



HAL
open science

A Comparative Study of Tools for Explicit Content Detection in Images

Adrien Dubettier, Tanguy Gernot, Emmanuel Giguët, Christophe Rosenberger

► **To cite this version:**

Adrien Dubettier, Tanguy Gernot, Emmanuel Giguët, Christophe Rosenberger. A Comparative Study of Tools for Explicit Content Detection in Images. 2023 International Conference on Cyberworlds (CW 2023), Oct 2023, Sousse, Tunisia. pp.464-471, 10.1109/CW58918.2023.00077 . hal-04179978

HAL Id: hal-04179978

<https://hal.science/hal-04179978v1>

Submitted on 10 Aug 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

A Comparative Study of Tools for Explicit Content Detection in Images

Adrien Dubettier, Tanguy Gernot, Emmanuel Giguet, Christophe Rosenberger
Normandie Univ, UNICAEN, ENSICAEN, CNRS, GREYC, 14000 Caen, France
firstname.surname@unicaen.fr

Abstract—Cyberworlds offer a vast quantity of knowledge and services on all topics for Internet users. The protection of children is an important issue on Internet and could be solved by detecting automatically explicit content. Another application is to facilitate digital forensic experts when analyzing media such as hard drives to detect child pornography content in criminal affairs. In this work, we focus on images and we study the efficiency of existing methods from the literature mainly based on machine learning and deep learning approaches. We apply a rigorous protocol with significant datasets in order to draw conclusions on the performance we can expect in real conditions. This study shows that this task is not really solved by existing tools. Moreover, the frontier of explicit content is also not always easy to define.

Index Terms—Explicit Content Detection, Digital investigation, Deep Learning, protection of children in Cyberworlds.

I. INTRODUCTION

While digital technologies have opened up new perspectives for children, at the same time, these technologies have also exposed them to threats and dangers far beyond traditional childhood violence. One of the most critical threats is the alarming growth in online child sexual exploitation and abuse (CSEA). Between 1-20% of children (12-17 years of age) were subjected to CSEA in 2020 across 13 countries in Eastern and Southern Africa and South East Asia.

The risk of unwanted exposure to online porn for children increases with the development of digital tools and the intensification of the consumption of online content. Access to pornography for children and youth is doubly facilitated by the unsupervised communication on the Internet and the free and unrestricted access to illegal content of porn sites. Parental control can be activated on devices and most search engines offer filtration for website search results and for images that may appear as the result of a query.

However, the diversity of sources, from websites to social networks, applications on smartphones and the diversity of the explicit contents make impossible the development of a total protection shield for children. Legally speaking, things are changing, slowly but surely. Governments take initiative to block porn sites to protect minors. For instance, the exposure of pornographic photos and videos to minors is banned under a French domestic violence law passed in July 2020 [1]. The legislation specifies that companies may not exonerate themselves of responsibility simply by asking the internet user if they are over 18. Nevertheless, this measure

does not guarantee that only an adult public can access to the pornographic content.

From research and industry, much effort is dedicated and aims at improving the detection and filtering of explicit content. An explicit content detection engine allows to detect Not Suitable For Work (nsfw_model) media (i.e., image/ video) content. An Explicit content engine is generally based on deep learning models which return a probability to indicate the explicitness in media content. The big players on Internet propose their solutions mainly to verify uploaded images contents, we can cite SafeSearch API for Google [2], Azure Analyze Image API for Microsoft [3] or Rekognition API for Amazon [4]. It is difficult to assess their performance as the probability of explicitness is not always given. Concerning academic works, many researchers considered this machine learning task with classical features such as skin detectors [5]–[7] or convolutional networks [8]–[11]. These tools have been proposed in the last then years but their relative performance is not well known.

In this study, we focus on detecting offending content in images such as pornography or child pornography. Two main applications are targeted. The first one concerns the detection of explicit content to avoid children exposure on the Web. The second one concerns its application for digital investigation on a hard disk to detect child pornography. Our main contribution in this paper lies in an independent benchmarking of tools dedicated to explicit content detection in images. We built for that significant datasets composed of different situations.

The paper is organized as follows. In section II, we list different approaches in the literature that allow to detect explicit content in images. Section III is dedicated to the experimental protocol we follow in this study. We show in section IV the relative performance of different tools for explicit content detection. We analyze in detail obtained results in order to identify trends for this application. We conclude and give some perspectives in section V.

