

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/318408211>

Adult Content Detection in Videos with Convolutional and Recurrent Neural Networks

Article in *Neurocomputing* · July 2017

DOI: 10.1016/j.neucom.2017.07.012

CITATIONS

33

READS

3,480

4 authors:



Jônatas Wehrmann

Pontifícia Universidade Católica do Rio Grande do Sul

27 PUBLICATIONS 307 CITATIONS

[SEE PROFILE](#)



Gabriel Simões

Pontifícia Universidade Católica do Rio Grande do Sul

16 PUBLICATIONS 77 CITATIONS

[SEE PROFILE](#)



Rodrigo C. Barros

Pontifícia Universidade Católica do Rio Grande do Sul

105 PUBLICATIONS 1,570 CITATIONS

[SEE PROFILE](#)



Victor F Cavalcante

Motorola Mobility

28 PUBLICATIONS 122 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Imitation learning [View project](#)



Development of Fully-Flexible Receptor (FFR) Models for Molecular Docking [View project](#)

Adult Content Detection in Videos with Convolutional and Recurrent Neural Networks

Jônatas Wehrmann, Gabriel S. Simões, Rodrigo C. Barros*

*Pontifícia Universidade Católica do Rio Grande do Sul
Av. Ipiranga, 6681, 90619-900, Porto Alegre - RS, Brazil*

Victor F. Cavalcante*

*Motorola Mobility, R&D - Brazil
Rodovia SP 340 Km 128.7, 13820-000, Jaguariuna - SP, Brazil
victorfc@motorola.com*

Abstract

The amount of adult content on the Internet grows daily. Much of the pornographic content is unconstrained and freely-available for all users, requiring parents to make use of parental control strategies for protecting their children. Current parental control devices depend on human intervention, and hence there is the need of computational approaches for automatically detecting and blocking pornographic content. Towards that goal, this paper proposes *ACORDE*, a novel deep learning architecture that comprises both convolutional neural networks and LSTM recurrent networks for adult content detection in videos. Experiments over the freely-available NPDI dataset show that *ACORDE* significantly outperforms the previous state-of-the-art approaches for this task, decreasing by half the number of false positives and by a third the number of false negatives.

Keywords: adult content detection, deep learning, convolutional neural networks, recurrent neural networks

*Corresponding author

Email address: {jonatas.wehrmann,gabriel.simoes.001}@acad.pucrs.br, rodrigo.barros@pucrs.br (Jônatas Wehrmann, Gabriel S. Simões, Rodrigo C. Barros)

1. Introduction

The automatic detection of adult (pornographic) content in images and videos is an important and challenging task, especially due to the huge amount of freely-available adult content on the web, whose spread has significantly increased with the massive adoption of mobile devices across the globe. A recent report¹ indicates that the Internet traffic to porn websites accounted for 8.5% of the total in the UK in June 2013, surpassing the traffic for shopping, news, business, and social networks.

Even though organizations such as MPAA² have developed rating systems to protect viewers from adult scenes in motion pictures, content available on the web is practically unconstrained and easy-to-access, motivating the development of computational approaches that are capable of automatically detecting pornography with the final goal of protecting sensitive populations (e.g., children under 18). The task of automatically identifying adult content, however, poses a greater challenge than other classification problems due to the degree of subjectivity and uncertainty surrounding the problem. For instance, it is hard even for human beings to properly assess degrees of sensuality in scenes where people wear swimsuits or underwear. Indeed, sometimes more than one image/frame is needed for contextualizing the scene in order to define whether it should be classified as adult content or not.

Earlier work on pornography identification focused on human skin detection [1, 2, 3, 4], in which the idea is that greater amounts of detected skin would lead to higher probabilities of nudity within the image or video, hence characterizing the content as pornographic. Nevertheless, these approaches suffer with a high rate of false positives, especially in the context of beaches or practice of aquatic sports. More recent studies [5, 6, 7, 8] approached the problem under the perspective of *Bag of visual Words* (BoW) and similar models (e.g., BossaNova [9, 8]) for aggregating (quantizing) sophisticated image descriptors.

¹<http://goo.gl/nG0s7n>.

²<http://www.mpa.org/>

For benchmarking the proposed approaches in the area in terms of both
 30 video and image detection, researchers have used the NPDI dataset [8]. The
 best results achieved in NPDI are described by [10], where the authors propose
 a video descriptor based on binary features (*BinBoost* [11]) which is used with
 the BoW/BossaNova representations. However, the very same approach reaches
 only 44.6% of mean average precision (mAP) in the well-known PASCAL VOC
 35 dataset [12], while recent deep learning approaches reach about 60% of mAP
 in that same dataset [13]. This is a clear indication that deep learning based
 approaches could be a good option for pornography detection in both images
 and videos.

Therefore, in this paper we propose a novel approach for adult content detec-
 40 tion in videos, namely *ACORDE* (Adult Content Recognition with Deep Neu-
 ral Networks). Its architecture makes use of a Convolutional Neural Network
 (ConvNet) as a feature extractor and of a Long Short-Term Memory (LSTM) to
 perform the final video classification. *ACORDE* extracts feature vectors from
 the video *keyframes* of NPDI, building a sorted set of semantic descriptors. This
 45 set is used to feed the LSTM that is responsible for analyzing the video in an
 end-to-end fashion. The proposed approach does not require fine-tuning nor
 re-training the ConvNet. Results show that *ACORDE* comfortably establishes
 the new state-of-the-art for adult content detection in NPDI, reducing by half
 the number of false positives and by a third the false negatives.

