

A mid-level video representation based on binary descriptors: A case study for pornography detection

Carlos Caetano^{a,b,*}, Sandra Avila^c, William Robson Schwartz^b, Silvio Jamil F. Guimarães^d, Arnaldo de A. Araújo^a

^aUniversidade Federal de Minas Gerais, NPDI — DCC/UFMG, Minas Gerais, Brazil

^bUniversidade Federal de Minas Gerais, SSIG Group — DCC/UFMG, Minas Gerais, Brazil

^cUniversity of Campinas, RECOD Lab — DCA/FEEC/UNICAMP, Campinas, Brazil

^dPontifical Catholic University of Minas Gerais, VIPLAB — ICEI/PUC Minas, Minas Gerais, Brazil

Abstract

With the growing amount of inappropriate content on the Internet, such as pornography, arises the need to detect and filter such material. The reason for this is given by the fact that such content is often prohibited in certain environments (e.g., schools and workplaces) or for certain publics (e.g., children). In recent years, many works have been mainly focused on detecting pornographic images and videos based on visual content, particularly on the detection of skin color. Although these approaches provide good results, they generally have the disadvantage of a high false positive rate since not all images with large areas of skin exposure are necessarily pornographic images, such as people wearing swimsuits or images related to sports. Local feature based approaches with Bag-of-Words models (BoW) have been successfully applied to visual recognition tasks in the context of pornography detection. Even though existing methods provide promising results, they use local feature descriptors that require a high computational processing time yielding high-dimensional vectors. In this work, we propose an approach for pornography detection based on local binary feature extraction and BossaNova image representation, a BoW model extension that preserves more richly the visual information. Moreover, we propose two approaches for video description based on the combination of mid-level representations namely BossaNova Video Descriptor (BNVD) and BoW Video Descriptor (BoW-VD). The proposed techniques are promising, achieving an accuracy of 92.40%, thus reducing the classification error by 16% over the current state-of-the-art local features approach on the Pornography dataset.

Keywords: Binary descriptors, mid-level representation, Bag-of-Words, BossaNova, pornography

1. Introduction

Due to the fast growth of images and videos publicly available on the Internet, the need for recognition of their contents arises. Besides the obvious need for methods related to image and video searches, it is also important to perform recognition or classification of contents that may be considered undesirable or offensive to allow the development of methods for filtering them.

The largest group of images and videos available on the Internet that people may find offensive is related to pornographic materials. A report by the ExtremeTech¹ technology site suggests that 30% of all Internet traffic is associated with pornography. They arrived at this number by estimating the traffic that a popular pornographic website generates every day and multiplied it by several other pornographic websites of similar size found on the Internet. According to their report, the

largest website provider of this type of content receives three times more pageviews than major news websites (about 4.4 billion pageviews per month) and the average time spent on this site can be five times higher than in news sites.

According to Short et al. [1], pornography can be considered as any sexually explicit material with the aim of sexual arousal or fantasy. However, this definition leads to many challenges when trying to detect pornographic content, such as the bounds of “explicit” for something to be considered as pornographic material. Some works in the literature deal with this issue by dividing the material into several classes [2], complicating even more the classification task. On the other hand, there are works that choose to deal with it by using a conceptually simple evaluation considering only two classes (pornographic and non-pornographic) [3, 4], the focus of this work.

Detecting and filtering pornographic visual content from the Internet is a concern in many environments (e.g., schools, workplaces). According to Lopes et al. [5], linked text tags to pictures and videos are clearly not sufficient, since inappropriate content can be maliciously attached to seemingly innocent texts. A typical situation would be, for example, the employment of search keywords commonly used by children attached to websites with pornographic content. In addition, adults may also not wish to be exposed to such contents, for instance, from

*Corresponding author.

Email addresses: carlos.caetano@dcc.ufmg.br (Carlos Caetano), sandra@dca.fee.unicamp.br (Sandra Avila), william@dcc.ufmg.br (William Robson Schwartz), sjamil@pucminas.br (Silvio Jamil F. Guimarães), arnaldo@dcc.ufmg.br (Arnaldo de A. Araújo)

¹<http://www.extremetech.com/computing/123929-just-how-big-are-porn-sites>

results received from search engines available on the web.

In recent years, several works in literature have been mainly focused on detecting pornographic images and videos based on visual content rather than textual information [2–4, 6–18]. Most of these works are based on skin color detection approaches since a large fraction of pixels that have colors related to the human skin [19]. Nevertheless, a shortcoming of these approaches is related to the high rate of false positives, since not all images with large areas of skin exposure are necessarily pornographic (pictures of people wearing swimsuits, or sports-related images). Furthermore, another issue to be considered is that grayscale pictures cannot be classified using color related features.

The pornography detection task can be interpreted as a visual recognition task in the context of object recognition [5]. Approaches based on local features in conjunction with Bag-of-Words models (BoW) have been successfully applied to visual classification tasks [20, 21]. In such approaches, images are represented as histograms constructed from a set of visual features. No explicit model of the object is needed and the variability of examples (related to rotation, shape scale or illumination) is treated by a training set that includes such variability. In view of that, approaches based on BoW models are suitable to the task of pornography detection.

Despite the existing methods based on BoW models produce promising results in the pornography detection context, these also make use of local feature descriptors that require a high computational processing time and generate high-dimensional real-valued vectors. For example, Avila et al. [4] made use of HueSIFT feature descriptor [22], a variant of SIFT descriptor [23] that includes color information, taking an average time of 2.5 seconds to densely extract the local features of an image generating a feature vector consisting of 165 floating point values. In fact, this is still not fast enough for real-time applications, that require a short response time. Moreover, the comparison between two extracted features would spend more computational time due to the high dimensionality. On the other hand, to satisfy the requirements of web pornographic image recognition both on precision and speed, Zhuo et al. [18] proposed a pornographic image recognition method based on the binary descriptor Oriented FAST and Rotated BRIEF (ORB), which is a low-complexity alternative. However, their work focused only on static images.

In this paper, we formalize a video descriptor approach to the visual recognition problem in the context of pornography detection in videos. The method is based on both a low-complexity alternative for feature extraction using binary descriptors and a combination of mid-level representations. We apply it to the classical BoW model generating the BoW Video Descriptor (BoW-VD). We also apply it to the BossaNova, a BoW model extension that preserves the visual information in a richer way, which generates the BossaNova Video Descriptor (BNVD). Our proposal has as advantage the fact that it does not depend on any skin detector or shape models to classify pornography; besides, according to the experimental results, it outperforms the state-of-the-art results on the Pornography dataset [4]. *To the best of our knowledge, ours is the best result reported to date on the*

Pornography dataset employing local feature descriptors.

The use of binary descriptors and the mid-level representation for videos were first introduced in our previous works [17] and [24]. This paper presents several new aspects in comparison with the previous ones. Those aspects are highlighted in the following:

- Formalization of BossaNova Video Descriptor (BNVD). In this work, we present a new formulation which generalizes the BNVD allowing the use of different aggregation functions.
- Proposal of BoW Video Descriptor (BoW-VD). In this work, we also propose a video descriptor by using aggregation functions for combining the traditional BoW mid-level representations.
- Improvement of the experimental results. In this work, we study the behavior of binary descriptors and the mid-level representation by using several parameter settings, including the use of global pooling for creating a video descriptor. Moreover, we show the computational times for generating our video descriptor.

The remainder of this paper is organized as follows. We start by explaining the classical non-binary local feature descriptors, the most recent binary feature descriptors and the BossaNova mid-level representation (Section 2). Next, we survey the recent works on pornography detection (Section 3). We then introduce the complete formalism of our video descriptor (Section 4). Afterwards, we analyze our experiments regarding the proposed video descriptor and we perform a comparison with state-of-the-art approaches (Section 5). Finally, we present our concluding remarks (Section 6).

2. Theoretical Background

The most common approach for visual recognition task consists of three distinct steps [25]: (i) extraction of local features; (ii) encoding of the local features in an intermediate representation (mid-level); and (iii) classification of the mid-level representation, usually based on machine learning techniques. Typically, the extracted local features tend to be invariant to some transformations caused by camera changes such as rotation, scale and illumination. To address these properties, local features such as SIFT [23] or SURF [26] descriptors are usually extracted. Regarding the mid-level representation, BoW models [27] are the most common approaches used to encode the extracted local features. Moreover, to improve the efficiency of retrieval and classification without sacrificing accuracy, there are several methods [28–32] that can be applied to the local features or mid-level representations in order to convert them into compact similarity-preserving binary codes. Finally, the purpose of the classification is to learn a function able to assign discrete labels to images/videos. To that end, most of the visual recognition works make use of machine learning techniques, such as Support Vector Machines (SVM).

2.1. Local Feature Descriptors

A local feature descriptor can be considered as a function applied to a region of the image to perform its description. The simplest way to describe a region is to represent all the pixels in this region in a single vector. However, depending on the information to be described, this would result in a high-dimensional vector leading also to a high computational complexity for a future recognition of this region [33].