II. LITERATURE REVIEW

Sexually explicit image can be defined as an image depicting nudity or depicting any person engaging in sexual conduct. There are several existing methods available in the

TABLE I
COMPARISON OF TOOLS FOR EXPLICIT CONTENT DETECTION.

Name	Considered tools				
	<i>nsfw_model</i>	<i>NudeNet</i>	<i>NuDetective</i>	<i>SkinDetection</i>	<i>DeepPornDetection</i>
Year	2020	2019	2008	2022	2018
Language	Python	Python	Java	Python	Python
Principle	Machine learning model	Ensemble of neural nets	nudity detection	skin detection	transfer learning
References	[8]	[10]	[9]	[12]	[11]

literature to detect explicit content in images. Classic ones are based on computer vision algorithms and color models for skin identification and nudity detection. More recently, deep learning techniques based on transfer learning have arisen to separate explicit and non-explicit pictures. We present these approaches and review some related open source GitHub packages and commercial solutions.

A. Computer Vision Approaches

Computer vision algorithms combined with pattern recognition is a classical approach to Explicit Content Detection, underlying the idea that pornographic pictures include naked people that could be characterized using skin identification models. Color models are at the core of skin detection algorithms. RGB color model is the most well-known color model since it is a standard to display images on electronic devices. However, alternative models which facilitate the color distance computations of are preferred in computer vision. HSV is a model based on the perception of color similarity: it separates the image intensity from the color information. HSV color model has been used in order to discriminate elements that are not human skin, and then to detect skin color zones images [5]. Results have been improved with the YCbCr color space which is more appropriate for analytical purposes. Indeed, using a transformation from the RGB color space to the YCbCr color space in order to get the percentage of skin in an image, Basilio et al. [6] obtain an accuracy of 88.8% explicit content detection. [7] also use the YCbCr color space to detect skin. It is then combined with a Linear Discriminant Analysis to identify of pornographic contents.

Choosing *a priori* the appropriate color space model may not be a relevant approach since many factors can affect the results, including the skin color, the variation in skin tone according to different races, the light conditions. . . To handle this, authors in [13] propose a combination of different color space models to try to activate the most appropriate model depending on the context. The use of the color space model alone is unreliable since several types of pictures may contain a high percentage of bare skin without being explicit ones, such as a face in close-up (selfies), or family pictures on the beach. Conversely, there are some pictures that can depict explicit content while showing a low percentage of naked skin, such as clothed porn or from a distance. Xiaoyin Wang et al. [14] aim at detecting whether there is a naked body or not in a picture. They regard

a human body as the combination of some key rectangles such as limbs, face, trunk with the navel as use navel detection and Forward Propagation neural network to detect.

B. Deep Learning Approaches

With a large dataset including explicit and non-explicit pictures, it is possible to train a model with a convolutional network to classify these two types of picture. That was done by G. Laborde [8] (*nsfw_model*) and used transfer learning and fine tuning in order to obtain a model able to achieve an accuracy of 93%. A. Q. Bhatti et al [15], on their own dataset, trained a model (called *NudeNet*) that is capable of telling if an image is explicit or not, and to what degree if that is the case. They use the Resnet-50 architecture 1 by 2 convolutional deep learning neural network to give to a picture a score between 0 and 1 where 1 is the highest most explicit. They have an accuracy of 95% when tested on their dataset. The efficiency of these solutions is interesting and the *nsfw_model* achieves an accuracy of 93%, Authors in [15] reach an accuracy of 91% and *nudenet* does a little bit lower with 90%. In 2018, Alex Lykesas [11] proposes another transfer learning approach for this task with 12 layers. The obtained accuracy is claimed to reach 98% on a dataset composed of 15000 images (we use this dataset in this comparative study).

It is difficult to assess the relative performance of proposed methods in the literature as they have been evaluated on different datasets. It could also be interesting to better understand their results and errors. We considered 5 tools from the literature (see Table I). In this work, we try to answer this question in the following through a rigorous experimental protocol.

III. EXPERIMENTAL PROTOCOL

In order to compare different tools for explicit content detection, we need some datasets, tools to benchmark, some evaluation metrics and an evaluation scenario.

A. Datasets

We used different datasets in this work. They have different number of images, resolutions and difficulties. Some used images have a very low resolution (128×128 pixels) making the explicit detection task challenging. The frontier between sexy and porn images is not very clear and could be cultural dependent. We plan in the future to study this point with a

subjective evaluation to assess possible bias related to gender, culture or age.