50 This paper is organized as follows. Section 2 briefly introduces the NPDI
 dataset as well as recent methods for pornographic classification of videos. Sec-
 tion 3 describes our proposed approach in detail. Section 4 presents how the
 experimental setup was organized for performing the empirical analysis, which
 is presented in Section 5. Finally, in Section 6 we detail our conclusions and
 55 future work directions.

2. Background

This section discusses earlier work that perform adult content detection, and also describes the NPDI dataset, which will be used to validate our novel approach.

2.1. NPDI dataset

Currently, the largest publicly-available pornographic dataset is NPDI [9], which comprises nearly 80 hours from 802 videos (half of them with adult content), all downloaded from the Internet. The non-adult class is further subdivided in 201 easy-to-classify videos and 200 hard-to-classify videos. The latter
65 were selected based on textual search queries like *beach*, *wrestling*, and *swimming*, in order to verify the ability of the proposed classifiers in scenarios of high skin-exposure. The adult class comprises 401 videos selected from adult content web sites. Figure 1 shows a sample of frames from the easy non-adult, hard non-adult, and adult classes.



Figure 1: Frames from the NPDI dataset.

70 As described in [9], a scene segmentation algorithm was employed to extract keyframes from the videos, resulting in a total of 16,727 images. Each video may contain 1 to 320 keyframes. The average amount of keyframes per class are: 15.6 for adult videos; 33.8 for easy non-adult videos; and 17.5 for hard non-adult videos. NPDI has a wide ethnic diversity with asian, black, white, and
75 multi-ethnic videos. Issues like one-keyframe videos and *anime*-style content considerably increase the challenge of NPDI.

2.2. Related Work

In the work of [9] and [10], the authors make use of both low and mid-level visual features extracted from the NPDI dataset. They use such features to build a final movie representation. The method is based in a low-complexity alternative for feature extraction using binary descriptors and a combination of mid-level representations. They aggregate the descriptors via the BoW model [14], generating the *BoW Video Descriptor* (BoW-VD). Also, they use the BossaNova method [15], which is an improved extension of the BoW model, generating the *BossaNova Video Descriptor* (BNVD). BNVD is a video descriptor that represents the median distance for each visual word of a given *codeword*³ for a *codebook*⁴.

The work of [16] is the first to use deep neural networks to address the pornography detection problem. That work proposes a method that requires fine-tuning two distinct ConvNets, namely *AlexNet* [17] and *GoogLeNet* [18]. The author performs the training phase by reusing models pre-trained over the ImageNet dataset [19] and fine-tunes them over NPDI. That approach requires the training of ten distinct models: one model per training fold (5) and per network (2). Keyframes were rescaled to 256×256 to allow the data augmentation process with crops of the size 224×224 randomly sampled from each image in order to avoid overfitting. To normalize the data, the author subtracted the mean image from all instances.

Unfortunately, several methodological aspects are not clearly detailed in the paper, such as: i) the stopping criteria adopted, ii) the usage of a validation set, ii) values of important hyper-parameters like learning rate, momentum, and regularization; and iii) the updated layers in each model. Note that the absence of a proper validation set may compromise the reliability of the results. For the test phase, each network predicts *benign* (non-adult) and *adult* probabilities for each keyframe. The probabilities from both models are averaged, and a

³A codeword is the centroid of a given visual words cluster.

⁴A codebook is a set of codewords.

105 video is classified as adult (benign) when most of its keyframes are predicted as belonging to the *adult (benign)* class.

3. *ACORDE*

In this paper we propose a novel method for adult content detection in images and videos, namely *ACORDE* (Adult Content Recognition with Deep
110 Neural Networks). The architecture of *ACORDE* comprises a convolutional neural network (ConvNet) [20] for feature extraction and a Long Short-Term Memory network (LSTM) for sequence learning [21]. ConvNets are the current state-of-the-art for many computer vision tasks such as image classification [22], object detection [13], video analysis [23, 24, 25, 26] and image segmentation [27].
115 LSTMs are well suited to learn representations of sequences such as videos and texts. The conjoint use of both algorithms has been used to solve problems in video analysis [28] and scene captioning [29].

A ConvNet is a deep learning strategy that combines three ideas to ensure some degree of shifting, scale, and distortion invariance regarding the image
120 content: local receptive fields (filters), shared weights, and spatial (or temporal) pooling [30]. The convolution operator is applied in order to replace fully-connected matrix multiplications, granting the two first mentioned ideas and considerably reducing the amount of parameters within a network. Convolutional filters are learned using the well-known backpropagation algorithm
125 [31]. This process can be seen as *representation learning*, i.e., the network acting as a feature extractor. Learning representations from images is vital for the success of the computer vision task at hand. Eq. 1 defines a convolution, where (x, y) is a position on the j^{th} feature map from the i^{th} network layer; m indexes the set of feature maps, b_{ij} is the corresponding bias value, w_{ijm}^{pq} is the weight's value at position (p, q) , and P_i and Q_i are the height and width of the filter, respectively. The ReLU (Rectifier Linear Unit) activation function [17] is
130 often used as a source of non-linearity, essentially thresholding values in zero:

