

## **Summarisation of short-term and long-term videos using texture and colour**

### Author

Carvajal, Johanna, McCool, Chris, Sanderson, Conrad

### Published

2014

### Conference Title

IEEE Winter Conference on Applications of Computer Vision

### Version

Accepted Manuscript (AM)

### DOI

[10.1109/wacv.2014.6836025](https://doi.org/10.1109/wacv.2014.6836025)

### Rights statement

© 2014IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

### Downloaded from

<http://hdl.handle.net/10072/395923>

### Griffith Research Online

<https://research-repository.griffith.edu.au>

# Summarisation of Short-Term and Long-Term Videos using Texture and Colour

Johanna Carvajal, Chris McCool, Conrad Sanderson

NICTA, GPO Box 2434, Brisbane, QLD 4001, Australia \*  
University of Queensland, School of ITEE, St Lucia, QLD 4072, Australia  
Queensland University of Technology (QUT), Brisbane, QLD 4000, Australia

## Abstract

*We present a novel approach to video summarisation that makes use of a Bag-of-visual-Textures (BoT) approach. Two systems are proposed, one based solely on the BoT approach and another which exploits both colour information and BoT features. On 50 short-term videos from the Open Video Project we show that our BoT and fusion systems both achieve state-of-the-art performance, obtaining an average F-measure of 0.83 and 0.86 respectively, a relative improvement of 9% and 13% when compared to the previous state-of-the-art. When applied to a new underwater surveillance dataset containing 33 long-term videos, the proposed system reduces the amount of footage by a factor of 27, with only minor degradation in the information content. This order of magnitude reduction in video data represents significant savings in terms of time and potential labour cost when manually reviewing such footage.*

## 1. Introduction

Video abstraction aims at providing concise representations of long videos. It has applications in browsing and retrieval of large volumes of videos [1] and also in improving the effectiveness and efficiency of video storage [21]. Video abstraction can be categorised into two general groups: video summarisation and video skimming [10, 21]. Video summarisation, also known as still image abstraction, static storyboard or static video abstract, is a compilation of representative frames selected from the original video [6]. Video skimming, also known as moving image abstraction or moving/dynamic storyboard, is a collection of short video clips [2, 10]. Both approaches should preserve the most important content from the video in order to present a comprehensible and understandable description for the end user.

In general, video skimming provides a more coherent and visually attractive result. It often retains a high-level of linguistic meaning due to its capacity to combine audio

and moving elements [14, 21]. However, video summarisation is easier to generate and is not constrained in terms of timing and synchronisation [2, 21].

Video summarisation is an active area of research within the computer vision community and it has been applied in various video categories such as Wildlife Videos [23], sports videos [15], TV documentaries [2], among others. In [1] the various approaches to video summarisation are divided into six techniques consisting of: feature selection, clustering algorithms, event detection methods, shot selection, trajectory analysis and the use of mosaics. Often a combination of techniques is used, for example one of the most common approaches is to combine feature selection with a form of clustering [2, 6, 13].

In [24] a video summary is obtained by extracting a feature vector from each frame and then clustering the resulting set of feature vectors. The smallest clusters are then removed. A keyframe – a frame that forms part of the video summary – is selected for each cluster centroid by taking the frame whose feature vector is closest to the centroid. Similar approaches are adopted in [2, 6, 7, 10] where the major difference is in the choice of feature vector used to represent each frame. Colour histograms are used in [2, 6], motion-based features are used in [7], and saliency maps are used in [10]. Each of the previously proposed feature vectors has its drawbacks. For instance, the colour histogram approach used in [2, 6] retains only coarse information about the frame. Motion-based features of [7] fail when the motion in the videos is too large. Finally, the saliency maps used in [10] perform poorly for cluttered and textured backgrounds. To date, limited work has been done on incorporating texture information to perform video summarisation.