1) *Dataset 200*: A first toy dataset has been created and is composed of 200 images. It has been built searching 100 images on Pixabay [16] for the safe part. Some adult websites provided us 100 images categorized as unsafe. Images contain 50,000 to 3 millions pixels and 10% of them are synthetic ones (for both classes). Figure 1 describes the content of the dataset.

2) *Dataset 4000*: A second larger dataset has been defined by retrieving content from different sites. Different websites such as [16], [17] were used for the unsafe part and [18], [19], [20] for the safe part. It is composed of 4000 images distributed as follows: 900 images containing nudity, 700 showing sexual intercourse, 400 sexual drawings or hentai, which brings the unsafe images to 2000, and 1000 pictures containing one or two persons, 400 pictures of a crowd, 300 faces and 100 pictures of something else, either landscape or animals. Image resolutions vary from 50,000 to 64 millions pixels and 10% of them are synthetic ones. Figure 2 describes the content of the dataset.

3) *Dataset 15000*: This dataset has been retrieved from GitHub [11]. It consists of 7500 images collected on tumbzilla for the unsafe part and 7500 pictures of people wearing bikini, collected on Google, for the safe part. No synthetic image is present in this dataset and images have all a low resolution (57600 pixels). Figure 3 describes the content of the dataset.

B. Tested tools

We present the 5 tested tools for explicit content detection in images. Table I summarizes the description of each tool.

1) *SkinDetection*: We wanted to test a "skin detection only" solution and used SkinDetection [12] that combines HSV and YCbCr filters in order to obtain the ratio of bare skin in an image. It is a simple solution used as baseline tool in this study.

2) *nsfw_model*: This model has been proposed by G. Laborde [8]. This approach is based on transfer learning and fine tuning. It returns the probability an image belongs to the following categories: drawing/neutral/sexy/porn/hentai. In order to compare it with the other tools, we had to reduce these five categories into two : safe and unsafe with the following formulas:

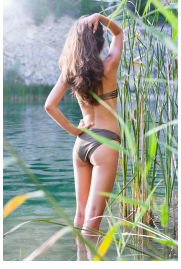

$$Safe = drawing + neutral + \frac{sexy}{2} \quad (1)$$

$$Unsafe = porn + hentai + \frac{sexy}{2} \quad (2)$$

We adopted this solution after doing some tests and we believe that it is the best solution due to the subjective nature of this category. Table II shows as illustration the output of this model

for 2 sexy images. The one on the left is defined as safe and the right one as unsafe. This illustration shows the difficulty to process particular situations (posture, little clothing...).

TABLE II
COMPARISON BETWEEN TWO "SEXY" PICTURES WITH NSFW_MODEL

		
<i>drawing</i>	0.079	0.099
<i>neutral</i>	0.084	0.078
<i>sexy</i>	0.715	0.616
<i>porn</i>	0.089	0.0045
<i>hentai</i>	0.033	0.162
<i>safe</i>	0.52	0.48
<i>unsafe</i>	0.48	0.52
<i>result</i>	safe	unsafe

3) *NudeNet*: The aim of the NudeNet Project [21] is to build an open source dataset for Nudity detection and to provide a pre-trained Deep Learning model for this task. The author also implements an Exposed Part Detection and Censoring module using Object Detection.

4) *DeepPornDetection*: This tool has been proposed by Alex Lykesas [11] in 2018. A 12-layers deep neural network has been trained (13,845,282 parameters) on the dataset 15000 we use in this work. Figure 4 shows an illustration of output of this tool on a set of images from the bikini dataset.

5) *NuDetective*: This Forensic Tool was developed in order to assist digital investigation examiners to conduct such analysis in a timely manner at the crime scene. This commercial tool performs the automatic detection of nudity in images and also performs analysis of file names [22]. NuDetective has been created in Brazil in 2010, following the amendment where possession of files containing child pornography is considered as a crime.

C. Evaluation metrics

As performance metrics, we consider the accuracy (correct recognition rate) given the ground truth. We also compute the confusion matrix showing the differences between predictions by a tool and the ground truth. We finally consider the computation time.