$relu(v) = \max(0, v)$.

$$v_{ij}^{xy} = \text{relu} \left(b_{ij} + \sum_m \sum_{p=0}^{P_i-1} \sum_{q=0}^{Q_i-1} w_{ijm}^{pq} v_{(i-1)m}^{(x+p)(y+q)} \right) \quad (1)$$

Recurrent Neural Networks (RNNs) [32] are especially designed to learn
 135 sequential or time-varying patterns. LSTMs [21], for instance, are a specific
 type of RNN that uses a basic unit called *memory cell* to allow the constant
 error flow through time. A memory cell comprises gates that are designed to
 protect the cell from irrelevant inputs and from irrelevant content within the
 cell.

As in [33], the LSTM implemented within *ACORDE* is defined by the fol-
 140 lowing components: block input (Eq. 2), input gate (Eq. 3), forget gate (Eq. 4),
 cell state (Eq. 5), output gate (Eq. 6) and block output (Eq. 7), where \mathbf{x}^t is the
 input vector at time t , \mathbf{W} are the weight matrices connected to the input, \mathbf{R}
 are the recurrent weight matrices, \mathbf{p} are peephole weight vectors, and \mathbf{b} are bias
 145 vectors. Non-linear functions are denoted by σ , g and h . The sigmoid function
 is used in the gates, whereas the hyperbolic tangent is used in the block input
 and output. For short, the LSTM internal state is denoted by A .

$$\mathbf{z}_t = g(\mathbf{W}_z \mathbf{x}_t + \mathbf{R}_z \mathbf{y}_{t-1} + \mathbf{b}_z) \quad (2)$$

$$\mathbf{i}_t = \sigma(\mathbf{W}_i \mathbf{x}_t + \mathbf{R}_i \mathbf{y}_{t-1} + \mathbf{p}_i \cdot \mathbf{c}_{t-1} + \mathbf{b}_i) \quad (3)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_f \mathbf{x}_t + \mathbf{R}_f \mathbf{y}_{t-1} + \mathbf{p}_f \cdot \mathbf{c}_{t-1} + \mathbf{b}_f) \quad (4)$$

$$\mathbf{c}_t = \mathbf{i}_t \cdot \mathbf{z}_t + \mathbf{f}_t \cdot \mathbf{c}_{t-1} \quad (5)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_o \mathbf{x}_t + \mathbf{R}_o \mathbf{y}_{t-1} + \mathbf{p}_o \cdot \mathbf{c}_t + \mathbf{b}_o) \quad (6)$$

$$\mathbf{y}_t = \mathbf{o}_t \cdot h(\mathbf{c}_t) \quad (7)$$

Figure 2 depicts *ACORDE*'s architecture. Let $\{\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_n\} \in \mathcal{F}$ be a
 keyframe sequence, \mathbf{x}_t be the set of features extracted from the t^{th} frame by the
 150 ConvNet, and \mathbf{y}_t the predicted probability of existing adult content at the t^{th}
 temporal iteration. The final prediction $y_{|\mathcal{F}|}$ given by *ACORDE* uses as features

the LSTM's internal state A from its final iteration. State A has recurrent changes given new inputs and the outputs of the internal gates. Note that in Figure 2 the recurrence is unrolled. Thus, when $t = 3$, there is no $\{\mathbf{x}_t, \mathbf{y}_t\} \forall (t > 3)$.

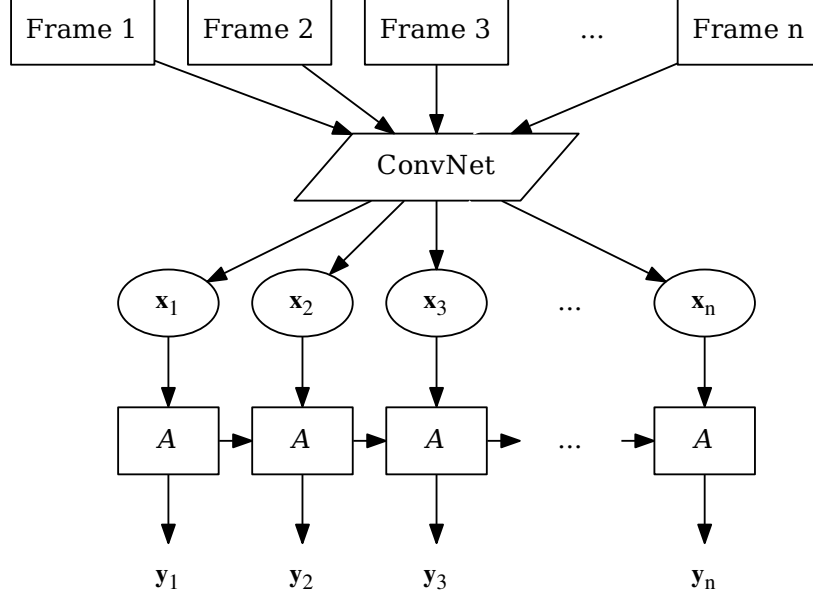


Figure 2: *ACORDE* architecture.