**Contributions.** In this paper we first propose the use of texture information to improve video summarisation. We propose the use of the computationally efficient and effective bag-of-textures approach; we conjecture that this will improve video summarisation as it has been successfully applied to a range of image processing tasks, such as matching and classification of natural scenes and faces [12, 19, 22]. The bag-of-textures model divides an image into small patches, extracts appearance descriptors from each patch, quantises each descriptor into a discrete “visual word”, and

---

\***Acknowledgements:** NICTA is funded by the Australian Government through the Department of Communications and the Australian Research Council through the ICT Centre of Excellence Program.

then computes a compact histogram representation [8], providing considerably different information than colour histograms. In addition, we propose a fusion based system for video summarisation, where both colour and texture information is exploited. This will allow us to overcome the shortcomings of either approach. Similar approaches have been shown to be advantageous in object classification tasks [11]. We show that our system may be applied not only to short-term videos but also to long-term videos, helping in the detection of the existence of a rare species of fish.

The layout of this paper is as follows. In Section 2 we describe in detail our proposed video summarisation method that exploits the benefits of using texture histograms based on the bag-of-textures model. In Section 3 we present our improved video summarisation method that fuses the visual information provided by both the colour and texture histograms. In Section 4 we describe how we evaluate the video summaries of short-term and long-term videos. In Section 5, we present experiments which show that the proposed methods obtain higher performance than existing methods based on colour histograms. Section 6 summarises the main findings.

## 2. Bag-of-Textures for Video Summarisation

This section describes our proposed bag-of-textures (BoT) approach. There are four main stages:

1. Pre-processing: The input video is sub-sampled after which each frame is filtered and rescaled.
2. BoT representation:
  - (i) *Local Texture Features*. Each frame is divided into small patches (blocks) and from each block we extract 2D-DCT features, which is an effective and compact representation [16].
  - (ii) *Dictionary Training*. A generic visual dictionary is trained to describe the most commonly occurring textures in an independent training set.
  - (iii) *Generation of BoT Histogram*. Each frame is represented by a histogram which is obtained by matching the feature vectors from each block to the dictionary.
3. Keyframe selection: Similar frames are grouped into an automatically determined number of clusters. One keyframe is selected per cluster.
4. Post-processing: In this final stage, we eliminate possible repetitive frames and create the static video summary.

Each of these stages is elucidated in the following sections.

### 2.1. Pre-processing

#### 2.1.1 Sampling and Rescaling

The original input video is re-sampled to one frame per second in order to reduce the number of video frames to be examined. Each frame is then converted into gray-scale and re-scaled to be a quarter of its original size, in order to reduce the computational cost of the following stages.

#### 2.1.2 Noise Filtering

There are often uninformative frames that appear at the beginning and/or the end of a segment that may affect the appearance of a video summary [6]. These frames are usually colour-homogeneous due to fade-in and fade-out effects, and have a small standard deviation of their pixel values. Frames with a standard deviation below a threshold are eliminated.

### 2.2. BoT Representation

#### 2.2.1 Local Texture Features

Each frame is divided into  $N$  overlapping blocks. To each block we apply the 2D discrete cosine transform (2D-DCT) to obtain a  $D$ -dimensional feature vector that represents the local texture information [16]. Thus, the local texture feature for the  $n$ -th block of the  $i$ -th frame is  $\mathbf{x}_{i,n}$ .

#### 2.2.2 Dictionary Training

The dictionary is trained using the  $k$ -means algorithm [3] by pooling the local texture features from a set of training frames. The resulting  $G$  cluster centers  $\{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_G\}$  represent the local textures (codewords) of the dictionary.

#### 2.2.3 Generation of BoT Histogram

In the BoT approach the  $i$ -th frame is represented by a histogram,  $\mathbf{h}_i^{\text{BoT}}$ . This  $G$ -dimensional histogram represents the relative frequency of the local texture features within the frame. The  $g$ -th dimension of  $\mathbf{h}_i^{\text{BoT}}$  is the relative frequency of the  $g$ -th local texture feature from the dictionary, similar to [5]. The histogram is normalised to sum to one. Thus, each local texture feature can be converted to a local histogram,  $\mathbf{h}_{i,n}^{\text{BoT}}$ , of dimension  $G$  where each dimension  $g$  is given by,

$$h_{g,i,n}^{\text{BoT}} = \begin{cases} 1 & \text{if } g = \arg \min_{k \in \{1, \dots, G\}} \|\mathbf{x}_{i,n} - \boldsymbol{\mu}_k\|_2 \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