In this section, we review local descriptors, which can be classified in two distinct ways [34]: (i) non-binary descriptors and (ii) binary descriptors. It is important to say that new approaches for local descriptors have been proposed in the literature, so the following list is not an exhaustive list. However, it can be considered as a representative group of the most relevant descriptors for our context.

2.1.1. Non-Binary Descriptors

SIFT – Scale Invariant Feature Transform

One of the most important descriptors used in the literature is the SIFT [23]. This descriptor performs a scale-space analysis leading to a great performance according to the scale invariance [35]. Although the author has developed the SIFT descriptor to be used on object recognition tasks, it has become the most widely used descriptor in several other applications. This is due to its high discriminative power and stability.

To describe each patch, an orientation α is assigned selecting the angle that represents the histogram of local gradients (calculated for each pixel around the keypoint). Then, the region of points around the keypoint, oriented by α , is divided into subregions composed by a grid of size $G \times G$. Next, a histogram of orientation consisting of B bins is created from the samples of each subregion. The descriptor is then obtained from the concatenation of the histograms of these subregions, composed of $G \times G \times B$ values. The default values for G and B are usually 4 and 8, respectively, resulting in a vector of 128 length. Finally, the descriptor is normalized turning it robust to illumination variations.

SURF – Speeded-Up Robust Features

To overcome the problem of high computational processing time of SIFT, Bay et al. [26] proposed a faster descriptor, called SURF. It can be seen as an approximation of SIFT and it has the same idea of using histograms based on local gradients. SURF is based on integral images to approximate convolutions, which provides a considerable improvement in efficiency (as compared to SIFT). Despite the approximations on the descriptor creation, there is no considerable loss in the rotation and scale invariance.

Similarly to SIFT, SURF assigns an orientation to each patch described: a circular region around the keypoint is described according to the distribution of responses received by a Haar-wavelet filter. The size of the wavelet region and the sampling parameter are dependent of a scale σ in which the keypoint is detected. Weighted with a Gaussian function around the keypoint, the filter responses are represented by vectors in a two-dimensional space and then summed up. The dominant orientation determines the orientation of the keypoint. Then, the

patch is divided into a grid consisting of 4×4 subregions. For each subregion, a feature vector composed of four dimensions is calculated using a Haar wavelet-filter and then a sum vector of orientations is calculated in each cell. Finally, the concatenation of the feature vectors of each subregion produces the SURF descriptor ($4 \times 4 \times 4 = 64$ dimensions).

2.1.2. Binary Descriptors

BRIEF – Binary Robust Independent Elementary Features

The BRIEF descriptor [36] is one the simplest of the binary descriptors and also the first published. By itself, BRIEF is neither scale nor rotation invariant and does not have an elaborate sampling pattern. Nevertheless, its performance is similar to a more complex local descriptor, the SURF, when compared to its robustness to illumination, blur, and perspective distortion.

The BRIEF descriptor is represented by a binary string in which each bit represents a simple comparison between two elements inside a patch. The bit is set to '1' if the first point is more intense than the other one, otherwise it is set to '0'. The most common strategy for choosing these points is based on a randomly way according to a Gaussian distribution with respect to the keypoint of the patch. The number of selected points leads to the descriptor size (e.g., 128, 256 and 512 bits).

ORB – Oriented FAST and Rotated BRIEF

The ORB descriptor [37] is similar to BRIEF. However, it is robust to noise and invariant to rotation. The invariance to rotation is achieved by estimating the patch rotation using the intensity centroid. Patch moments are used to compute the intensity centroid, which outperform gradient-based approaches.

The sampling pattern is steered estimating the orientation and usual binary tests are used for computing the descriptor. Furthermore, for selecting a couple of points, a k -nearest neighborhood strategy based on error-prone is performed. The random sampling has been replaced by a sampling scheme that uses machine learning for decorrelating BRIEF features under rotational invariance. Unlike BRIEF, the ORB descriptor size is fixed in 256 bits.

BRISK – Binary Robust Invariant Scalable Keypoints

The approach used by BRISK [38] is very similar to BRIEF once the descriptor is computed based on simple binary tests comparing pixel intensities. However, BRISK presents three main differences: (i) it takes into account the rotation of the point to be described; (ii) it makes use of a scale-space theory to maximum adapt the sampling pattern in scale space; and (iii) it uses a special pattern for the binary tests, instead of a random probability distribution. In this way, BRISK becomes invariant to rotation and scale.

As illustrated in Figure 1(a), the BRISK descriptor uses a pattern of points p_i equally distributed in concentric circles around the keypoint to be described. To compare points, the authors defined two distinct sets of pairs of points, long distance pairs and short pairs. Long distance pairs are composed of (i, j) where $\|p_i - p_j\| > \delta_{min}$ and they are used to estimate the orientation of

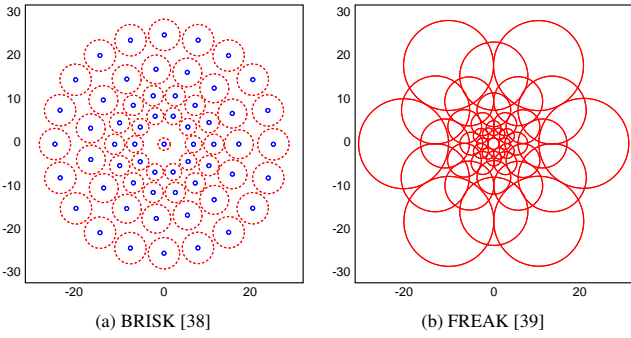


Figure 1: Sampling patterns of two local binary descriptor. (a) Illustrates the sampling pattern of BRISK descriptor which is based on 60 points: the small blue circles denote the sampling locations; the bigger red dashed circles are drawn at a radius, which corresponds to the standard deviation of the Gaussian kernel used to smooth the intensity values at the sampling points. (b) Illustrates the sampling pattern of FREAK descriptor in which each circle represents a receptive field where the image is smoothed with its corresponding Gaussian kernel.

the keypoint using a gradient mean. Then, a Gaussian smoothing is applied to the concentric rings used for the sampling pattern and 512 short pairs (whose distance is less than a threshold δ_{max}) are used to build the descriptor with simple comparisons between the points.

FREAK – Fast Retina Keypoint

To extend the BRISK descriptor, Alahi et al. [39] proposed a descriptor called FREAK. Similarly to BRISK, the FREAK descriptor has its sampling pattern based on Gaussians but its sampling point distribution is biologically inspired by the retina pattern of the human eye.

FREAK has two main important differences compared to BRISK. First, it consists of an allocation of concentric distributions with an exponential growth over the distance of the keypoint. The second difference is based on the fact that the sampling pattern overlaps on different concentric circles, as shown in Figure 1(b). The overlap between sampling regions adds some redundancy that increases the discriminative power of the descriptor. The authors mention that this redundancy is also present in the receptive fields of the human retina.

The sampling pattern used in the original implementation of FREAK descriptor consists of 43 “receptive fields”, leading to 903 possible comparison tests. Thus, for the construction of the final descriptor, FREAK uses a similar method to the ORB descriptor using a greedy approach to select the less correlated comparison tests making it more discriminative. To achieve a maximum performance, 512 binary tests are used.

BinBoost

Trzcinski et al. [40] proposed a new framework with the aim of creating a binary descriptor extremely compact and highly discriminative. BinBoost descriptor is robust to changes in lighting and viewpoint.

Unlike the aforementioned binary descriptors which compute the feature vector based on simple binary tests comparing the intensity of pixels, each bit generated by BinBoost is computed

by using a binary hash function, the same way as the AdaBoost classifier does [41]. This function is based on weak learners that take into account the orientation intensity of gradients on the patch to be described. The hash function is optimized iteratively, i.e., at each iteration, incorrect samples are assigned to a greater weight while the weight of the correct samples is decreased. In this way, the next bits to be calculated will tend to correct the error of their predecessors.

2.2. Mid-level Representation: BossaNova

We now overview the BossaNova mid-level image representation, which introduces a density-based pooling strategy by computing the histogram of distances between the local descriptors and the codewords. More details can be found in [4, 42].

Let \mathcal{X} be an unordered set of binary descriptors extracted from an image. $\mathcal{X} = \{\mathbf{x}_j\}$, $j \in [1, N]$, where $\mathbf{x}_j \in \{0, 1\}^D$ is a binary descriptor vector and N is the number of binary descriptors in the image. Let \mathcal{C} be a visual codebook² obtained by the k -medians algorithm. $\mathcal{C} = \{\mathbf{c}_m\}$, $m \in [1, M]$, where $\mathbf{c}_m \in \{0, 1\}^D$ is a codeword and M is the number of visual codewords. \mathbf{z} is the final vectorial BossaNova representation of the image used for classification.