Since *ACORDE* makes use of a RNN to analyze each frame, we can obtain probability y_i iteratively. This is useful for detecting and blocking videos in real-time applications. The work of [16] has difficulties in classifying video streams because of its majority voting strategy, which harms predictions performed in real-time. For instance, a video containing only a few pornographic scenes will most probably be misclassified. On the other hand, methods based on histogram representations (e.g., [10]) depend on the training of the final classifier (often a non-linear SVM) to perform accurate intermediate predictions. Training a classifier based on the entire videos' histograms introduces a bias, in which the induced model may expect a large amount of explicit content at test time. Such an assumption is probably false, especially for the initial stages of the adult

video.

3.1. Image feature extraction

Training deep models in large datasets such as ImageNet [19] generates dis-
 170 criminative models capable of encoding semantic information from the images.
ACORDE makes use of pre-trained ConvNets to extract representative feature
 vectors \mathbf{x}_i from images/frames. We have experimented with both GoogleNet [18]
 and ResNet [22] architectures. GoogleNet is a compact ConvNet with the same
 prediction power of a VGG-19 [34], whereas ResNet holds the actual state-of-
 175 the-art for the ImageNet challenge. Common ResNet incarnations are composed
 by 52, 101, and 152 layers.

ACORDE extracts 1024 features from the last convolutional layer from
 GoogleNet and 2048 features from ResNets 52, 101, and 152. It normalizes the
 images by subtracting the RGB mean from the available pornographic dataset
 180 (in our experiments, NPDI). For generating more robust features, *ACORDE* em-
 ploys 10 crops from each original frame: top-left, top-right, bottom-left, bottom-
 right, and center (and then the same crops but horizontally mirrored). The final
 features are defined by the average of the vectors extracted from the 10 crops.

Table 1 shows the amount of parameters in well-known ConvNet architec-
 185 tures. The model proposed by [16] is composed by both an AlexNet and a
 GoogleNet, totalizing $\approx 72M$ of parameters, much more than the networks ex-
 plored in *ACORDE*. In addition, *ACORDE* requires only the forward pass of
 the architectures, since it does not train a ConvNet model over the available
 pornographic data.

Table 1: Amount of parameters in well-known ConvNets.

| Network | <i>#Parameters</i> |
|-----------------|--------------------|
| AlexNet [17] | $\approx 66M$ |
| GoogleNet [18] | $\approx 6M$ |
| ResNet-50 [22] | $\approx 25M$ |
| ResNet-101 [22] | $\approx 44M$ |
| ResNet-152 [22] | $\approx 60M$ |

190 To evaluate the discriminative power of the extracted features, we trained a linear SVM ($C = 1$) at image level for the NPDI data. We have labeled each keyframe based on the respective video class. Note, however, that not every frame in an adult video has adult content. Therefore, the results presented in Table 2 should be interpreted as the capability of each feature extractor to
195 recognize frames from adult movies (and not of recognizing adult content per se). The results presented in Table 2 are the average accuracy obtained per fold. Note that the multiple crops improve the resulting accuracy in all cases. Results show that the high-level features extracted by the deep networks have similar discriminative power, but all of them outperform the well-known SIFT (BoW)
200 descriptor by $\approx 20\%$. SIFT (BOF) was generated by using 100,000 randomly sampled points clustered in 128 keywords.

Table 2: Average accuracy (%) \pm standard deviation from different feature extractors.

| Feature Extractor | 1 Crop | 10 Crops |
|-------------------|---------------|---------------------------------|
| GoogleNet | 86.89 ± 1 | 89.34 ± 1 |
| ResNet-50 | 89.07 ± 1 | 90.01 ± 1 |
| ResNet-101 | 88.68 ± 1 | 90.27 ± 1 |
| ResNet-152 | 87.91 ± 0 | 89.28 ± 1 |
| SIFT (BOF) | 67.68 ± 3 | — |

3.2. Video-based learning

ACORDE's architecture is trainable in an end-to-end fashion and it accepts variable-length inputs. Whereas one does not need to re-train the ConvNet over
205 the pornographic dataset, the LSTM must be trained in the available data (e.g., NPDI). The parameter updating process of the LSTM is performed by using the full gradient through time [35], which means the model can learn complex long-term relationships among frames.

Let $\mathcal{F} = \{\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_n\}$ be a video composed by n RGB frames, and
210 $\mathbf{x}_t = \phi(\mathcal{F}_t)$ the features extracted from frame t after the last convolutional layer's activation. The process performed by ϕ is time-independent and can be massively parallelized. The final prediction step is then represented by

$y_n = \varphi(\mathbf{x}_n, \varphi(\mathbf{x}_{n-1}))$, where φ indicates the LSTM recurrence. Note that the forward pass of the n^{th} frame is equivalent to a very deep neural network with
215 n layers performing several nonlinearities.

The use of an LSTM network as a sequence-learner causes values of y_n to be generated by using information from all previous frames $\mathcal{F}_t \forall t \in \{1, 2, \dots, n-1\}$. The assumption here is that the LSTM hidden state A that has been used to compute y_n contains information collected from the video as a whole.

220 To summarize, the advantages of using an LSTM within *ACORDE* are twofold: i) allows for the learning of fixed-sized representations of variable-length inputs in an end-to-end fashion; and ii) allows for the mapping of long-term dependencies within frames.