D. Testing scenarios

For all each tested tool, we compute the accuracy and the confusion matrix. Most tools return a class probability,

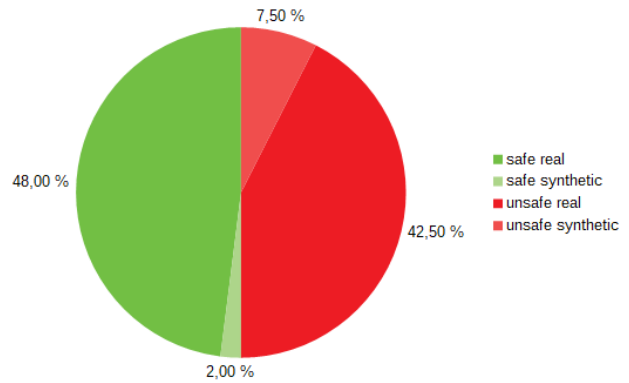
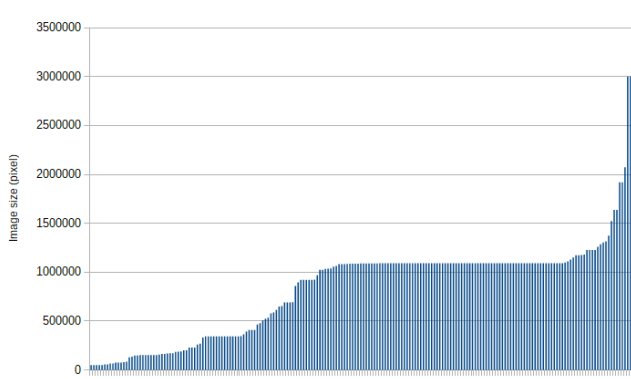


Fig. 1. Dataset 200: Distribution of image size and composition.

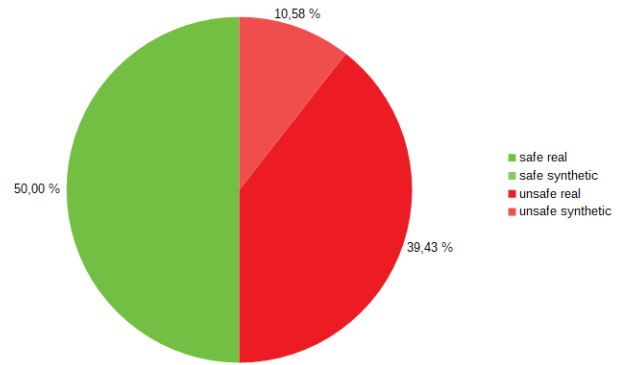
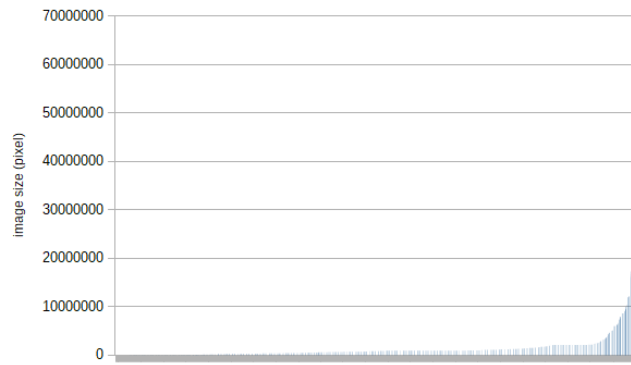


Fig. 2. Dataset 4000: Distribution of image size and composition.

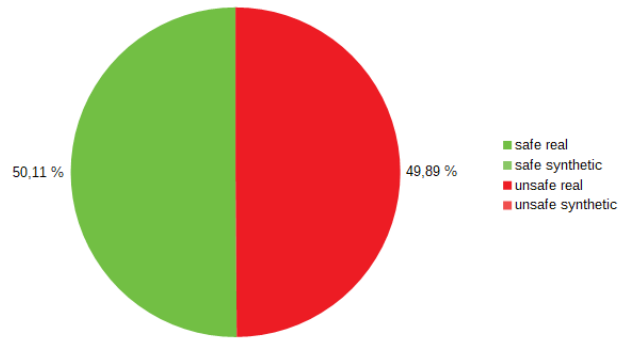
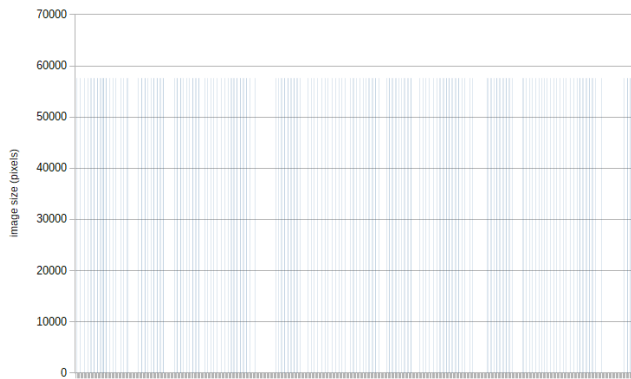