These  $N$  local histograms can then be summed and normalised to produce the final BoT histogram,

$$\mathbf{h}_i^{\text{BoT}} = \frac{1}{N} \sum_{n=1}^N \mathbf{h}_{i,n}^{\text{BoT}}. \quad (2)$$

### 2.3. Keyframe Selection

To obtain a set of keyframes we adopt an approach similar to that of [6]. A keyframe is a frame that forms part of the video summarisation. The  $k$ -means algorithm is used to cluster similar frames into  $K$  segments, and the resultant centroids are then used to select the keyframes.

Initially, the frames are grouped consecutively, assuming that sequential frames share similar content. To automatically determine the number of clusters,  $K$ , we calculate the Euclidean distance between two consecutive frames. If the distance is greater than a threshold  $\tau$  then  $K$  is incremented. For each cluster centroid the frame whose BoT histogram is closest is selected as a keyframe. A total of  $K$  keyframes is then reached.

### 2.4. Post-processing

Having obtained the initial set of  $K$  keyframes we then attempt to discard those keyframes which are too similar. This is achieved by comparing all keyframes against each other. If the Euclidean distance between the BoT histograms of the keyframes is smaller than a threshold  $\tau$  then one of the two keyframes under consideration is discarded. This gives the final static video summary that consists of  $N_{as}$  keyframes, where  $N_{as} \leq K$ , with  $as$  standing for automatic summary.

Lastly, the static video summary is obtained after organising the resulting keyframes in temporal order.

### 3. Fusion of Colour and BoT

In this section, we present a hybrid system that fuses colour histograms [6] and BoT texture information, termed as CaT (for Colour and Texture). The proposed CaT approach to video summarisation has the same 4 stages as our proposed BoT video summarisation approach, but with additions in order to obtain colour histograms. We describe these additions below.

1. Pre-processing: The input video is processed in two independent ways. First, we obtain the BoT histograms as described in Section 2.1. Second, to obtain the colour histograms we extract the Hue component, from the HSV colour space, of the unscaled input frame similar to [6]. In both cases we remove uninformative frames by employing the noise filtering process described in Section 2.1.
2. Texture and Colour Histogram: The BoT histogram is the same as explained in Section 2.2. The colour histogram,  $\mathbf{h}_i^{\text{hue}}$ , of the  $i$ -th frame is computed using only the Hue component as in [6].
3. Keyframe Selection: The BoT and colour histograms are clustered using  $k$ -means. This stage is similar to

Section 2.3. The difference lies in the distance measure used to compare all frames against each other.

- (i) To select the number of keyframes  $K$  we combine the information from the BoT and colour histograms. When calculating the distance between frame  $a$  and  $b$  we use the weighted summation of Euclidean distances:

$$\alpha \|\mathbf{h}_a^{\text{BoT}} - \mathbf{h}_b^{\text{BoT}}\|_2 + \beta \|\mathbf{h}_a^{\text{hue}} - \mathbf{h}_b^{\text{hue}}\|_2 \quad (3)$$

under the constraints  $\alpha + \beta = 1$ ,  $\alpha \geq 0$ ,  $\beta \geq 0$ .

- (ii) Each keyframe is selected by finding the frame which is closest to each cluster centroid. For the CaT approach the distance between a frame and a centroid is calculated as a weighted summation of the Euclidean distances, as per (3).

4. Post-processing: To eliminate similar frames we use the procedure described in Section 2.4 but replace the Euclidean distance with the weighted summation of the Euclidean distances, as per (3).

## 4. Datasets and Evaluation Metrics

To evaluate the performance of video summarisation we use two datasets consisting of short- and long-term video data. The short-term data is obtained from the Open Video Project<sup>1</sup>. The long-term data is a new dataset that consists of 14 hours of underwater video surveillance which monitors the behaviour of marine wildlife.

### 4.1. Short-Term Videos

We use the 50 videos from the Open Video Project which contain ground truth [6]. Each ground truth consists of the summary provided by  $P = 5$  users. The users provided the summaries under no restrictions upon length nor appearance of the summaries.