To keep more information than the BoW during the pooling step, the BossaNova pooling function, g , estimates the probability density function of α_m : $g(\alpha_m) = \text{pdf}(\alpha_m)$, by computing the following histogram of distances $z_{m,b}$:

$$\begin{aligned}
 g : \mathbb{R}^N &\longrightarrow \mathbb{R}^B, \\
 \alpha_m &\longrightarrow g(\alpha_m) = z_m, \\
 z_{m,b} &= \text{card} \left(\mathbf{x}_j \mid \alpha_{m,j} \in \left[\frac{b}{B}; \frac{b+1}{B} \right] \right), \\
 \frac{b}{B} &\geq \alpha_m^{\min} \text{ and } \frac{b+1}{B} \leq \alpha_m^{\max}, \quad (1)
 \end{aligned}$$

where B denotes the number of bins of each histogram z_m , $\alpha_{m,j}$ represents a dissimilarity (i.e., a distance) between \mathbf{c}_m and \mathbf{x}_j , and $[\alpha_m^{\min}; \alpha_m^{\max}]$ limits the range of distances for the descriptors considered in the histogram computation.

In addition to the pooling strategy, Avila et al. [4] also proposed a localized soft-assignment coding that considers only the k -nearest codewords for coding a local descriptor. After computing a local histogram z_m for all the c_m centers, the BossaNova image representation \mathbf{z} [4] is given by:

$$\mathbf{z} = \left[[z_{m,b}], s t_m \right]^T, \quad (m, b) \in \{1, \dots, M\} \times \{1, \dots, B\}, \quad (2)$$

where \mathbf{z} is a vector of size $M \times (B + 1)$, s is a nonnegative constant and t_m is a scalar value for each codeword, counting the number of binary descriptors \mathbf{x}_j close to that codeword.

The idea of enriching BoW representations with extra knowledge from the set of local descriptors has been explored on several approaches [43, 44]. However, those works opt by parametric models that lead to very high-dimensional image representations. By using a simple histogram of distances to capture

²The codebook is usually built by clustering a set of local descriptors. It can be defined by the set of codewords (or visual words) corresponding to the centroids of clusters.

the relevant information, BossaNova remains flexible and keeps the representation compact. For these reasons, we decided to employ it in this work for mid-level features.

In short, the BossaNova vector is defined by three parameters: (i) the number of codewords M ; (ii) the number of bins B in each histogram; and (iii) the range of distances $[\alpha_m^{min}, \alpha_m^{max}]$. As in [4], we set up the bounds as $\alpha_m^{min} = \lambda_{min} \cdot \sigma_m$ and $\alpha_m^{max} = \lambda_{max} \cdot \sigma_m$, where σ_m is the standard deviation of each cluster c_m obtained by k -medians clustering algorithm.

3. Related Works

According to Ries and Lienhart [19], works involving pornography detection can be divided into three main groups: (i) approaches based on skin color that exploit the assumption that pornographic images/videos generally have large areas of skin color; (ii) approaches based on shape information; and (iii) approaches based on local features in conjunction with Bag-of-Words (BoW) models.

In this section, we survey the literature on the pornography detection. In Section 3.1, we cover works based on skin color detection as well as works based on shape information, since all shape based approaches presented in this paper also employ a step of finding pixels that have skin related colors. Section 3.2 presents the works that make use of local features and BoW models.

3.1. Skin Color and Shape Information Based Approaches

Most of the pornography detection works are based on skin color detection or shape information. This is partly because the most obvious property in pornographic images is the large fraction of pixels presenting skin related colors and that most of the pornographic images share some characteristic shapes [19].

The approach proposed by Forsyth and Fleck [8] begins by finding areas with skin color in the image. First, they transformed each pixel value into intensity value and two tone values. After that, decision rules are applied to find regions with skin color. Upon skin detection, a corner detector and a Hough transform are applied to find candidates for human limbs. These candidates are iteratively combined according to a set of constraints to model the geometry of the human body. Finally, if it is possible to assemble the limbs in a geometrically reasonable way, the image is classified as pornographic.

Jones and Rehg [9] constructed a 3D histogram of 256 bins for each color channel. From these histograms, five different features are extracted, such as the percentage of pixels related to the skin or the number of connected areas of skin. Finally, a decision tree is trained based on these characteristics. However, the authors presented results suggesting that histogram with less bins are sufficient, since the latter can cause overfitting.

Inspired by the color histogram used by Jones and Rehg [9], Rowley et al. [11] generated a skin color based map to determine connected components. Then, skin based and connected component features are extracted from the map, such as mean and standard deviations. The authors also employed other color characteristics (e.g., the edge pixels within the skin regions).

Finally, these characteristics are used as input to a SVM classifier.

The approach presented by Lee et al. [45] employed a learning scheme based on the skin color distribution of the image, using a neural network to learn and classify whether the input image contains skin exposure. Furthermore, a feature is used to detect textures with roughness to reject non-skin objects. Then, three types of features related to the targeted form (area size, aspect ratio and location) are extracted and sent to an AdaBoost classifier. Finally, a face detection algorithm is applied to filter out false candidates related to face photos. Using the latter technique, Lee et al. [46] separated the image colors in skin and non-skin groups. Next, a texture analysis is applied to verify if the likelihood of the area is composed by skin and a face detection algorithm is applied to eliminate face photos. In addition, the authors verified the presence of “holes” in the binary images to detect photos related to swimsuits. For the remaining images, features regarding the position of the skin region and morphological characteristics are extracted to train a SVM classifier.

A method to estimate skin regions was proposed by Yu and Han [16] using simple operations in the HSV color space plus an additional post-processing to reduce noise. The method is fast and accurate enough to filter “easy” pornographic images before a more robust identification process. Basically, it is a threshold used in the Hue component to select pixels related to skin. Then, an edge density map is computed to remove incorrectly detected regions. Assuming that the density of edges is low in skin regions, the edge pixels that have a high density are removed. Morphological operations are also used to reduce possible noise. Finally, mean and standard deviation of skin regions are calculated and another threshold is used to decide whether the image is pornographic or not.

Zaidan et al. [47] proposed an anti-pornography system based in two different stages: (i) skin detection and (ii) pornography classification. In the first stage, to detect skin regions from the image they combined the Bayesian method, a grouping histogram technique and back-propagation neural network based on the YCbCr and RGB color spaces. In the second stage, the features from the skin are extracted and classified, which determines if an image is pornographic.

Although there are many approaches based on skin color detection to classify pornographic content, these works often have the disadvantage of a high rate of false positives since not all images with large areas of skin exposure are necessarily pornographic (for example, pictures with people wearing swimsuits, or sports-related images). Furthermore, another hurdle is the diversity of human skin color, making the classification process even more complicated. Another issue is that grayscale images cannot be classified using color related features [19]. Shape-based approaches present the same problem, since they also make use of skin color information. It is important to mention that these are image-based approaches, and to be applied to video, a voting scheme should be considered.

3.2. Local Features and Bag-of-Words Based Approaches

Another widely used approach in the literature of pornography detection is the employment of local feature extraction.

Most of the works applied BoW models as an intermediate representation to encode the extracted local features followed by a classification step.

Deselaers et al. [2] were the first to use local features with BoW models to classify pornographic content. They proposed an approach to filter and classify pornography in different categories. The SIFT detector was used to detect interesting points. They did not use any descriptor for describing the detected regions, claiming the advantage of patches as they provide color information. Next, each patch was reduced using Principal Component Analysis (PCA) and SVM was used as the final classifier. In a similar way, Lopes et al. [5] performed image classification using the SIFT detector, HueSIFT descriptors and SVM classifier. A comparison is also made between the SIFT descriptors and HueSIFT applied to pornography, showing that the combination of color information and local features perform better. In [13], the same authors extended their work to detect nudity in videos. To that end, they performed the same approach for selected video frames, but a majority voting scheme is held over the frames to set the video final classification.

To pre-filter pornographic images for isolating features of interest, Steel [15] proposed a variant of the SIFT descriptor (Mask-SIFT) that uses a Gaussian pre-filter to remove all pixels of an image which are not related to skin. The image is then processed using a median filter to fill missing pixels and eliminate noise by creating a “mask image”. Once this “mask image” is created, the SIFT descriptor is used to extract features from human related areas. For classification, the author developed a cascade based classifier filtering the images based on skin, shape and local features to determine whether an image is pornographic.

A study regarding the impact of moving patterns was made by Souza et al. [48]. Color information is incorporated into the space-time interest point detector (STIP). They incorporated color information using the normalized-RGB color system at the detection phase called ColorSTIP. The local space-time regions are described in terms of histograms of Hue and combined with the default HOG-HOF feature histogram used by STIP, which they named as HueSTIP. BoW models are applied to encode the spatial-temporal features and SVM as final classifier.

Avila et al. [4] also presented an approach for pornographic video classification. First, the videos are segmented into shots and then the central frame of each shot is taken as keyframes to represent the video. After that, HueSIFT descriptors are extracted on a dense spatial grid. However, different from previous approaches, they used a BoW model extension — the BossaNova mid-level image representation — to encode the local features. A SVM classifier is then applied to classify the keyframes extracted from each video shot and a majority voting scheme is employed to predict the video class.