4. Experimental Methodology

225 In this section, we present the experimental methodology that we employed for evaluating *ACORDE*'s performance. We describe the baseline algorithms that are compared with *ACORDE* in Section 5.2, the evaluation measures that assess the quality of the classifiers in Section 4.2, and the hyper-parameters required for running *ACORDE* in Section 4.3

230 4.1. Baseline Algorithms

We compare *ACORDE* with five classification strategies that have been previously applied for adult video classification: (i) BossaNova-HueSIFT [8]; (ii) BossaNova-BRISK [9], (iii) BNVD [9]; (iv) BoW-VD [10]; and (v) AGbNet [16].

235 4.2. Evaluation Measures

As discussed in Section 3, the outputs of *ACORDE* for each class are probability values, and the same is true for the baseline algorithms. Hence, the final predictions are often generated after thresholding these probability values in 50%. However, this choice is arbitrary and defining an optimal threshold is
240 difficult and subjective. Hence, we avoid choosing thresholds by employing the

Receiver Operating Characteristic (ROC) curve as the evaluation criterion for comparing the different approaches. For generating a ROC curve for a given classification method, one must select a predefined number of different thresholds within $[0,1]$ to be applied over the outputs of each method. Finally, the true positive rate is plotted in function of the false positive rate for different cut-off points. The interpolation of these points generates a ROC-curve, and then we use the area under such a curve ($AU(ROC)$) as the indication of quantitative performance obtained by each method. In addition, we also show the values of accuracy (thresholding the probabilities in 50%) since the classes in NPDI are properly balanced.

4.3. Hyper-Parameters Settings

Table 3 shows the setup that has been used in the ImageNet-based ConvNet training for the GoogleNet and ResNet architectures, as well as the parameters for the training of the LSTM over the NPDI dataset. Note that these parameters are the default of the original papers and that no effort has been made in order to tune them.

Table 3: Hyper-parameters setup.

| Hyper-parameter | GoogleNet | ResNets | LSTM |
|--------------------------------|--------------------|-------------------------|--------------------|
| Optimizer | SGD | SGD | Adam |
| Learning Rate (α) | 1×10^{-2} | 1×10^{-1} | 1×10^{-3} |
| Decay of α (γ) | 4% every 8 epochs | 10× when error plateaus | — |
| Momentum | 0.9 | 0.9 | — |
| Weight Decay | 2×10^{-4} | 1×10^{-4} | 1×10^{-4} |
| Weight Initialization | [36] | [36] | [36] |

Regarding the LSTM training, we use separated folds for validation and testing. *ACORDE* learns the model by using three out of the five folds. The training stops when the loss function plateaus for ten consecutive epochs. The best model is chosen by using the validation set as a proxy of the test set.

5. Experiments and Discussion

In this section we present the experimental analysis that was performed in order to evaluate *ACORDE*. In Section 5.1, we evaluate the importance of coupling the LSTM within *ACORDE*'s architecture, whereas in Section 5.2 we
265 compare *ACORDE* with the current state-of-the-art in the NPDI dataset. In all experiments, we show the performance of four distinct versions of *ACORDE*: *ACORDE-GN* (based on the GoogleNet architecture), *ACORDE-50* (based on the ResNet-50 architecture), *ACORDE-101* (based on the ResNet-101 architecture), and *ACORDE-152* (based on the ResNet-152 architecture).

270 5.1. Evaluating *ACORDE*'s Architecture

In order to show the gains provided by *ACORDE* when coupling the LSTM in the output of the convolutional neural network, we compare the performance of *ACORDE* within the NPDI dataset with convolutional-only approaches that employ either average pooling or max pooling over the features of the videos'
275 keyframes. In addition, we also compare *ACORDE* with the use of only the convolutional network with per-frame predictions, in which the final prediction is simply given by majority voting of the per-frame predictions. Table 4 shows the results of such an analysis. Note that all *ACORDE*'s versions outperform their corresponding convolutional-only counterparts, clearly indicating
280 that *ACORDE*'s architecture with the coupled LSTM is indeed beneficial for the video classification problem.

5.2. State-of-the-Art Comparison

In this section, we show the results obtained by comparing the predictive performance of the following algorithms: BossaNova-HueSIFT, BossaNova-BRISK,
285 BNVD, BoW-VD, AGbNet, and all versions of *ACORDE*.

Table 5 shows the results regarding both accuracy and $AU(ROC)$. For simplicity, the table is split into three sections: i) baselines results; ii) *ACORDE*'s results with no cropping; and iii) *ACORDE*'s results with cropping.