To evaluate the performance on the short-term video data we use the ‘‘Comparison of User Summaries’’ (CUS) method [6]. This method compares the automatic video summarisation and ground truth by exhaustively calculating the distance between the frames from the automatic summarisation and the ground truth. Two frames are similar if the distance between their respective feature vectors (histograms) is less than an evaluation threshold  $\delta$ . If the frames match they are removed from the next iteration of the comparison process. For performance evaluation, the distance measure used for the BoT approach is the Euclidean distance, however, to be consistent with prior work [6], the distance measure for the colour histograms is the  $L_1$ -norm. Therefore, the distance measure used for CaT

<sup>1</sup>Open Video Project: <http://www.open-video.org>

is the weighted summation of the Euclidean distance for the BoT histograms and the  $L_1$ -norm for the colour histograms:

$$\alpha \|\mathbf{h}_a^{\text{bof}} - \mathbf{h}_b^{\text{bof}}\|_2 + \beta \|\mathbf{h}_a^{\text{hue}} - \mathbf{h}_b^{\text{hue}}\|_1. \quad (4)$$

Various evaluation metrics exist to measure the quality of an automatic video summary. We use three evaluation metrics so that we can compare our proposed approaches with two state-of-the-art methods [6, 2]. To compare with [6] we use accuracy ( $acc$ ) and error ( $err$ ), and to compare with [2] we use the  $F$ -measure.

To calculate  $acc$  and  $err$ , each frame in the automatic video summary is compared with all frames in the user summary and then the number of matching frames ( $N_m$ ) and non-matching frames ( $N_{nm}$ ) are calculated:

$$acc = \frac{N_m}{N_u}, \quad err = \frac{N_{nm}}{N_u} \quad (5)$$

where  $N_{as}$  and  $N_u$  are the total number of frames from the automatic and user summary, respectively.

The  $F$ -measure, defined as

$$F = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \quad (6)$$

is used to provide a single number that balances precision =  $N_m/N_{as}$  and recall =  $N_m/N_u$ .

The evaluation metrics are presented as an average. First, we take the average from the  $P$  users to obtain  $acc_P$ ,  $err_P$ , and  $F_P$ ; for each video there are  $P = 5$  users. Then we take the average across all of the videos to obtain  $\overline{acc}$ ,  $\overline{err}$ , and  $\overline{F}$ . In terms of  $\overline{acc}$  it is desirable to have a high value as it measures the number of matching frames. In terms of  $\overline{err}$  it is desirable to have a small value as it measures the number of non-matching frames. With regards to  $\overline{F}$  it is desirable to obtain a high value, which occurs when the precision and recall are large.

## 4.2. Long-Term Videos

The long-term videos consist of 14 hours of underwater footage from 33 videos which are on average 25 minutes in duration. This data was obtained from the NSW-DPI<sup>2</sup>, courtesy of David Harasti. Example images are shown in Figure 1. In each video there is always at least one segment where a rare species of fish, the black cod, is within view. Normally these videos would be inspected by a human expert to determine if there is an instance of the rare fish within. We propose that video summarisation can be used to reduce the amount of footage to be viewed in order to detect the existence of this rare species of fish.

Using ground truth which provides time-stamps when this rare species is within view, we examine the effectiveness of video summarisation to provide at least one



**Figure 1.** Example images from the long-term underwater surveillance videos; the added red ellipsoids highlight the rare species of interest.

keyframe in each static video summary with the rare species of interest within view. This is useful as it presents a way to reduce the time and cost of manually viewing a large amount of video data.

To calculate the performance of long-term videos we present results in terms of detection accuracy and the average compression ratio ( $R_c$ ). Detection accuracy refers to whether an instance of the rare species is among any of the chosen keyframes for a static video summary; 75% would mean that there is at least 1 keyframe of the rare species in 75% of the static video summaries.

