matches stay almost constant when using 100 hashes per frame instead of 50. The reason for this is that the biggest key points are the most robust, allowing the first 50 or so hashes to be the most repeatable. The test was also done with flv (flv1) and mpg (mpeg1) files and the results reflect the same information, although the data is not shown in this paper.

Detection was also tested while using 25 hashes per frame. The results show that the averages of the 25 hashes per frame are basically half of the 50 hashes per frame averages, thus 50 hashes per frame will be used to increase the chance of hash detection and allowing more space to choose a detection threshold.

After close inspection of the results, it was seen that frames that have more than 10% hash matches of the hashes used per frame are positive matches, almost 100% of the time. There may still be exceptions, so it is safe to say that the minimum hashes detection threshold can be set to 20%. So, for detection with 50 hashes per frame, any frame match that has more than 10 hash matches can be considered a positive match.

### III. VIDEO FINGERPRINTING

The video fingerprinting technique discussed in this paper was developed with advertisement tracking as its main objective. For this a real time system is needed to monitor television stations constantly.

To upgrade the frame detecting algorithm to a video detecting algorithm, a set of frames from the original video must be fingerprinted so they can be detected. Fingerprinting every single frame is very inefficient, so a key frame detector was needed that could detect the same frames in videos with the same content. Key frame detectors have received a lot of attention in research and many attempts have been made to create a robust key frame detector [11] [12]. An existing video fingerprinting systems that uses a key frame detector is [4]. The use of key frames improve the fingerprinting system drastically, as it reduces the amount of data saved in the database as well as minimising processing time. In Figure 3 a basic diagram is shown of the system.

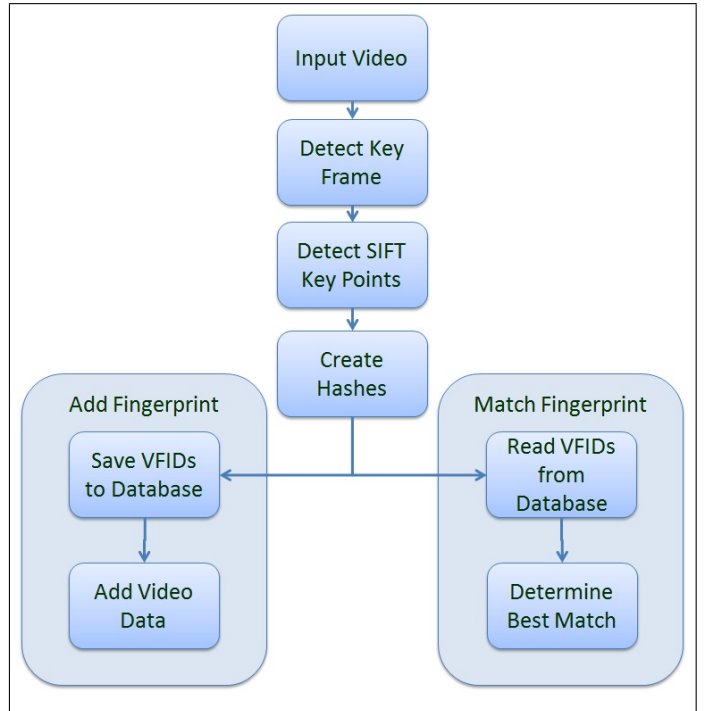


Figure 3: Diagram of a video fingerprinting system

#### A. Key Frame Detector

There are a lot of existing video shot and key frame detectors these days, but none of them retrieve shots and key frames perfectly. For this research a key frame detector is needed that can be run in real-time and it should have a high recall rate, even with some distortion or codecs changes present. A good frame difference measurement to use in this instance would be SIFT, as used in the key frame detector in [11]. The reason for this is that the fingerprinting algorithm also uses SIFT, but the key frame detector can't be run in real time and therefore it can not be used for detection in the system.

The key frames in this instance where detected by using Jensen Shannon Divergence (JSD), proposed by [12]. To quantise the difference of consecutive frames, the JSD of the two frames are calculated. Spikes in the JSD data represent shot boundaries in the video and it can thus be used to split the video into its different shots. Once a shot is detected a key frame can be selected from that shot. The first frame of the shots is selected as the key frame. In [4] it is mentioned that the first frame can be unreliable, but as the focus of this fingerprinting system is currently on advertisement tracking on television it won't be a problem.

To detect a shot boundary, the spikes in the JSD data was detected by comparing a point to the average of the points around it. If the point's value is bigger than the average, multiplied by a certain amount, it is detected as a key frame.

*Test:* Key frame detection is very important for the successful matching of a video's frames. The same frames have to be selected as key frames every time a video is run through the algorithm, so matches will be found in the database. In this section the robustness of the key frame detector is shown.