Table 4: Comparison between convolutional-only networks and *ACORDE* for video classification in NPDI.

| CNN | Aggregation Strategy | Accuracy (%) |
|-------------------|----------------------|--------------------------------|
| GoogleNet CNN | Average Pooling | 88.4 ± 3 |
| GoogleNet CNN | Max Pooling | 90.8 ± 1 |
| ResNet-50 CNN | Average Pooling | 92.1 ± 3 |
| ResNet-50 CNN | Max Pooling | 93.8 ± 3 |
| ResNet-101 CNN | Average Pooling | 92.1 ± 3 |
| ResNet-101 CNN | Max Pooling | 93.9 ± 1 |
| ResNet-152 CNN | Average Pooling | 92.9 ± 2 |
| ResNet-152 CNN | Max Pooling | 93.9 ± 2 |
| GoogleNet CNN | Majority Voting | 90.5 ± 3 |
| ResNet-50 CNN | Majority Voting | 92.8 ± 3 |
| ResNet-101 CNN | Majority Voting | 92.8 ± 2 |
| ResNet-152 CNN | Majority Voting | 91.6 ± 2 |
| <i>ACORDE-GN</i> | | 92.8 ± 1 |
| <i>ACORDE-50</i> | | 94.1 ± 1 |
| <i>ACORDE-101</i> | | 94.0 ± 1 |
| <i>ACORDE-152</i> | | 94.5 ± 1 |

Table 5: Results for video classification. (*) denotes the use of 10 crops during feature extraction. In bold the results that outperform the current state-of-the-art, and the best result achieved in the NPDI dataset are underlined.

| Method | Accuracy (%) | <i>AU(ROC)</i> |
|-----------------------|--------------------------------|----------------|
| BossaNova-HueSIFT [8] | 89.5 ± 1 | 0.954 |
| BossaNova-BRISK [9] | 88.6 ± 2 | 0.960 |
| BNVD [9] | 92.0 ± 1 | 0.973 |
| BoW-VD [10] | 92.4 ± 2 | 0.976 |
| AGbNet [16] | 94.1 ± 2 | - |
| <i>ACORDE-GN</i> | 92.8 ± 1 | 0.978 |
| <i>ACORDE-50</i> | 94.1 ± 1 | 0.986 |
| <i>ACORDE-101</i> | 94.0 ± 1 | 0.985 |
| <i>ACORDE-152</i> | 94.5 ± 1 | 0.986 |
| <i>ACORDE-GN*</i> | 93.3 ± 2 | 0.981 |
| <i>ACORDE-50*</i> | 94.8 ± 2 | 0.988 |
| <i>ACORDE-101*</i> | <u>95.6 ± 1</u> | <u>0.990</u> |
| <i>ACORDE-152*</i> | 95.3 ± 1 | <u>0.990</u> |

Regarding the baselines, the maximum accuracy that is reached is $94.1\% \pm 2$
 290 for AGbNet, which employs two ConvNets for classifying images and a majority voting scheme for classifying the videos. Both second and third placed methods are from [10], namely BNVD and BoW-VD, reaching $\approx 92\%$ of accuracy. Note that even *ACORDE*'s weakest approach, namely *ACORDE-GN*, outperforms all non-network baselines. It is outperformed by AGbNet, but recall that *ACORDE-GN* contains $\approx 12\times$ less parameters than AGbNet, as well
 295 as a more stable behaviour. Overall, five out of the eight *ACORDE* variations outperform all baselines for all evaluation measures.

ACORDE-101 and *ACORDE-152* achieve the largest accuracy values, $\approx 95.5\%$. These results surpass the best baseline approach by $\approx 1.5\%$ and the
 300 second-best by $\approx 3.2\%$. The *ACORDE* variations that employ residual connections achieve an $AU(ROC)$ of around 0.99. Given the insufficient description of AGbNet's parameters, we could not reproduce it in order to generate its respective $AU(ROC)$.

The use of the multiple-crop strategy (denoted by *) improves the predictive performance at video level for all *ACORDE*'s variations. The improvement
 305 ranges between 0.5% and 1.6%, with a greater impact for the deeper architectures. Shallower architectures such as *ACORDE-GN* and *ACORDE-50* presented higher variability when using multiple crops.

Figure 3 presents randomly selected keyframes from the videos misclassified
 310 by *ACORDE-101*. The first three keyframes (Figure 3-(a)) are examples of the non-adult class and the last three (Figure 3-(b)) of the adult class. Videos of breastfeeding babies usually present the exposure of women's breasts, which can be misleading to the classifier. *ACORDE* generates a probability of $y = 54\%$ of the first video being pornographic, leading to the incorrect prediction. The second video presents a game in which women in underwear have to perform certain
 315 moves and positions, though it is labeled as non-adult. *ACORDE* classifies that video as adult with $y = 83\%$ of probability. The last false positive is regarding a video with a single tricky keyframe. We highlight the fact that *ACORDE-101* perfectly classified all beach and sumo videos, which are often misclassified by

the baseline approaches. We also noticed that misclassified videos with large probability scores are very rare.



Figure 3: *ACORDE*-101 misclassified videos.

The lack of exposure of intimate body parts is the most frequent characteristic in the false negatives. For instance, both fourth and fifth videos are 1-keyframe videos in which the intimate body parts are mostly hidden. We do believe that training a model with all movie frames could minimize those classification errors. The last keyframe shown in Figure 3 is from a long *anime*-style video with very few explicit scenes. *Anime* content is much less frequent in the training data and greatly differs from the rest of the dataset, hence affecting the generalization performance of *ACORDE*.