To calculate the average compression ratio we first note that because we have long-term videos then for each video there might be many hundreds of keyframes. To present all of these keyframes effectively to the user we re-encode them into a static video summary by presenting each keyframe for 0.25 seconds. This gives the user time to effectively view the keyframe. Thus the  $t$ -th long-term video  $\mathbf{V}_t$  is converted to a static video summary  $\mathbf{S}_t$  with a compression ratio given by:

$$R_{c,t} = 4 \times \frac{\text{Duration}(\mathbf{V}_t)}{\text{Duration}(\mathbf{S}_t)} \quad (7)$$

where Duration is the duration of a video and the factor of 4 is introduced as there are 4 keyframes per second of the shortened video.

<sup>2</sup>New South Wales Department of Primary Industries, Australia.



## 5. Experiments

An important part of both the BoT and CaT approaches is the training of the dictionary to obtain the texture histograms. To train this dictionary we use 10 frames randomly selected from videos taken from the Open Video Project that have no user summaries, ensuring they are independent of the evaluation dataset. In addition, the frames selected to train the dictionary look significantly different to the ground truth provided by the users.

To obtain the proposed local texture features we divide each frame into a set of overlapping blocks. Similar to [19] we use a block size of  $8 \times 8$  with an overlap margin of 6 pixels, and represent each block as a  $D = 15$  dimensional feature vector containing 2D-DCT coefficients. We extract the first 16 2D-DCT coefficients, which represent low-frequency information [16], and omit the first coefficient as it is the most sensitive to illumination changes. With regards to the colour histogram, we quantise the Hue component into 16 bins as per [6]. These parameters are the same for all experiments.

The values for the threshold  $\tau$ , fusion weight  $\alpha$  and evaluation threshold  $\delta$  were determined experimentally. For all of the experiments we search for the optimal fusion parameter  $\alpha = \{0.0, 0.1, \dots, 1.0\}$ . Our proposed methods were implemented using the OpenCV [4] and Armadillo [18] C++ libraries.

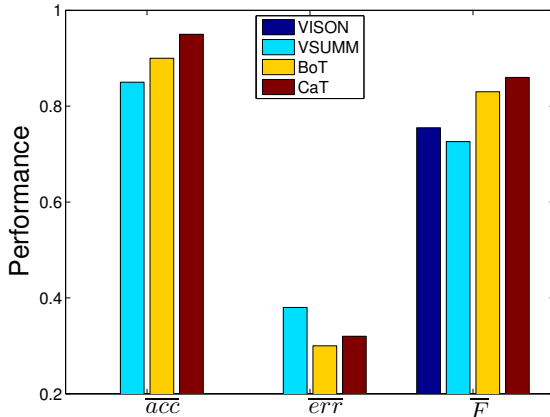
### 5.1. Short-Term Videos

We compare the performance against two baseline systems from literature, VSUMM [6] and VISON [2]. The two baseline systems both use colour information as their primary feature. VSUMM uses colour information by retaining only the Hue component of HSV and generating a histogram of 16 bins. VISON is a state-of-the-art approach and consists of a histogram of the HSV representation of each frame. It combines the HSV information in a compressed form such that the Hue component is treated with greater importance and results in a histogram of 256 bins.

An initial set of experiments were performed to find the optimal number of components for the dictionary of our proposed texture features. Using a fixed number of components  $G = \{8, 16, 32\}$  and a fixed number of thresholds  $\tau = \{0.05, 0.10, \dots, 0.5\}$ , we found that using just  $G = 8$  components provided optimal performance. We kept the number of components constant for the remainder of our experiments.

In Figure 2 we present a summary of the average performance for 50 short-videos of our proposed systems, BoT and CaT, and the two baselines. Two interesting results can be seen from this figure.

First, it can be seen that the texture-only BoT system performs better than either the VSUMM or VISON approaches which primarily use colour information. The BoT system



**Figure 2.** Comparative evaluation of our proposed methods with VSUMM [6] and VISON [2]. Lower values of  $\overline{err}$  as well as higher values of  $\overline{acc}$  and  $\overline{F}$  are desired.

obtains an average  $F$ -measure of  $\overline{F} = 0.83$ , which is a relative improvement of 9% when compared to VISON,  $\overline{F} = 0.76$ . Furthermore, the  $\overline{acc}$  and  $\overline{err}$  of the BoT system shows that it produces a more accurate summarisation than VSUMM and also has the lowest  $\overline{err}$  of any system<sup>3</sup>. This suggests that texture information is either equally or more important than colour information for the task of video summarisation.