Fig. 3. Dataset 15000: Distribution of image size and composition.

we consider the maximal probability value for predicting the safe/unsafe class. The NuDetective solution does not give a confidence index but instead it marks as suspect the images it considers so.

IV. EXPERIMENTAL RESULTS

We show the performance of this comparative study by answering some questions.

A. What is the performance of tested tools?

Table III shows the confusion matrix for each tool on all datasets. Most tools generate errors similarly on each class

except SkinDetection and DeepPornDetection. SkinDetection tends to badly recognize unsafe images and DeepPornDetection safe ones.

Table IV presents the accuracy values for each tool on all datasets. We can see clearly that 3 tools provide poor results namely SkinDetection, DeepPornDetection and NuDetective on the 2 first datasets. DeepPornDetection provides very good results on the last dataset but this solution has been trained with this dataset, the result is thus biased. The accuracy on the Dataset 15000 is lower as the image resolution is low and

TABLE III
CONFUSION MATRIX FOR EACH TOOL ON ALL DATASETS.

Datasets		Predicted class			
		<i>safe</i>	<i>unsafe</i>		
Dataset 200	nsfw_model	<i>safe</i>	42%	8%	
		<i>unsafe</i>	9%	41%	
	NudeNet	<i>safe</i>	43.5%	6.5%	
		<i>unsafe</i>	10%	40%	
	NuDetective	<i>safe</i>	31.5%	18.5%	
		<i>unsafe</i>	17%	33%	
	SkinDetection	<i>safe</i>	46.5%	3.5%	
		<i>unsafe</i>	37.5%	12.5%	
	DeepPornDetection	<i>safe</i>	11.5%	38.5%	
		<i>unsafe</i>	0.5%	49.5%	
	Dataset 4000	nsfw_model	<i>safe</i>	45.75%	4.25%
			<i>unsafe</i>	10.7%	39.3%
NudeNet		<i>safe</i>	46.3%	3.7%	
		<i>unsafe</i>	5.9%	44.1%	
NuDetective		<i>safe</i>	32.1%	17.9%	
		<i>unsafe</i>	16.7%	33.3%	
SkinDetection		<i>safe</i>	39.3%	10.7%	
		<i>unsafe</i>	35.6%	14.4%	
DeepPornDetection		<i>safe</i>	5%	45%	
		<i>unsafe</i>	0%	50%	
Dataset 15000		nsfw_model	<i>safe</i>	26.6%	23.4%
			<i>unsafe</i>	4.5%	45.5%
	NudeNet	<i>safe</i>	13.6%	36.4%	
		<i>unsafe</i>	1.6%	48.4%	
	NuDetective	<i>safe</i>	11.2%	38.8%	
		<i>unsafe</i>	3.7%	46.3%	
	SkinDetection	<i>safe</i>	49.5%	0.5%	
		<i>unsafe</i>	23.5%	26.5%	
	DeepPornDetection	<i>safe</i>	48.7%	1.3%	
		<i>unsafe</i>	3.1%	46.9%	