Figure 4 shows a 2-dimensional *t-distributed stochastic neighbor* plotting (*t*-SNE) [37] for further analysis of *ACORDE* and the baseline methods. *t*-SNE is a dimensionality reduction technique for high-dimensional data visualization. Figure 4-(a) presents the plotting of the video-based features extracted from *ACORDE*'s LSTM. Figures 4-(b) and (c) present, respectively, BNVD's and BoW-VD's generated video features. Note that *ACORDE* generates discriminative features that practically allows a linear separation of porn and non-porn videos. An interesting fact is also the possibility of separating the easy and hard non-adult subclasses, which are automatically recognized by *ACORDE* without explicitly training with these categories.

Figures 5, 6, and 7 show the confusion matrix for the best *ACORDE* version (*ACORDE* -101) and the baseline non-network methods. Note that *ACORDE*-101 decreases by half the number of false positives of BNVD and BOW-VD, clearly indicating that the high-level features from the ConvNet and the sequence learning of the LSTM are more robust in identifying adult content in

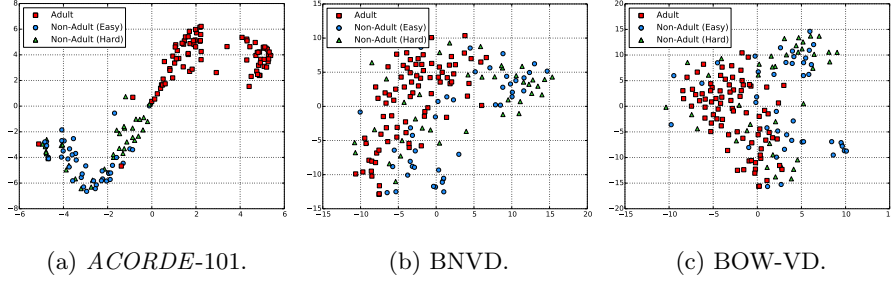


Figure 4: t -SNE plots.

345 scenarios with large skin-exposure. In addition, *ACORDE* -101 also decreases by a third the number of false negatives of the baseline approaches (20 versus 33), proving to be safer to use in a parental control device than the previous state-of-the-art approaches.

| | Predicted Class | |
|--------------|-----------------|-----------|
| | Adult | Non-Adult |
| Actual Class | Adult | 381 |
| | Non-Adult | 16 |

Figure 5: *ACORDE*-101 confusion matrix.

| | Predicted Class | |
|--------------|-----------------|-----------|
| | Adult | Non-Adult |
| Actual Class | Adult | 370 |
| | Non-Adult | 33 |

Figure 6: BNVD confusion matrix.

6. Conclusions and Future Work

350 In this paper we proposed *ACORDE*, a novel deep neural network architecture that comprises both a ConvNet and an LSTM for adult video classification. To the best of our knowledge, *ACORDE* is the first method in the literature that makes use of such a type of architecture for detecting adult content in

| Actual Class | Predicted Class | |
|--------------|-----------------|-----------|
| | Adult | Non-Adult |
| | Adult | Non-Adult |
| Adult | 373 | 28 |
| Non-Adult | 33 | 368 |

Figure 7: BOW-VD confusion matrix.

videos. We performed several experiments in the NPDI pornography dataset
 355 for verifying the best design choices for *ACORDE*. After a thorough empirical
 analysis, most of the *ACORDE*'s variations were capable of outperforming the
 current state-of-the-art methods for adult content detection.

As future work, we intend to make available a novel image-based adult
 dataset, which will be the largest dataset publicly-released to date. Moreover,
 360 we would like to verify which is the best transfer learning strategy for using
 in the pornography context. Finally, we want to develop a novel method for
 real-time segmentation of body parts in adult videos.

Acknowledgment

We would like to thank Motorola Mobility and the Brazilian research agen-
 365 cies CAPES and CNPq for funding this work. In addition, we gratefully ac-
 knowledge the support of NVIDIA Corporation with the donation of the Tesla
 K40 GPUs used for this research.

- [1] M. M. Fleck, D. A. Forsyth, C. Bregler, Finding naked people, in: 4th
 European Conference on Computer Vision, 1996, pp. 593–602.
- 370 [2] M. J. Jones, J. M. Rehg, Statistical color models with application to skin
 detection, in: IEEE Computer Society Conference on Computer Vision and
 Pattern Recognition, Vol. 1, 1999, p. 280 Vol. 1.
- [3] 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, part Special Issue on Visual Information Pro-
 375 cessing.

- [4] H. Zuo, W. Hu, O. Wu, Patch-based skin color detection and its application to pornography image filtering, in: Proceedings of the 19th International Conference on World Wide Web, WWW '10, ACM, New York, NY, USA, 2010, pp. 1227–1228.
- [5] T. Deselaers, L. Pimenidis, H. Ney, Bag-of-visual-words models for adult image classification and filtering, in: International Conference on Pattern Recognition, 2008, pp. 1–4.
- [6] A. P. B. Lopes, S. E. F. de Avila, A. N. A. Peixoto, R. S. Oliveira, A. A. Araújo, A bag-of-features approach based on hue-sift descriptor for nude detection, in: European Signal Processing Conference, 2009.
- [7] A. P. B. Lopes, S. E. F. d. Avila, A. N. A. Peixoto, R. S. Oliveira, M. d. M. Coelho, A. d. A. Araújo, Nude detection in video using bag-of-visual-features, in: XXII Brazilian Symposium on Computer Graphics and Image Processing, 2009, pp. 224–231.
- [8] S. Avila, N. Thome, M. Cord, E. Valle, A. D. A. Araújo, Pooling in image representation: The visual codeword point of view, Computer Vision and Image Understanding 117 (5) (2013) 453–465.
- [9] C. Caetano, S. Avila, S. Guimaraes, A. d. A. Araújo, Pornography detection using bossanova video descriptor, in: 2014 22nd European Signal Processing Conference (EUSIPCO), IEEE, 2014, pp. 1681–1685.
- [10] C. Caetano, S. Avila, W. R. Schwartz, S. J. F. Guimarães, A. de A. Araújo, A mid-level video representation based on binary descriptors: A case study for pornography detection, Neurocomputing (2016) –.
- [11] T. Trzcinski, M. Christoudias, P. Fua, V. Lepetit, Boosting binary keypoint descriptors, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2013, pp. 2874–2881.