With the aim of a fast feature detection and extraction, Zhuo et al. [18] proposed a pornographic image recognition method based on the ORB binary descriptor. Their approach is divided in two parts: coarse detection and fine detection. The coarse detection identifies the non-pornographic images using a skin color detector conducted in YCbCr color space. For the remain-

ing images containing much more skin-color regions, a fine detection is conducted. To that end, ORB descriptors are extracted from the skin-color regions and encoded by a BoW model. After that, they combined it with a 72-dimensional HSV color histogram of the whole image and trained a classification model using SVM for image recognition. Another fast approach based on binary descriptors was proposed by Yaghoubyan et al. [49]. First, they detected a region of interest (ROI) per input image by a threshold according to the ratio of skin color area in the SKN color space. Images with higher probability proceed to the further stages, otherwise they are labeled as non-pornographic. Then, keypoints are detected in the ROI eliminating noisy keypoints placed in the image background. FREAK binary descriptors are encoded by a BoW representation and a SVM classifier is applied to recognition. We also explored in [17] the benefit of using several binary descriptors and mid-level representation for videos to identify pornographic content with a SVM classifier and a majority voting scheme.

In contrast to previous approaches, which employed descriptors as local features, Ulges and Stahl [14] used the low frequency coefficients of the Discrete Cosine Transform (DCT) in YUV color space. As in previous approaches, BoW technique is also used as well as SVM for classification. To compare, the authors also implemented an approach based on color similar to Jones and Rehg [9] and concluded that the BoW-based approach outperforms approaches based only on color. Zhang et al. [50] employed visual attention model to find regions of interest (ROI) composed by pornographic content. Given an image, a face detection algorithm is applied to remove the face or ID photo from some non-pornographic images. Then, a visual salient map in compressed domain is computed to construct visual attention model, according to the large number of exposed skin regions. After, four features of color, texture, intensity and skin are extracted from the pornographic regions and a BoW model is applied. Finally, SVM is employed for classification.

In general, approaches based on local features have shown more satisfactory results than the ones based on color information [19]. An important advantage of using local features is that they can be computed independent of color information. Moreover, another advantage is that these models compact the regions of the image in a fixed-size vector, making easy the comparison of image regions and/or the image as a whole. However, extraction of local features may prove to be more time-consuming than to examine image characteristics related to color.

4. BossaNova Video Descriptor (BNVD)

Multimedia understanding is related to visual recognition in which we are interested in identifying, for example, pornography content and action recognition, in multimedia collections. Usually, the approaches to cope with those problems consider (i) extraction of local image descriptors; (ii) encoding of the local features in a mid-level representation; and (iii) classification and/or search of the image descriptor. When a static image is considered, the visual recognition is based on, for instance, classification of a mid-level representation computed from the

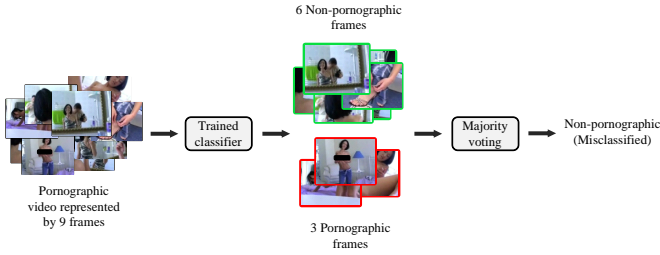


Figure 2: Example of video containing nine frames, six non-pornographic and three pornographic. The video will be misclassified to non-pornographic by the majority voting scheme, even been pornographic.

local image descriptors. A simple extension of this approach for video classification is to compute a mid-level representation by extracting local image descriptor from the video frames. Unfortunately, this simple and naive approach presents poor performance since it is susceptible to noise and it is not so discriminative for representing the video content.

To bypass this aforementioned problem, some methods are based on majority voting in which a binary classification is performed over the images. In a majority voting scheme to decide about the classes A and B, for example, the video is classified as belonging to the class A if the number of images classified as class A is greater than the number of images classified as class B. In literature, this approach is used in several applications, including pornography detection [4, 13, 17]. Despite the good performance achieved on those works, some issues must be considered: (i) if a video contains few frames of a class A, the video will be classified as class B (as illustrated in Figure 3), however in some applications such as video pornography, the existence of a few pornography frames must characterize the video as pornography; (ii) the presence of noise could exert influence over the number of frames of a specified class; and (iii) the probability of correct and wrong classifications are ignored since the evaluation is done by a static frame.

To cope with these issues, the classification of a video must be done from a descriptor which represents all video content. In this way, we propose a strategy that aggregates the information of the mid-level representations of all video frames into a single representation. We present our video descriptor as follows.

Let \mathcal{V} be a video sequence. $\mathcal{V} = \{f^i\}$, $i \in [1, N]$, where f^i is the keyframe³ of the shot i and N is the number of keyframes. Let $\mathcal{Z} = \{z^i\}$, $i \in [1, N]$ be a set of BossaNova vectors computed for the video \mathcal{V} in which z^i is a BossaNova vector extracted for the keyframe f^i . Let \mathbf{O} and \mathbf{P} be two functions for aggregating the information of BossaNova and the Bag of Visual Words. The BossaNova Video Descriptor (BNVD) can be modeled by a function \mathbf{W} as follows:

$$\begin{aligned} \mathbf{O} : \mathbb{R}^B &\longrightarrow \mathbb{R}^B, \\ \mathbf{P} : \mathbb{R}^M &\longrightarrow \mathbb{R}^M, \end{aligned}$$

³A keyframe is a frame that represents the content of a logical unit, like a shot or scene, for example.

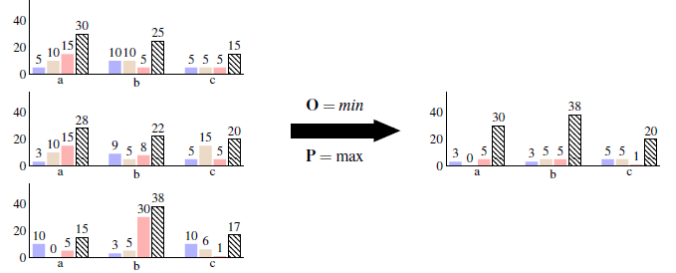


Figure 3: An example of computing the video descriptor from three image descriptors. The functions *min* and *max* are used for aggregating the set of BossaNovas.

$$\begin{aligned} \mathbf{W} : \mathbb{R}^Z &\longrightarrow \mathbb{R}^Z, \\ \mathcal{Z} &\longrightarrow \mathbf{W}(\{z^i\}) = [[o_{m,b}], p_m]^T, \\ o_{m,b} &= \mathbf{O}(\{z_{m,b}^i\}), \\ p_m &= \mathbf{P}(\{t_m^i\}), \end{aligned} \quad (3)$$

where $Z \subset \{1, \dots, M\} \times \{1, \dots, B\}$, and $z^i = [[z_{m,b}^i], t_m^i]^T$.

Intuitively, this new video descriptor represents a relation for each codeword to the codebook, since each BossaNova representation contains information regarding the distance-to-codeword distribution. The main goal of applying the functions \mathbf{O} and \mathbf{P} to the BossaNova vectors is to employ a filtered-like operation to the entire video content which is represented by this mid-level representation. In this work, we study the behavior of our video descriptor by using the following functions: *median*, *mean*, *min* and *max*.

We also propose another video descriptor by using the classical BoW representation, called Bag-of-(Visual)-Words Video Descriptor (BoW-VD), that can be seen as a simplification of the BNVD since the distance to the codewords (z_m) are ignored.

Figure 3 presents an example of applying *min* and *max* functions to the combination of three image descriptors. In this case, the *min* function is applied to the local histograms z_m (represented by the colored bins) and the *max* function is applied to t_m (represented by the dashed bin).

5. Experimental Results

In this section, we present the results obtained through the evaluation of the proposed video descriptors applied to the pornography detection task. After describing our experimental setup (Section 5.1) and giving details regarding Pornography dataset [4] (Section 5.2), we report the results achieved using the proposed video descriptors. Those results are organized in two groups: (i) quantitative (Section 5.3) and (ii) comparative analysis (Section 5.4).