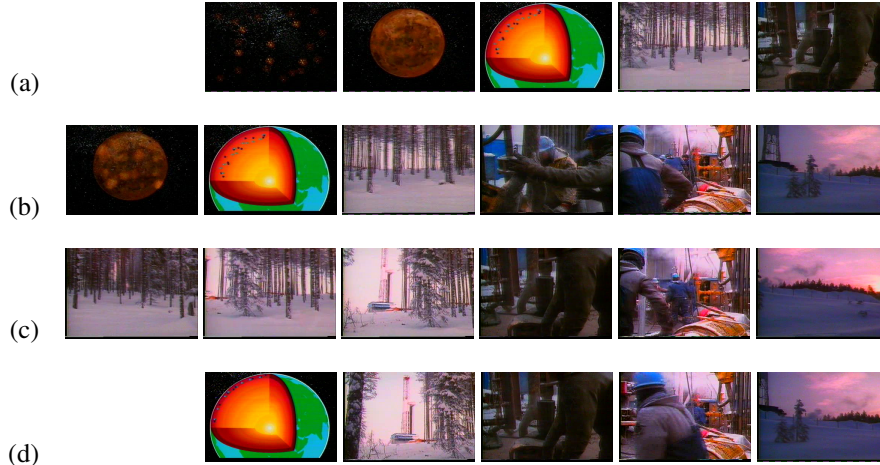
Second, the proposed CaT system (fusing colour histograms and the proposed texture histograms) performs better than the two baseline systems and the proposed texture-only BoT system. The CaT system has an average  $F$ -measure of  $\overline{F} = 0.86$ , which is a relative improvement of 13% when compared to VISON  $\overline{F} = 0.76$ , the previous state-of-the-art approach.

Figure 3 shows the qualitative results for the automatic summarisation provided by VSUMM and VISON as well as our proposed BoT and CaT systems. It can be seen that VSUMM (Figure 3a) with  $F_P = 0.83$ , VISON (Figure 3b) with  $F_P = 0.78$ , and our proposed BoT (Figure 3c) with  $F_P = 0.74$  contain some keyframes that may not be of interest and/or are repetitive. In contrast, the proposed CaT system (Figure 3d) provides the most consistent video summary with  $F_P = 0.86$ .

### 5.2. Long-Term Videos

In this section we present results on 33 long-term videos which last on average for 25 minutes. We examine the applicability of video summarisation to long-term videos to efficiently detect a rare species of fish and measure performance in terms of detection accuracy and compression rate (see Section 4.2).

<sup>3</sup>No results in terms of  $\overline{acc}$  and  $\overline{err}$  were supplied for VISON in [2].



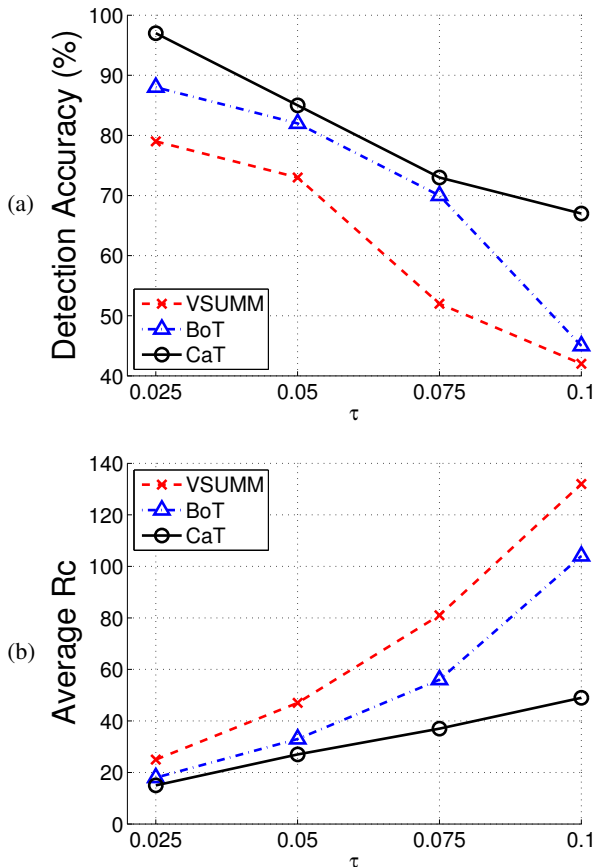
**Figure 3.** Static video summary for “the future of energy gases - segment 09”, using (a) VSUMM, (b) VISON, (c) proposed BoT, and (d) proposed CaT.