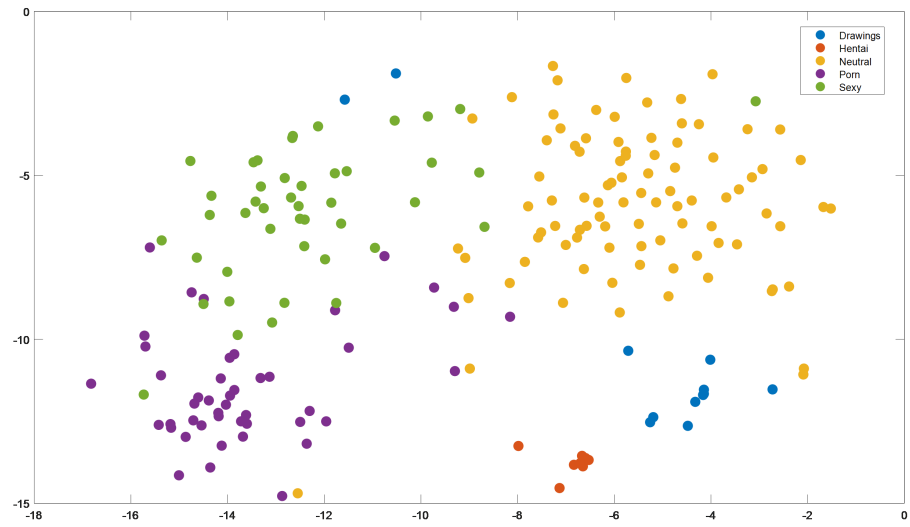
Fig. 4. Illustration of DeepPorndetection results on the subset of bikini dataset [11].

safe images contain many bikini images that are sometimes suggestive.

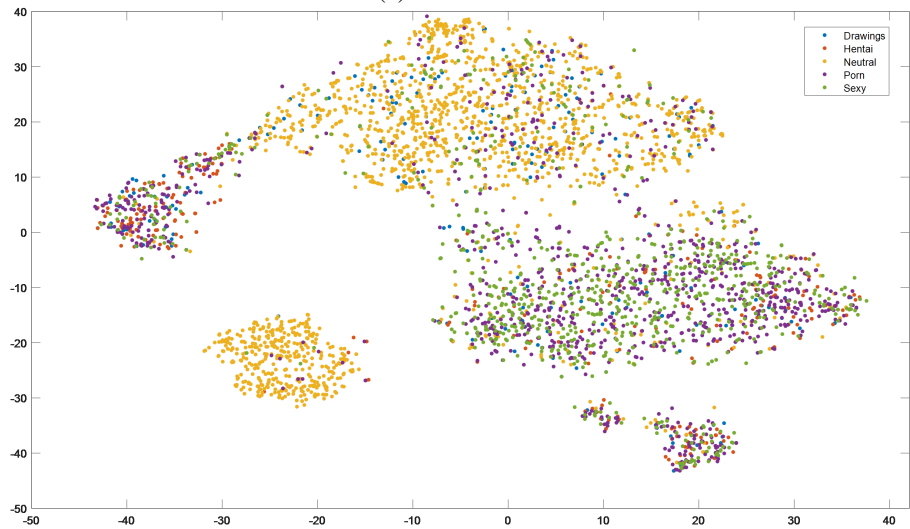
B. What is the discriminancy of deep features ?

As the nsfw_model is an open CNN architecture, it is possible to analyse the discriminancy of associated deep features. We considered the more precise categories (porn, hentai, sexy, neutral, drawings) used in the dataset 15000. We can obtain 1001 nsfw_model features, we plot in Figure 5 their projection in 2D for the three datasets considered in this study. This projection has been obtained by using the t-Distributed Stochastic Neighbor Embedding algorithm [23].

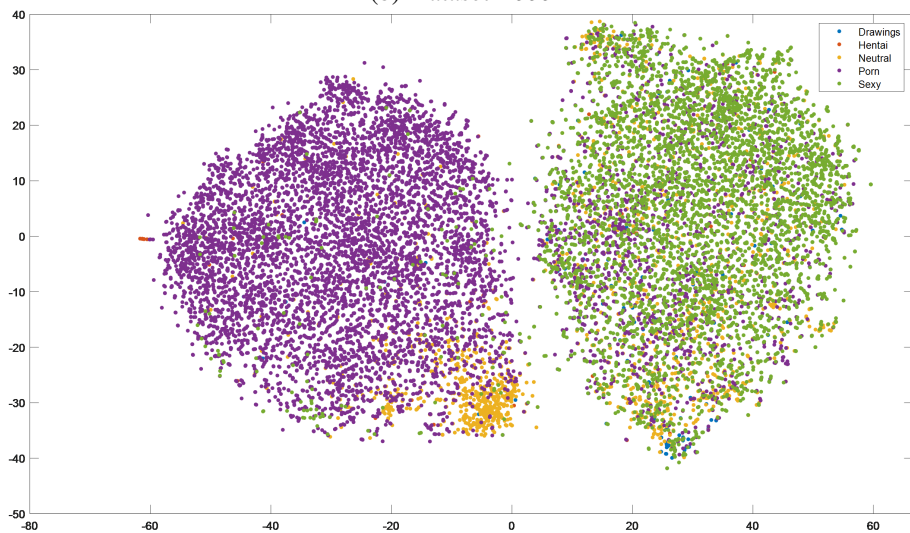
For the dataset 200 (see Figure 5 (a)), we can see that categories are well separated. We can see that in some cases, it is difficult to distinguish sexy from porn categories. It is more complex to distinguish all categories in the dataset