- [12] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, A. Zisserman, The pascal visual object classes challenge: A retrospective, *International Journal of Computer Vision* 111 (1) (2015) 98–136.
- [13] S. Ren, K. He, R. Girshick, J. Sun, Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks., *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2016) 1–1.
- [14] J. Yang, Y.-G. Jiang, A. G. Hauptmann, C.-W. Ngo, Evaluating bag-of-visual-words representations in scene classification, in: *Proceedings of the international workshop on Workshop on multimedia information retrieval*, ACM, 2007, pp. 197–206.
- [15] S. Avila, N. Thome, M. Cord, E. Valle, A. d. A. Araújo, Bossa: Extended bow formalism for image classification, in: *2011 18th IEEE International Conference on Image Processing*, IEEE, 2011, pp. 2909–2912.
- [16] M. Moustafa, Applying deep learning to classify pornographic images and videos, in: *7th Pacific-Rim Symposium on Image and Video Technology (PSIVT)*, 2015.
- [17] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [18] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, A. Rabinovich, Going deeper with convolutions, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1–9.
- [19] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, L. Fei-Fei, ImageNet Large Scale Visual Recognition Challenge, *International Journal of Computer Vision* 115 (3) (2015) 211–252.

- 430 [20] Y. Le Cun, B. Boser, J. S. Denker, R. E. Howard, W. Hubbard, L. D. Jackel, D. Henderson, Handwritten digit recognition with a back-propagation network, *Advances in neural information processing systems* 2.
- [21] S. Hochreiter, J. Schmidhuber, Long short-term memory, *Neural computation* 9 (8) (1997) 1735–1780.
- 435 [22] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [23] S. Ji, W. Xu, M. Yang, K. Yu, 3d convolutional neural networks for human action recognition, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35 (1) (2013) 221–231.
- 440 [24] K. Simonyan, A. Zisserman, Two-stream convolutional networks for action recognition in videos, in: *Advances in Neural Information Processing Systems*, 2014, pp. 568–576.
- [25] J. Wehrmann, G. Simões, B. Rodrigo, T. Paula, D. Ruiz, (deep) learning from frames, in: *Proceedings of the Brazilian Conference on Intelligent System*, 2016.
- 445 [26] G. S. Simões, J. Wehrmann, R. C. Barros, D. D. Ruiz, Movie Genre Classification with Convolutional Neural Networks, in: *International Joint Conference on Neural Networks (IJCNN 2016)*, 2016.
- 450 [27] G. Lin, C. Shen, I. D. Reid, A. van den Hengel, Efficient piecewise training of deep structured models for semantic segmentation, in: *IEEE Computer Vision and Pattern Recognition (CVPR) 2016*, 2016.
- [28] A. Rohrbach, M. Rohrbach, B. Schiele, The long-short story of movie description, in: *German Conference on Pattern Recognition*, Springer, 2015, pp. 209–221.
- 455

- [29] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, T. Darrell, Long-term recurrent convolutional networks for visual recognition and description, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 2625–2634.
- [30] I. Goodfellow, Y. Bengio, A. Courville, Deep learning, book in preparation for MIT Press (2016).
URL <http://www.deeplearningbook.org>
- [31] D. E. Rumelhart, G. E. Hinton, R. J. Williams, Learning internal representations by error propagation, Tech. rep., DTIC Document (1985).
- [32] L. Fausett (Ed.), Fundamentals of Neural Networks: Architectures, Algorithms, and Applications, Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1994.
- [33] K. Greff, R. K. Srivastava, J. Koutník, B. R. Steunebrink, J. Schmidhuber, Lstm: A search space odyssey, IEEE Transactions on Neural Networks and Learning Systems PP (99) (2016) 1–11.
- [34] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, in: International Conference on Learning Representations, 2015.
- [35] A. Graves, J. Schmidhuber, Framewise phoneme classification with bidirectional lstm and other neural network architectures, Neural Networks 18 (5) (2005) 602–610.
- [36] K. He, X. Zhang, S. Ren, J. Sun, Delving deep into rectifiers: Surpassing human-level performance on imagenet classification, in: The IEEE International Conference on Computer Vision (ICCV), 2015, pp. 1026–1034.
- [37] L. v. d. Maaten, G. Hinton, Visualizing data using t-sne, Journal of Machine Learning Research 9 (Nov) (2008) 2579–2605.