The accuracy and average compression ratio of the algorithm for various thresholds,  $\tau = \{0.025, 0.05, \dots, 0.1\}$ , is presented in Figure 4. It can be seen in Figure 4a that the CaT algorithm consistently outperforms the BoT and VSUMM algorithms. We attribute this to the fact that the background in these videos is relatively stable and so the colour histograms used in VSUMM do not change as often compared to the short-term videos used in [6]. In Figure 4b it can be seen that while using the VSUMM algorithm provides better average compression ratio than either the BoT or CaT approaches, it comes at the cost of accuracy. In general the proposed fusion approach provides the most consistent trade-off between accuracy and average compression ratio.

We take the optimal system at the threshold  $\tau = 0.05$  as this provides a high degree of detection accuracy, 85%, and a good average compression ratio of 27. This system will allow a user to see the fish of interest in 85% of the summarised videos while reducing the amount of video data to view by 27 times, more than an order of magnitude. Such an approach would reduce the 14 hours of video data to just 31 minutes, thus enabling significantly more efficient reviewing of the data.

## 6. Summary and Future Work

In this paper, we have proposed the novel use of textures to perform video summarisation. We proposed to use a visual-bag-of-textures (BoT) in two ways. First, a BoT system which uses only texture features is proposed and it is shown to outperform two state-of-the-art systems which use colour only, VSUMM and VISON. Second, a fused system that combines Colour and Texture (CaT) is proposed and it is shown to provide further improvements.



**Figure 4.** Demonstration of the trade-off between (a) the detection accuracy and (b) the average compression ratio  $R_c$  for the 33 long-term videos using the CaT, BoT and VSUMM approaches.

Both of our proposed systems outperform two state-of-the-art approaches, VSUMM and VISON, which use colour features. Experiments on 50 short-term videos, obtained from the Open Video Project, show that our proposed texture-only system (BoT) obtains an  $F$ -measure of 0.83, which is better than either VSUMM or VISON which obtain an average  $F$ -measure of 0.73 and 0.76, respectively. Furthermore, our fused system (CaT) demonstrates that combining colour and texture features yields state-of-the-art performance with an average  $F$ -measure of 0.86.

We have also shown that video summarisation can be applied effectively to long-term videos. Using 33 long-term surveillance videos, in our case underwater surveillance footage, we have shown that video summarisation can be used to significantly reduce the amount of footage to view, by up to a factor of 27, with only a minor degradation in the information content.

Future work should examine alternative features and application settings with a particular emphasis for long-term videos. For instance, emphasising the importance of foreground objects [17] should be explored, as well as explicit modelling of movement (or actions) of such objects [9, 20]. Moreover, the applicability of video summarisation to CCTV surveillance footage should also be considered.