5.1. Experimental Setup

In Figure 4, we illustrate the overview of our experimental approach. Note that this methodology is adapted from [4, 17], in which the binary image descriptors are computed

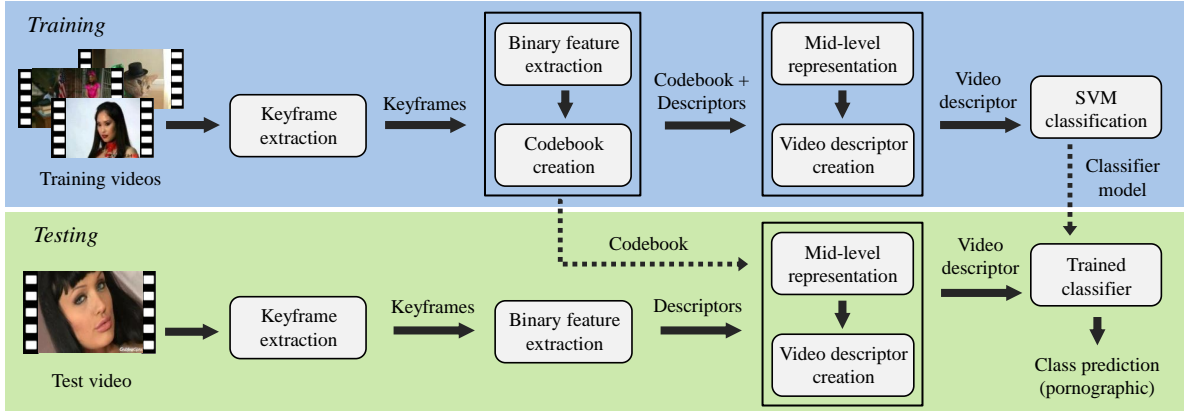


Figure 4: Overview of our methodology for pornography detection with the proposed video descriptor.

followed by the computation of a mid-level representation for each keyframe. Our proposed video descriptors are computed according to Equation 3 for representing each video that will be classified.

For comparison purposes, we applied the same experimental setup described in [17] to the binary feature extraction. We extracted the most representative binary descriptors — BRIEF [36], ORB [37], BRISK [38], FREAK [39] and BinBoost [40] — from patches of 16×16 pixels sampled regularly using a step size of 6 pixels. For BRIEF, ORB, BRISK, and FREAK, we obtained the implementations from OpenCV repository [51]; for BinBoost, we used the code available on the authors website⁴ and we assessed the length of the descriptor which may vary among 8, 16 and 32 bytes.

Codebook learning is performed by a k -medians clustering algorithm with Hamming distance over one million randomly sampled descriptors. The motivation to use such clustering technique is the fact that it is able to deal with the binary string nature of the binary descriptors.

For BossaNova mid-level feature extraction, we kept the parameter values the same as in [4, 17]: $B = 10$, $\lambda_{min} = 0$, $\lambda_{max} = 3$ and $s = 10^{-3}$, except for the number of visual code-words M . For comparison, we also extracted BoW mid-level features, obtained with hard-assignment coding and average pooling.

For binary classification (pornography vs. non-pornography), we used support vector machines (SVM) with a nonlinear kernel. Kernel matrices are computed as $\exp(-\gamma d(x, x'))$ in which d is the distance and γ is set to the inverse of the pairwise mean distances. The SVM penalty parameter C is set to 10, which works best for the Pornography dataset.

All experiments were conducted on a 64-bit Ubuntu Linux machine powered by Intel® Xeon® CPU X5670 @ 2.93 GHz with 24 cores and 70 GB RAM.

5.2. Pornography Dataset

We evaluate our approach on the challenging Pornography dataset [4]. It is composed of 400 pornographic and 400 non-

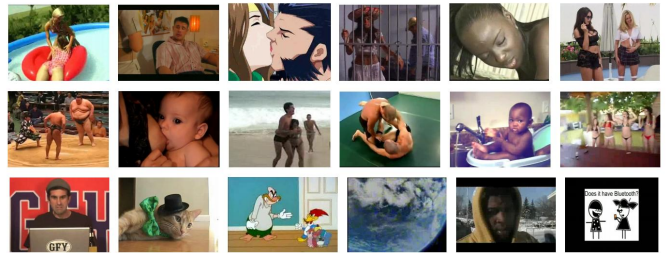


Figure 5: Example frames from the Pornography dataset [4], illustrating the diversity of pornographic videos (top row) and the non-pornographic videos (difficult on the middle row, a group with high skin exposure that is quite challenging for the detector, and easy at the bottom row).

pornographic videos resulting in nearly 80 hours. The non-pornographic class is divided into two sub-classes: (i) “easy” with 200 randomly selected videos from the Internet; and (ii) “difficult”, with 200 videos selected from text search queries as “beach”, “wrestling” and “swimming”, which is a group of high skin exposure becoming a quite challenging dataset. The video shots are already separated and a keyframe is selected to summarize the content of each shot into a static image (16,727 video keyframes). The experimental protocol of the dataset consists of a classical 5-fold cross-validation (640 videos for training and 160 for testing, on each fold). The video classification performance is reported by accuracy rate, where the final video label is obtained by the classification of the video descriptor, as described in Section 4.

Figure 5 depicts the diversity of pornographic videos and the non-pornographic videos (difficult and easy). The dataset is available upon request and the sign of a license agreement⁵.

5.3. Quantitative Analysis

In this section, we present discussions involving different aggregation functions used for combining the mid-level representation. We also evaluate the impact of codebook size on the proposed video descriptors. Furthermore, we present the computational cost of different strategies to represent the data.

⁴<http://www.cvlab.epfl.ch/research/detect/binboost>

⁵<https://sites.google.com/site/pornographydatabase/>

Table 1: Video classification (%) results (and standard deviations) of the proposed video descriptors, using Bag-of-Words (BoW), BossaNova, and different aggregation functions on Pornography dataset. For each method, we extracted a codebook size of $M = 256$, as suggested in [4, 17].

Approach	<i>max</i>	<i>min</i>	<i>mean</i>	<i>median</i>
BoW-VD (BRIEF)	85.55 ± 4	68.94 ± 5	86.79 ± 2	87.16 ± 2
BoW-VD (ORB)	87.54 ± 3	77.18 ± 3	87.29 ± 2	88.04 ± 2
BoW-VD (BRISK)	86.55 ± 2	71.70 ± 5	88.04 ± 3	87.04 ± 3
BoW-VD (FREAK)	85.91 ± 4	80.30 ± 3	86.55 ± 3	86.42 ± 3
BoW-VD (BinBoost $d = 8$)	86.66 ± 2	72.44 ± 3	86.78 ± 1	87.78 ± 1
BoW-VD (BinBoost $d = 16$)	87.54 ± 2	73.32 ± 3	87.78 ± 1	87.77 ± 3
BoW-VD (BinBoost $d = 32$)	87.28 ± 2	72.82 ± 4	88.28 ± 2	88.53 ± 2
BNVD (BRIEF)	88.16 ± 3	80.92 ± 4	88.66 ± 2	89.03 ± 1
BNVD (ORB)	89.28 ± 3	84.42 ± 3	89.28 ± 3	89.02 ± 1
BNVD (BRISK)	87.16 ± 2	82.30 ± 4	88.66 ± 0.5	89.27 ± 1
BNVD (FREAK)	87.66 ± 2	85.66 ± 1	88.04 ± 1	89.66 ± 2
BNVD (BinBoost $d = 8$)	89.28 ± 3	82.05 ± 3	87.66 ± 1	90.77 ± 2
BNVD (BinBoost $d = 16$)	87.79 ± 2	82.42 ± 3	88.54 ± 2	90.90 ± 1
BNVD (BinBoost $d = 32$)	88.28 ± 2	81.92 ± 2	88.15 ± 1	89.41 ± 2

5.3.1. Aggregation Function

To compare different aggregation functions for the proposed video descriptors (as mentioned in Section 4), we applied four different functions to $z_{m,b}^i$ and t_m^i : (i) *max*, which selects the maximum existing value; (ii) *min*, as the opposite of *max*, selects the minimum existing value; (iii) *mean*, which computes the average values; and (iv) *median*, which selects the median values. Table 1 presents these experiments for both BossaNova and BoW. It can be seen that the *median* function outperforms the others, except for the approach using the ORB descriptor for BNVD and BRISK, FREAK and BinBoost ($d = 16$) for BoW. In such cases, *mean* function presented a slightly higher accuracy value, but trebling the standard deviation on BNVD when compared to the *median* function. Since the *median* function presented the best results when compared to the other functions, from now on we use the results obtained by it as the main comparative method.

5.3.2. Codebook Size

To provide a more comprehensive analysis of the proposed video descriptors, we evaluated their behavior by varying the codebook size M in the classification performance. These experiments are shown in Figure 6. We can observe that the approach using BinBoost ($d = 16$) for both BoW-VD and BNVD achieved 92.40% of accuracy (with a codebook size $M = 4096$) and 92.02% of accuracy (with a codebook size $M = 1024$), respectively. The results obtained correspond to the mid-level parameters used from [4], so these parameters were not tuned for the proposed descriptors. Therefore, our approaches can achieve further results, as can be seen in Figure 6.