The recall and precision measurements are used to evaluate information retrieval systems, and can be defined as:

$$\text{precision} = \frac{t_p}{t_p + f_p} \quad (1)$$

$$\text{recall} = \frac{t_p}{t_p + f_n} \quad (2)$$

Where  $t_p$  is the true positives (correct detections),  $f_p$  the false positives (wrong detections) and  $f_n$  the false negatives (missed detections). After running the test on the avi xvid files the following results were found:

Wrapper	Codec	$t_p$	$f_p$	$f_n$	Precision	Recall
avi	XviD	3594	93	89	0.975	0.976
mpg	mpeg1	3568	78	116	0.979	0.968

Table III  
KEY FRAME TEST RESULTS

### B. Video Detection

To match a video to the database, the unknown video stream is analysed for key frames. Once a key frame is found the frame is processed and compared to the database. If the frame was detected successfully, and the number of hashes matched was greater than 30% the video is detected, but if the number of hashes are between 10% and 30%, two consecutive frame matches are needed to confirm a video match.

*Test:* To test the video detection, a day's video was recorded from a television channel. All the advertisements were then fingerprinted and added to the database using 50 hashes per key frame and 16 bit hashes. The video data was then run through the detector. Of the 128 advertisements, 124 were detected successfully, some being detected more than once, when they were repeated in other time slots. In Table IV the results are shown, with the first line only referring to the advertisements and if they were detected or not during the day, whereas the second line represents all the advertisements shown during the day, repeats included.

	$t_p$	$f_n$	$f_p$	Precision	Recall
Advertisements	124	4	0	1	0.969
Advertisements (repeats included)	188	9	0	1	0.954

Table IV  
VIDEO FINGERPRINTING

The advertisements that were missed was because the key frame detector failed to detect sufficient key frames in each of the videos. This is because the key frame detector detects abrupt changes and these advertisements didn't have any.

## IV. CONCLUSION

In this paper a new method for video fingerprinting was proposed, based on the SIFT technique and a Shazam-like hashing method. The test results show the frame fingerprinting algorithms' detection characteristics and meaningful decisions was made that improves the detection effectiveness of the algorithm. A key frame detector algorithm is also used to use frame fingerprinting to in video fingerprinting. The system does function correctly and can detect videos, but future work will include improving the key frame detector and making the frame fingerprinting more robust towards other types of distortions.

## REFERENCES

- [1] "http://www.techrepublic.com/blog/tech-news/youtubes-video-fingerprinting-receives-mixed-reactions/1389."
- [2] D. Pereira and L. Loyola, "Robust video fingerprinting system," in *Communications and Electronics (ICCE), 2010 Third International Conference on*, pp. 141–146, 2010.
- [3] X. Su, T. Huang, and W. Gao, "Robust video fingerprinting based on visual attention regions," in *Acoustics, Speech and Signal Processing, 2009. ICASSP 2009. IEEE International Conference on DOI - 10.1109/ICASSP.2009.4959886*, pp. 1525–1528, 2009.
- [4] C. De Roover, C. De Vleeschouwer, F. Lefebvre, and B. Macq, "Robust video hashing based on radial projections of key frames," *Signal Processing, IEEE Transactions on*, vol. 53, no. 10, pp. 4020–4037, 2005.
- [5] S. Lee and C. Yoo, "Video fingerprinting based on centroids of gradient orientations," in *Acoustics, Speech and Signal Processing, 2006. ICASSP 2006 Proceedings. 2006 IEEE International Conference on DOI - 10.1109/ICASSP.2006.1660364*, vol. 2, pp. II–II, 2006.
- [6] S. Lee and C. Yoo, "Robust video fingerprinting for content-based video identification," *Circuits and Systems for Video Technology, IEEE Transactions on DOI - 10.1109/TCSVT.2008.920739*, vol. 18, no. 7, pp. 983–988, 2008.
- [7] A. Ferman, A. Tekalp, and R. Mehrotra, "Robust color histogram descriptors for video segment retrieval and identification," *Image Processing, IEEE Transactions on DOI - 10.1109/TIP.2002.1006397*, vol. 11, no. 5, pp. 497–508, 2002.
- [8] A. Massoudi, F. Lefebvre, C.-H. Demarty, L. Oisel, and B. Chupeau, "A video fingerprint based on visual digest and local fingerprints," in *Image Processing, 2006 IEEE International Conference on DOI - 10.1109/ICIP.2006.312834*, pp. 2297–2300, 2006.
- [9] A. Wang, "An industrial strength audio search algorithm," in *ISMIR*, pp. 7–13, 2003.
- [10] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 60, p. 91, Nov. 2004.
- [11] G. Liu, X. Wen, W. Zheng, and P. He, "Shot boundary detection and keyframe extraction based on scale invariant feature transform," in *Computer and Information Science, 2009. ICIS 2009. Eighth IEEE/ACIS International Conference on*, pp. 1126–1130, 2009.
- [12] Q. Xu, P. Wang, B. Long, M. Sbert, M. Feixas, and R. Scopigno, "Selection and 3d visualization of video key frames," in *Systems Man and Cybernetics (SMC), 2010 IEEE International Conference on*, pp. 52–59, 2010.

**R. Moolman** is currently pursuing a Master's degree in Computer and Electronic Engineering at the North-West University. He received his Bachelor's degree in Computer and Electronic Engineering in 2010. His current research interests are video fingerprinting, digital signal processing and software design.