## References

- [1] M. Ajmal, M. Ashraf, M. Shakir, Y. Abbas, and F. Shah. Video summarization: techniques and classification. In *Lecture Notes in Computer Science, Vol. 7594*, pages 1–13. 2012.
- [2] J. Almeida, N. J. Leite, and R. da S. Torres. VISON: Video Summarization for ONLINE applications. *Pattern Recognition Letters*, 33(4):397–409, 2012.
- [3] C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- [4] G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000.
- [5] N. Dardas, Q. Chen, N. D. Georganas, and E. Petriu. Hand gesture recognition using bag-of-features and multi-class support vector machine. In *IEEE Int. Symp. Haptic Audio-Visual Environments and Games (HAVE)*, pages 1–5, 2010.
- [6] S. E. F. de Avila, A. P. B. Lopes, A. da Luz Jr., and A. de Albuquerque Araújo. VSUMM: a mechanism designed to produce static video summaries and a novel evaluation method. *Pattern Recognition Letters*, 32(1):56–68, 2011.
- [7] A. Divakaran, K. A. Peker, and H. Sun. Video summarization using motion descriptors. In *Proc. SPIE Conf. on Storage and Retrieval from Multimedia Databases*, 2001.
- [8] K. Grauman and B. Leibe. Visual object recognition. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 5(2):1–181, 2011.
- [9] M. T. Harandi, C. Sanderson, S. Shirazi, and B. C. Lovell. Kernel analysis on Grassmann manifolds for action recognition. *Pattern Recognition Letters*, 34(15):1906–1915, 2013.
- [10] Q.-G. Ji, Z.-D. Fang, Z.-H. Xie, and Z.-M. Lu. Video abstraction based on the visual attention model and on-line clustering. *Signal Processing: Image Communication*, 28(3):241–253, 2013.
- [11] Z. Li, Y. Liu, R. Hayward, and R. Walker. Color and texture feature fusion using kernel PCA with application to object-based vegetation species classification. In *IEEE International Conference on Image Processing (ICIP)*, pages 2701–2704, 2010.
- [12] Z. Lin and J. Brandt. A local bag-of-features model for large-scale object retrieval. In K. Daniilidis, P. Maragos, and N. Paragios, editors, *European Conference on Computer Vision (ECCV)*, volume 6316 of *Lecture Notes in Computer Science*, pages 294–308. Springer Berlin Heidelberg, 2010.
- [13] P. Mundur, Y. Rao, and Y. Yesha. Keyframe-based video summarization using Delaunay clustering. *International Journal on Digital Libraries*, 6(2):219–232, 2006.
- [14] J. Oh, Q. Wen, J. Lee, and S. Hwang. Video abstraction. In *Video Data Management and Information Retrieval*, pages 321–346. Idea Group Inc. and IRM Press, 2004.
- [15] J.-Q. Ouyang and R. Liu. Ontology reasoning scheme for constructing meaningful sports video summarisation. *IET Image Processing*, 7(4):324–334, 2013.
- [16] W. B. Pennebaker and J. L. Mitchell. *JPEG Still Image Data Compression Standard*. New York: Van Nostrand Reinhold, 1993.
- [17] V. Reddy, C. Sanderson, and B. C. Lovell. Improved foreground detection via block-based classifier cascade with probabilistic decision integration. *IEEE Transactions on Circuits and Systems for Video Technology*, 23(1):83–93, 2013.
- [18] C. Sanderson. Armadillo: an open source C++ linear algebra library for fast prototyping and computationally intensive experiments. Technical report, NICTA, 2010.
- [19] C. Sanderson and B. C. Lovell. Multi-region probabilistic histograms for robust and scalable identity inference. In *Lecture Notes in Computer Science (LNCS), Vol. 5558*, pages 199–208, 2009.
- [20] A. Sanin, C. Sanderson, M. Harandi, and B. C. Lovell. Spatio-temporal covariance descriptors for action and gesture recognition. In *Workshop on the Applications of Computer Vision (WACV)*, pages 103–110, 2013.
- [21] B. T. Truong and S. Venkatesh. Video abstraction: a systematic review and classification. *ACM Trans. Multimedia Comput. Commun. Appl.*, 3(1), Feb. 2007.
- [22] J. Yang, K. Yu, Y. Gong, and T. Huang. Linear spatial pyramid matching using sparse coding for image classification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1794–1801, 2009.
- [23] S.-P. Yong, J. Deng, and M. Purvis. Key-frame extraction of wildlife video based on semantic context modeling. In *International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, 2012.
- [24] Y. Zhuang, Y. Rui, T. Huang, and S. Mehrotra. Adaptive key frame extraction using unsupervised clustering. In *IEEE International Conference on Image Processing*, volume 1, pages 866–870, 1998.