(a) Dataset 200



(b) Dataset 4000



(c) Dataset 15000

Fig. 5. 2D TSNE projection of nsfw_model features on the 3 tested datasets.

TABLE IV
ACCURACY ON FOR EACH TESTED TOOL ON ALL DATASETS.

*NOTE THAT DEEPPORNDTECTION HAS BEEN TRAINED ON THE DATASET 15000, THE OBTAINED PERFORMANCE IS BIASED.

Name	Considered tools in this study				
	<i>nsfw_model</i>	<i>NudeNet</i>	<i>NuDetective</i>	<i>skinDetection</i>	<i>DeepPornDetection</i>
Dataset 200	83%	83.5%	64.5%	59%	61%
Dataset 4000	85%	90.4%	65.4%	53.7%	55%
Dataset 15000	72.1%	62%	57.5%	76%	95.6%*

TABLE V
AVERAGE COMPUTATION TIME EXPRESSED IN SECOND.

<i>nsfw_model</i>	<i>NudeNet</i>	<i>NuDetective</i>	<i>SkinDetection</i>	<i>DeepPornDetection</i>
0.28	0.14	0.02	0.06	0.16

4000. Figure 5 (c) related to the dataset 15000 shows how challenging is this set of images composed mainly of sexy and porn ones. The frontier is quite clear.

C. Is it always easy to predict the correct class?

We tried to analyze qualitatively the obtained results. In some cases, the classification of an image in the safe/unsafe class is not easy even by a human. Unusual positions of the body, such in yoga or gymnastics, tends to mislead the classifier into qualifying a rather safe image as unsafe. We have to deal cases such as situations depicted more suggestive than explicit (see Figure 4).

All of them contain nudity, but any private parts (breast, genitalia,...) are hidden by an arm or a leg, another one with a body covered by tattoos and a Rio carnival dancer. Table VI presents the output of the three most efficient tools. Image a) in Figure 4 is strongly considered as unsafe by NudeNet and nsfw_model, the woman is naked but the image does not reveal any intimate part. The second image b) is highly considered as safe by NudeNet and SkinDetection but the image content could be shocking for a child. The last image c) is considered as safe by all tools as probably expected by a male but not necessary by a female. It clearly shows that the frontier between these two classes safe/unsafe remains tricky.

D. Computation time

Table V shows the average computation time for each tested tool. We can see that the computation time is very low for all of them (maximum equals to 300ms).

V. CONCLUSION AND PERSPECTIVES

The proposed comparative study on challenging datasets put into abviousness the performance of nsfw_model and NudeNet. There is also room of improvement as efficiency is far to be perfect. Images at the frontier of the safe/unsafe classes are difficult to classify.

Is there a unique frontier? based on gender, religion, culture, this frontier could be different. It could be interesting to develop a fine solution where the decision threshold could

be dependent of external factors. It would concern a very interesting perspective of this work.