5.3.3. Computational Cost

Table 2 presents a computational average time required to (i) extract the descriptors, (ii) create the codebook, (iii) generate the mid-level representations, and (iv) compute the proposed video descriptors. We can see how fast binary descriptors are

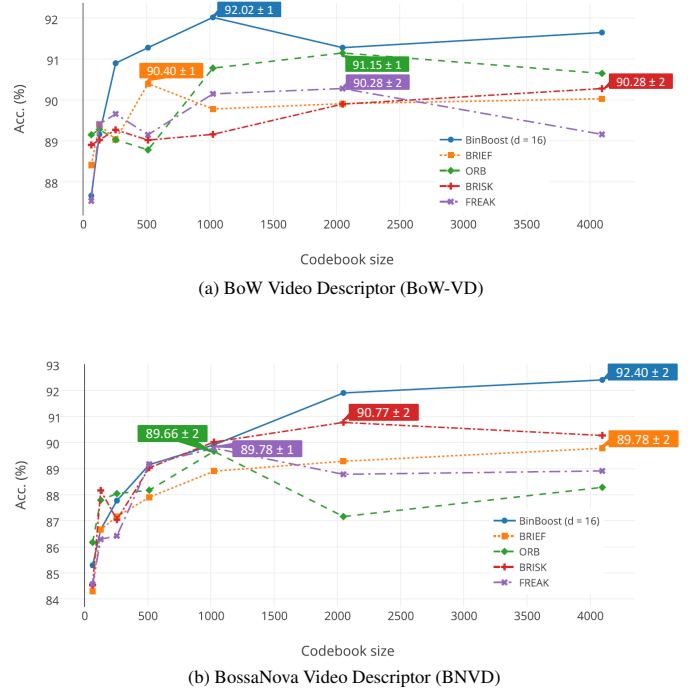


Figure 6: Video classification (%) results (and standard deviations) of video descriptors varying the codebook size M on Pornography dataset.

compared to the HueSIFT descriptor used by [4]. BRIEF and ORB descriptors presented an average extraction 10 times faster than HueSIFT and BinBoost ($d = 32$), which is the slowest binary descriptor, becoming twice as fast.

It is also possible to observe in Table 2 how the computational times to create the codebook differs from our approaches to [4]. In this evaluation, Avila et al. [4] approach showed the best time. This is partly because our approaches employed the k -medians algorithm while Avila et al. [4] used the k -means algorithm. However, it is important to emphasize that the codebook creation step is performed during the training phase, i.e.,

Table 2: Computational time required to: (i) extract the descriptors, (ii) create the codebook, (iii) generate the mid-level representations, and (iv) compute the proposed video descriptor.

Approach	Descriptor	Codebook	BoW	BossaNova	BoW-VD	BNVD
HueSIFT [4]	2.54	1.61×10^3	0.12	1.33	–	–
BRIEF	0.24	4.75×10^3	0.05	0.50	5.88	15.26
ORB	0.23	6.51×10^3	0.04	0.49	5.66	15.84
BRISK	0.64	10.23×10^3	0.08	0.83	14.74	30.20
FREAK	0.31	9.48×10^3	0.08	0.84	7.96	23.70
BinBoost $d = 16$	0.81	16.12×10^3	0.02	0.23	17.12	21.44

an offline phase. Moreover, Table 2 presents the computational time for creating the proposed video descriptor. The calculation is done by adding the average time values of extracting the descriptors with the average time for creating the mid-level representation. This value is multiplied by the average number of keyframes per video. Finally, we sum up the combination time between the mid level representations. Considering a codebook size of $M = 256$, on average, this time is 0.001 and 0.012 milliseconds for BoW-VD and BNVD, respectively.

5.4. A Comparative Analysis

In this section, we perform a comparison with classical and state-of-the-art methods, including both experiments with methods we have reimplemented ourselves and published results reported in the literature.

5.4.1. Comparison to Classical Methods

In Table 3, we present a comparison of our proposed approaches to three different methods: (i) a implemented classical BoW method; (ii) our previous work [17]; and (iii) an approach using video global pooling. We can notice a considerable improvement achieved by our proposed video descriptor in both approaches, BoW-VD and BNVD. BoW-VD reached 88.0% of accuracy with ORB descriptor outperforming the basic BoW (BRISK) approach. Moreover, BNVD reached 90.9% of accuracy with BinBoost descriptor ($d = 16$). The comparison with our previous work [17] is particularly relevant because we employed the same binary descriptors (BRIEF, ORB, BRISK and FREAK with default parameters). We note an absolute improvement of (by up to) 2.8% from BossaNova to BNVD, reducing the classification error by more 21%.

Moreover, in Table 3, we present a comparison of our proposed approaches with an approach using video global pooling, i.e., creating just one mid-level representation for the video using all local features extracted from the keyframes. It can be seen that both BNVD and BoW-VD offered superior accuracy values when compared to a simple approach such as the global pooling of all the local features.

5.4.2. Comparison to State-of-the-art Methods

In Table 4, we present our results in comparison to the previous published methods that were also assessed in the Pornography dataset [4]. We can notice the improvement obtained by our proposed video descriptor in both approaches, BoW-VD and

BNVD. BoW-VD reached 92.4% of accuracy with BinBoost descriptor ($d = 16$) outperforming the basic BoW approaches from Souza et al. [48] and Avila et al. [4]. Moreover, BNVD reached 92.0% of accuracy with BinBoost descriptor ($d = 16$) also outperforming the BossaNova approaches from Avila et al. [4] and Caetano et al. [17]. Further, our results outperformed all the published local feature approaches which, as far as we know, used HueSIFT [4] and spatio-temporal features [48].

Tables 5–8 show the confusion matrices for Souza et al. [48], Avila et al. [4], BoW-VD (BinBoost $d = 16$) and BNVD (BinBoost $d = 16$), respectively. We can notice a higher value of true positives achieved by BNVD and BoW-VD (92.3% and 93.0%, respectively) against Souza et al. BoW (ColorSTIP) [48] and Avila et al. BossaNova (HueSIFT) [4] (90.0% and 88.2%, respectively). Regarding the true negative values, our approaches achieved better when compared to Avila et al. [4] (91.8% against 90.8%) and a slightly close value when compared to Souza et al. [48] (91.8% against 92.0%).

Table 4 also presents results for Moustafa et al. [52] convolutional neural network (CNN) approach. We note that, while we do not have greater accuracy than method proposed by Moustafa et al. [52], we remain close. Despite this, it is important to mention that CNN approaches can be computationally expensive requiring much cost in time and computational resources. On the other hand, our approach is much more simple employing a low-complexity alternative for feature extraction using binary descriptors and linear aggregation functions to BoW models.

In Figure 7, we illustrate the ROC curves and the area under curve (AUC) for the Pornography dataset. Since 5-fold cross-validation protocol was applied, Figures 7 (a) and (b) presents the ROC curves for each fold and the mean curve for BoW-VD (BinBoost $d = 16$) and BNVD (BinBoost $d = 16$), respectively. Figure 7 (c) presents curves for our approaches in comparison to Avila et al. [4] and Caetano et al. [17]. Although the methods presented very similar curves, we can see that our proposed approaches achieved better AUC values. *To the best of our knowledge, ours is the best result reported to date on Pornography dataset employing local feature descriptors.*

We also investigated the cases where our method failed. The misclassified non-pornographic videos correspond to very challenging cases, such as breastfeeding sequences, sequences of children being bathed, and beach scenes (as illustrated in Figure 8(a)). In addition, the analysis of the misclassified

Table 3: Video classification (%) results (and standard deviations) of our approaches, implemented classical BoW methods, our previous work [17], and global pooling results on Pornography dataset [4]. For each method, we extracted a codebook size of $M = 256$, as suggested in [4, 17].

	Approach	Acc. (%)	Approach	Acc. (%)
Implemented methods	BoW (HueSIFT) [4]	83.0 ± 3	BossaNova (HueSIFT) [4]	89.5 ± 1
	BoW (BRIEF)	85.0 ± 3	BossaNova (BRIEF) [17]	86.3 ± 3
	BoW (ORB)	85.8 ± 2	BossaNova (ORB) [17]	86.5 ± 3
	BoW (BRISK)	87.0 ± 1	BossaNova (BRISK) [17]	88.6 ± 2
	BoW (FREAK)	85.8 ± 3	BossaNova (FREAK) [17]	86.9 ± 3
	BoW (BinBoost $d = 16$)	86.7 ± 3	BossaNova (BinBoost $d = 16$)	89.4 ± 5
Our results	BoW-VD (BRIEF)	87.2 ± 2	BNVD (BRIEF)	89.0 ± 1
	BoW-VD (ORB)	88.0 ± 2	BNVD (ORB)	89.0 ± 1
	BoW-VD (BRISK)	87.0 ± 3	BNVD (BRISK)	89.3 ± 1
	BoW-VD (FREAK)	86.4 ± 3	BNVD (FREAK)	89.7 ± 2
	BoW-VD (BinBoost $d = 16$)	87.8 ± 3	BNVD (BinBoost $d = 16$)	90.9 ± 1
Global pooling	BoW (BRIEF)	72.7 ± 1	BossaNova (BRIEF)	80.3 ± 3
	BoW (ORB)	72.2 ± 2	BossaNova (ORB)	78.8 ± 3
	BoW (BRISK)	73.1 ± 2	BossaNova (BRISK)	79.9 ± 2
	BoW (FREAK)	73.3 ± 3	BossaNova (FREAK)	79.6 ± 2
	BoW (BinBoost $d = 16$)	71.7 ± 2	BossaNova (BinBoost $d = 16$)	77.3 ± 1

Table 4: Video classification (%) results (and standard deviations) of our approaches and published results on Pornography dataset [4].

	Approach	Acc. (%)
Published results	Souza et al. [BoW (STIP)] [48]	89.6 ± –
	Souza et al. [BoW (HueSTIP)] [48]	90.0 ± –
	Souza et al. [BoW (ColorSTIP)] [48]	91.0 ± –
	Avila et al. [BoW (HueSIFT)] [4]	83.0 ± 3
	Avila et al. [BOSSA (HueSIFT)] [42]	87.1 ± 2
	Avila et al. [BossaNova (HueSIFT)] [4]	89.5 ± 1
	Caetano et al. [BossaNova (BRISK)] [17]	88.6 ± 2
	Moustafa et al. [CNN (AGbNet)] [52]	94.1 ± 2
Our results	BoW-VD (BinBoost $d = 16$)	92.4 ± 2
	BNVD (BinBoost $d = 16$)	92.0 ± 1

pornographic videos revealed that the method presented difficulties with poor quality videos or when the clip is borderline pornographic, with few explicit elements (as illustrated in Figure 8(b)). The same difficulty was also reported by Avila et al. [4].

6. Conclusions

The task of detecting and filtering pornographic visual content from the Internet is a well known concern in many environments, from schools to workplaces. In view of that, we presented a method for video pornography detection. The proposed method integrated the advantages of concepts presented in reference works in pornography detection area, contributing to improve the state of the art.

Specifically, our work focused on the supervised classification of pornographic content in videos based on local feature descriptors coded on mid-level representations. The description

Table 5: Average confusion matrix for Souza et al. [BoW (ColorSTIP)] [48] approach.

		Video was labeled as	
		porn	non-porn
Video class	porn	90.0%	10.0%
	non-porn	8.0%	92.0%

Table 6: Average confusion matrix for Avila et al. [BossaNova (HueSIFT)] [4] approach.

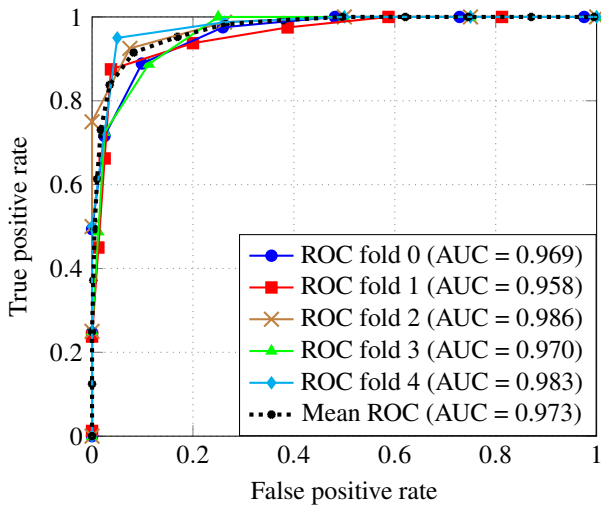
		Video was labeled as	
		porn	non-porn
Video class	porn	88.2%	11.8%
	non-porn	9.2%	90.8%

Table 7: Average confusion matrix for our BoW-VD (BinBoost $d = 16$) approach.

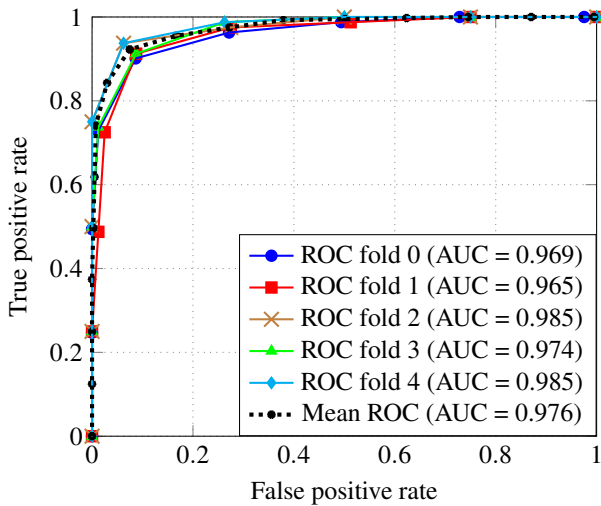
		Video was labeled as	
		porn	non-porn
Video class	porn	93.0%	7.0%
	non-porn	8.2%	91.8%

Table 8: Average confusion matrix for our BNVD (BinBoost $d = 16$) approach.

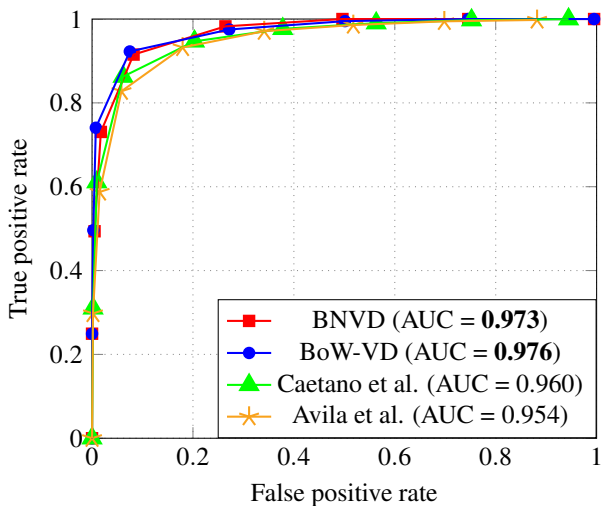
		Video was labeled as	
		porn	non-porn
Video class	porn	92.3%	7.7%
	non-porn	8.2%	91.8%



(a) ROC curves for BNVD (BinBoost $d = 16$)



(b) ROC curves for BoW-VD (BinBoos $d = 16$)

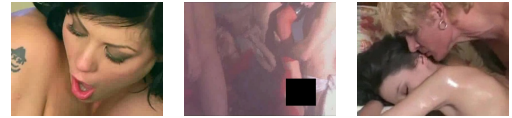


(c) Mean ROC curves

Figure 7: ROC curves and the AUC for the Pornography dataset.



(a) Misclassified non-pornographic videos examples.



(b) Misclassified pornographic videos examples.

Figure 8: Cases where our method fails. (a) Misclassified non-pornographic videos. (b) Misclassified pornographic videos.

of the local features was performed using binary descriptors, a low complexity alternative, and a mid-level representation called BossaNova, an extension of the Bag-of-Words model that richly preserves the visual information.

Our approach is based on a combination of mid-level representations, yielding the Bag-of-Words Video Descriptor (BoW-VD) and BossaNova Video Descriptor (BNVD). The main goal is to apply aggregation functions to combine the collections of mid-level representations, thus creating a filtered-like operation on all the video content which are represented by this mid-level representation. The results validated the proposed video descriptors in the context of pornography detection, outperforming the state of the art by more than two percentage points, reducing the classification error by over 16% (from 9.1% to 7.6%).

In order to provide more comprehensive analysis of our video descriptors, we propose evaluating their behavior on other video classification problems.

7. Acknowledgments

The authors are thankful to CNPq, CAPES and FAPEMIG, Brazilian research and development agencies for the support to this work.

References

- [1] M. B. Short, L. Black, A. H. Smith, C. T. Wetterneck, D. E. Wells, A review of internet pornography use research: Methodology and content from the past 10 years, *Cyberpsychology, Behavior, and Social Networking* 15 (1) (2012) 13–23.
- [2] T. Deselaers, L. Pimenidis, H. Ney, Bag-of-visual-words models for adult image classification and filtering, in: *International Conference on Pattern Recognition (ICPR)*, 2008, pp. 1–4.
- [3] E. Valle, S. Avila, F. Souza, M. Coelho, A. de A. Araújo, Content-based filtering for video sharing social networks, in: *Brazilian Symposium on Information and Computer System Security (SBSeg)*, 2012, pp. 625–638.
- [4] S. Avila, N. Thome, M. Cord, E. Valle, A. de A. Araújo, Pooling in image representation: the visual codeword point of view, *Computer Vision and Image Understanding (CVIU)* 117 (5) (2013) 453–465.
- [5] A. Lopes, S. Avila, A. Peixoto, R. Oliveira, A. de A. Araújo, A bag-of-features approach based on Hue-SIFT descriptor for nude detection, in: *European Signal Processing Conference (EUSIPCO)*, 2009, pp. 1552–1556.