REFERENCES

- [1] LegiFrance, "Law no. 2020-936 of july 30, 2020 to protect victims of domestic violence." [Online]. Available: <https://www.legifrance.gouv.fr/jorf/id/JORFTEXT000042176652>
- [2] Google, "Safesearch." [Online]. Available: <https://cloud.google.com/vision/docs/detecting-safe-search>
- [3] Microsoft, "Asure analyze image." [Online]. Available: <https://learn.microsoft.com/en-us/azure/cognitive-services/computer-vision/concept-detecting-adult-content>
- [4] Amazon, "Rekognition." [Online]. Available: <https://docs.aws.amazon.com/rekognition/latest/dg/procedure-moderate-images.html>
- [5] J. A. M. Basilio, G. A. Torres, G. S. Pérez, L. K. T. Medina, H. M. Pérez Meana, and E. E. Hernandez, "Explicit content image detection," *Signal & Image Processing: An International Journal (SIPIJ) Vol.*, vol. 1, 2010.
- [6] J. A. M. Basilio, G. A. Torres, G. S. Pérez, L. K. T. Medina, and H. M. Pérez Meana, "Explicit image detection using YCbCr space color model as skin detection," in *Proceedings of the 2011 American Conference on Applied Mathematics and the 5th WSEAS International Conference on Computer Engineering and Applications*, ser. AMERICAN-MATH'11/CEA'11. Stevens Point, Wisconsin, USA: World Scientific and Engineering Academy and Society (WSEAS), 2011, p. 123–128.
- [7] I. G. P. S. Wijaya, I. Widiartha, K. Uchimura, and G. Koutaki, "Pornographic image rejection using eigenporn of simplified lda of skin rois images," in *2015 International Conference on Quality in Research (QIR)*, 2015, pp. 77–80.
- [8] G. Laborde, "Deep nn for nsfw detection." [Online]. Available: https://github.com/GantMan/nsfw_model
- [9] Nudetective forensic tool. [Online]. Available: <http://www.eleuterio.com/nudetective.html>
- [10] P. Bedapudi. (2019) Nudenet: Neural nets for nudity classification, detection and selective censoring. [Online]. Available: <https://github.com/notAI-tech/NudeNet/>
- [11] A. Lykesas. Deep learning porn detection. [Online]. Available: <https://github.com/alexookah/Deep-Learning-Porn-Detection>
- [12] Skindetection. [Online]. Available: <https://github.com/CHEREF-Mehdi/SkinDetection>
- [13] S. Kolkur, D. Kalbande, P. Shimpi, C. Bapat, and J. Jatakia, "Human skin detection using RGB, HSV and YCbCr color models," in *Proceedings of the International Conference on Communication and Signal Processing 2016 (ICCASP 2016)*. Atlantis Press, 2017. [Online]. Available: <https://doi.org/10.2991%2Ficcas-16.2017.51>
- [14] X. Wang, C. Hu, and S. Yao, "An adult image recognizing algorithm based on naked body detection," *2009 ISECS International Colloquium on Computing, Communication, Control, and Management*, 2009. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/5267781>



Fig. 6. Examples of suggestive images at the frontier of explicit content.

TABLE VI
OUTPUTS OF THE 3 MOST EFFICIENT TESTED TOOLS FOR EXPLICIT CONTENT DETECTION FOR THE THREE IMAGES IN FIGURE 6

Image	<i>nsfw_model</i>	<i>NudeNet</i>	<i>skinDetection</i>
a)	unsafe (95%)	unsafe (99%)	safe (59%)
b)	unsafe (60%)	safe (96%)	safe (100%)
c)	safe (66%)	safe (78%)	safe (90%)

- [15] A. Q. Bhatti, M. Umer, S. H. Adil, M. Ebrahim, D. Nawaz, and F. Ahmed, "Explicit content detection system: An approach towards a safe and ethical environment," *Applied Computational Intelligence and Soft Computing*, 2018. [Online]. Available: <https://www.hindawi.com/journals/acisc/2018/1463546/>
- [16] pixabay. [Online]. Available: <https://pixabay.com/fr/>
- [17] nsfw_data_source_urls. [Online]. Available: https://github.com/EBazaro/nsfw_data_source_urls/tree/master/raw_data
- [18] Mpii human pose database. [Online]. Available: <http://human-pose.mpi-inf.mpg.de>
- [19] S. Shao, Z. Zhao, B. Li, T. Xiao, G. Yu, X. Zhang, and J. Sun, "Crowdhuman: A benchmark for detecting human in a crowd," *arXiv preprint arXiv:1805.00123*, 2018.
- [20] Utkface. [Online]. Available: <https://susanqq.github.io/UTKFace/>
- [21] P. Bedapudi. (2019, Mar.) Nudenet: An ensemble of neural nets for nudity detection and censoring. [Online]. Available: <https://praneethbedapudi.medium.com/nudenet-an-ensemble-of-neural-nets-for-nudity-detection-and-censoring-d9f3da721e3>
- [22] M. de Castro Polastro and P. M. da Silva Eleuterio, "Nudetective: A forensic tool to help combat child pornography through automatic nudity detection," in *2010 Workshops on Database and Expert Systems Applications*, 2010, pp. 349–353.
- [23] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne." *Journal of machine learning research*, vol. 9, no. 11, 2008.