- [6] D. Forsyth, M. Fleck, Identifying nude pictures, in: *IEEE Workshop on Applications of Computer Vision (WACV)*, 1996, pp. 103–108.
- [7] D. A. Forsyth, M. M. Fleck, Body plans, in: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 1997, pp. 678–683.
- [8] D. A. Forsyth, M. M. Fleck, Automatic detection of human nudes, *International Journal of Computer Vision (IJCV)* 32 (1) (1999) 63–77.
- [9] M. Jones, J. Rehg, Statistical color models with application to skin detection, *International Journal of Computer Vision (IJCV)* 46 (1) (2002) 81–96.
- [10] H. Zheng, M. Daoudi, B. Jedynak, Blocking adult images based on statistical skin detection, *Electronic Letters on Computer Vision and Image Analysis (ELCVIA)* 4 (2) (2004) 1–14.
- [11] H. Rowley, Y. Jing, S. Baluja, Large scale image-based adult-content filtering, in: *International Conference on Computer Vision Theory and Applications (VISAPP)*, 2006, pp. 290–296.
- [12] W. Hu, O. Wu, Z. Chen, Z. Fu, S. Maybank, Recognition of pornographic web pages by classifying texts and images, *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* 29 (6) (2007) 1019–1034.
- [13] A. Lopes, S. Avila, A. Peixoto, R. Oliveira, M. Coelho, A. de A. Araújo, Nude detection in video using bag-of-visual-features, in: *Brazilian Symposium on Computer Graphics and Image Processing (SIBGRAPI)*, 2009, pp. 224–231. doi:10.1109/SIBGRAPI.2009.32.
- [14] A. Ulges, A. Stahl, Automatic detection of child pornography using color visual words, in: *IEEE International Conference on Multimedia and Expo (ICME)*, 2011, pp. 1–6.
- [15] C. Steel, The mask-sift cascading classifier for pornography detection, in: *World Congress on Internet Security (WorldCIS)*, 2012, pp. 139–142.
- [16] J.-J. Yu, S.-W. Han, Skin detection for adult image identification, in: *International Conference on Advanced Communication Technology (ICACT)*, 2014, pp. 645–648.
- [17] C. Caetano, S. Avila, S. Guimarães, A. de A. Araújo, Representing local binary descriptors with BossaNova for visual recognition, in: *Symposium On Applied Computing (ACM SAC)*, 2014, pp. 49–54. doi:10.1145/2554850.2555058.
- [18] L. Zhuo, Z. Geng, J. Zhang, X. G. Li, ORB feature based web pornographic image recognition, *Neurocomputing* 173 (3) (2016) 511–517.
- [19] C. Ries, R. Lienhart, A survey on visual adult image recognition, *Multimedia Tools and Applications (MTA)* 69 (3) (2014) 661–688.
- [20] S. Agarwal, A. Awan, D. Roth, Learning to detect objects in images via a sparse, part-based representation, *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* 26 (11) (2004) 1475–1490.
- [21] J. Yang, Y.-G. Jiang, A. G. Hauptmann, C.-W. Ngo, Evaluating bag-of-visual-words representations in scene classification, in: *International Workshop on Workshop on Multimedia Information Retrieval (MIR)*, 2007, pp. 197–206.
- [22] K. E. A. Van de Sande, T. Gevers, C. G. M. Snoek, Evaluating color descriptors for object and scene recognition, *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* 32 (9) (2010) 1582–1596.
- [23] D. Lowe, Distinctive image features from scale-invariant keypoints, *International Journal of Computer Vision (IJCV)* 60 (2) (2004) 91–110.
- [24] C. Caetano, S. Avila, S. Guimarães, A. de A. Araújo, Pornography detection using BossaNova video descriptor, in: *European Signal Processing Conference (EUSIPCO)*, 2014, pp. 1681–1685.
- [25] K. Chatfield, V. Lempitky, A. Vedaldi, A. Zisserman, The devil is in the details: an evaluation of recent feature encoding methods, in: *British Machine Vision Conference (BMVC)*, 2011, pp. 1–12.
- [26] H. Bay, A. Ess, T. Tuytelaars, L.-V. Gool, Speeded-up robust features (SURF), *Computer Vision and Image Understanding (CVIU)* 110 (3) (2008) 346–359.
- [27] J. Sivic, A. Zisserman, Video Google: a text retrieval approach to object matching in videos, in: *International Conference on Computer Vision (ICCV)*, 2003, pp. 1470–.
- [28] M. Datar, N. Immorlica, P. Indyk, V. S. Mirrokni, Locality-sensitive hashing scheme based on p-stable distributions, in: *20th Annual Symposium on Computational Geometry (SCG)*, 2004, pp. 253–262.
- [29] Y. Gong, S. Lazebnik, Iterative quantization: A procrustean approach to learning binary codes, in: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2011, pp. 817–824.
- [30] F. Shen, C. Shen, Q. Shi, A. van den Hengel, Z. Tang, H. T. Shen, Hashing on nonlinear manifolds, *IEEE Transactions on Image Processing (TIP)* 24 (6) (2015) 1839–1851.
- [31] F. Shen, W. Liu, S. Zhang, Y. Yang, H. T. Shen, Learning binary codes for maximum inner product search, in: *IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 4148–4156.
- [32] C. E. dos Santos, E. Kijak, G. Gravier, W. R. Schwartz, Learning to hash faces using large feature vectors, in: *International Workshop on Content-Based Multimedia Indexing (CBMI)*, 2015, pp. 1–6.
- [33] K. Mikolajczyk, C. Schmid, A performance evaluation of local descriptors, *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* 27 (10) (2005) 1615–1630.
- [34] A. Canclini, M. Cesana, R. A., M. Tagliasacchi, J. Ascenso, C. R., Evaluation of low-complexity visual feature detectors and descriptors, in: *International Conference on Digital Signal Processing (DSP)*, 2013, pp. 1–7.
- [35] J. Morel, G. Yu, Is sift scale invariant?, *Inverse Problems and Imaging* 5 (1) (2011) 115–136.
- [36] M. Calonder, V. Lepetit, C. Strecha, P. Fua, BRIEF: binary robust independent elementary features, in: *European Conference on Computer Vision: Part IV (ECCV)*, 2010, pp. 778–792.
- [37] E. Rublee, V. Rabaud, K. Konolige, G. Bradski, ORB: An efficient alternative to SIFT or SURF, in: *International Conference on Computer Vision (ICCV)*, 2011, pp. 2564–2571.
- [38] S. Leutenegger, M. Chli, R. Siegwart, BRISK: Binary robust invariant scalable keypoints, in: *International Conference on Computer Vision (ICCV)*, 2011, pp. 2548–2555.
- [39] A. Alahi, R. Ortiz, P. Vanderghenst, FREAK: Fast retina keypoint, in: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 510–517.
- [40] V. L. T. Trzcinski, M. Christoudias, P. Fua, Boosting Binary Keypoint Descriptors, in: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013, pp. 2874–2881.
- [41] Y. Freund, R. E. Schapire, A decision-theoretic generalization of on-line learning and an application to boosting, *Journal of Computer and System Sciences* 55 (1) (1997) 119–139.
- [42] S. Avila, N. Thome, M. Cord, E. Valle, A. de A. Araújo, BOSSA: Extended BoW formalism for image classification, in: *International Conference on Image Processing (ICIP)*, 2011, pp. 2909–2912.
- [43] J. Sánchez, F. Perronnin, T. Mensink, J. Verbeek, Image classification with the fisher vector: Theory and practice, *International Journal of Computer Vision (IJCV)* 105 (3) (2013) 222–245.
- [44] X. Zhou, K. Yu, T. Zhang, T. Huang, Image classification using super-vector coding of local image descriptors, in: *European Conference on Computer Vision (ECCV)*, 2010, pp. 141–154.
- [45] J.-S. Lee, Y.-M. Kuo, P.-C. Chung, E.-L. Chen, Naked image detection based on adaptive and extensible skin color model, *Pattern Recognition* 40 (8) (2007) 2261–2270.
- [46] J.-S. Lee, F.-S. Yu, K.-Y. Huang, Pornography detection based on morphological features, *International Journal of Computer, Consumer and Control (IJ3C)* 2 (2013) 56–64.
- [47] A. Zaidan, N. Ahmad, H. Karim, M. Larbani, B. Zaidan, A. Sali, On the multi-agent learning neural and Bayesian methods in skin detector and pornography classifier: An automated anti-pornography system, *Neurocomputing* 131 (5) (2014) 397–418.
- [48] F. Souza, E. Valle, G. Cámara-Chávez, A. de A. Araújo, An evaluation on color invariant based local spatiotemporal features for action recognition, in: *25th Conference on Graphics, Patterns and Images (SIBGRAPI)*, 2012.
- [49] S. H. Yaghoubyan, M. A. Maarof, A. Zainal, M. F. Rohani, M. M. Og-haz, Fast and effective bag-of-visual-word model to pornographic images recognition using the freak descriptor, *Journal of Soft Computing and Decision Support Systems* 2 (6) (2015) 27–33.
- [50] J. Zhang, L. Sui, L. Zhuo, Z. Li, Y. Yang, An approach of bag-of-words based on visual attention model for pornographic images recognition in compressed domain, *Neurocomputing* 110 (13) (2013) 145–152.
- [51] G. Bradski, The OpenCV Library, *Dr. Dobb's Journal of Software Tools*.
- [52] M. Moustafa, Applying deep learning to classify pornographic images and videos, *ArXiv e-prints arXiv:1511.08899*.