

A Generic Approach to detect the forgeries in Color images using ELM method

J. Angel Sajani¹, Mr. A. Muthumari M.E².

¹M.E Final year, Dept. of Computer Science and Engg, University College of Engineering, Nagercoil.

²Assistant Professor, University College of Engineering, Nagercoil, Tamilnadu.

¹angelsajani@gmail.com

²muthu_ru@yahoo.com

Abstract –

The main objective of this project is to detect forgeries. We use a forgery detection method that exploits subtle inconsistencies in the color of the illumination of the image. The approach is machine based. We incorporate information from physics and statistical based illuminant estimators on image region for similar material. To handle the illumination estimation the modified histogram based illumination estimation technique is used. In this project the human face is detected based on the skin color segmentation techniques. After the face extraction, the features are also extracted. Then these features are combined to create the paired wise features. Then these combined features are classified by using the ELM (Extreme Learning Machine) method. ELM is used to increase the accuracy rate. It yields detection rates of 90% on a new benchmark dataset consisting of 200 images, and 83% on 50 images that were collected from the Internet.

Index Terms—Color constancy, illuminant color, image forensics, machine learning, spliced image detection, texture and edge descriptors.

I. INTRODUCTION

EVERY day, millions of digital documents are produced by a variety of devices and distributed by newspapers, magazines, websites and television. In all these information channels, images are a powerful tool for communication. Unfortunately, it is not difficult to use computer graphics and image processing techniques to manipulate images. Image composition (or splicing) is one of the most common image manipulation operations. When assessing the authenticity of an image, forensic investigators use all available sources of tampering evidence. Among other telltale signs, illumination inconsistencies are potentially effective for splicing detection: from the viewpoint of a manipulator, proper adjustment of the illumination conditions is hard to achieve when creating a composite image. In this work, we make an important step towards minimizing user interaction for an illuminant-based tampering decision-making. We propose a new semiautomatic method that is also significantly more reliable than earlier approaches. Quantitative evaluation shows that the proposed method achieves a detection rate of 86%, while existing illumination-based work is slightly better than guessing. We exploit the fact that local illuminant estimates are most discriminative when comparing objects of the same (or similar) material. Thus, we focus on the automated comparison of human skin color, and

more specifically faces, to classify the illumination on a pair of faces as either consistent or inconsistent. User interaction is limited to marking bounding boxes around the faces in an image under investigation. In the simplest case, this reduces to specifying two corners of a bounding box.

When assessing the authenticity of an image, forensic investigators use all available sources of tampering evidence. Among other telltale signs, illumination inconsistencies are potentially effective for splicing detection: from the viewpoint of a manipulator, proper adjustment of the illumination conditions is hard to achieve when creating a composite image [1].

In this spirit, Riess and Angelopoulou [2] proposed to analyze illuminant color estimates from local image regions. Unfortunately, the interpretation of their resulting so-called *illuminant maps* is left to human experts. As it turns out, this decision is, in practice, often challenging. Moreover, relying on visual assessment can be misleading, as the human visual system is quite inept at judging illumination environments in pictures [3], [4]. Thus, it is preferable to transfer the tampering decision to an objective algorithm.

In this work, we make an important step towards minimizing user interaction for an illuminant-based tampering decision making. We propose a new semiautomatic method that is also significantly more reliable than earlier approaches. Quantitative evaluation shows that the proposed method achieves a detection rate of 86%, while existing illumination-based work is slightly better than guessing. We exploit the fact that local illuminant estimates are most discriminative when comparing objects of the same (or similar) material. Thus, we focus on the automated comparison of human skin, and more specifically faces, to classify the illumination on a pair of faces as either consistent or inconsistent. User interaction is limited to marking bounding boxes around the faces in an image under investigation. In the simplest case, this reduces to specifying two corners (upper left and lower right) of a bounding box.

In summary, the main contributions of this work are:

- Interpretation of the illumination distribution as object texture for feature computation.
- A novel edge-based characterization method for illuminant maps which explores edge attributes related to the illumination process.
- The creation of a benchmark dataset comprised of 100 skillfully created forgeries and 100 original photographs

II. METHODOLOGY

To overcome the drawback of the existing system the proposed system is used. Illumination-based methods for forgery detection are either geometry -based or color-based. Geometry -based methods focus at detecting inconsistencies in light source positions between specific objects in the scene [5]–[11]. Color-based methods search for inconsistencies in the interactions between object color and light color [2], [12], [13].

A. Modules

- Dense Local Illuminant Estimation (IE)
- Face Extraction
- Computation of Illuminant Features
- Paired Face Features
- Classification

Dense Local Illuminant Estimation (IE)

The input image is segmented into homogeneous regions. Per illuminant estimator, a new image is created where each region is colored with the extracted illuminant color. This resulting in-intermediate representation is called illuminant map (IM). we

briefly examine the illuminant maps generated by the method of Riess and Angelopoulou [2]. To handle the illumination estimation the modified histogram based illumination estimation technique is used. Its local illuminant color estimate yields a so-called *illuminant map*. A human expert can then investigate the input image and the illuminant map to detect in-consistencies.

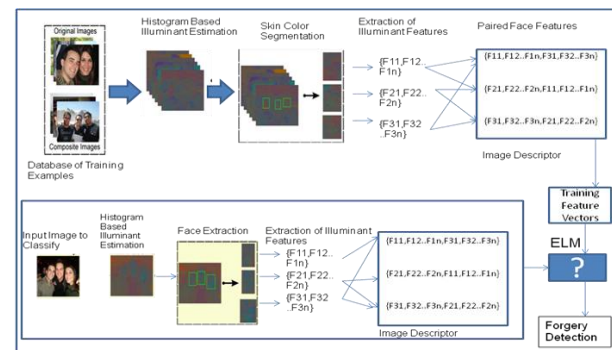


Fig. 1 summarizes these steps. In the remainder of this section, we present the details of these components.

In this module perform three operations are

- Get the Input Image
- Apply super pixel Segmentation
- Find the Inverse Intensity chromaticity Estimates.

Face Extraction

This is the only step that may require human interaction. An operator sets a bounding box around each face (e.g., by clicking on two corners of the bounding box) in the image that should be investigated. Alternatively, an automated face detector can be employed. We then crop every bounding box out of each illuminant map, so that only the illuminant estimates of the face regions remain. We require bounding boxes around all faces in an image that should be part of the investigation. For obtaining the bounding boxes, we could in principle use an automated algorithm, e.g., the one by Schwartz *et al.* [8]. However, we prefer a human operator for this task for two main reasons: a) this minimizes false detections or missed faces; b) scene context is important when judging the lighting situation.

For instance, consider an image where all persons of interest are illuminated by flashlight. The illuminants are expected to agree with one another. Conversely, assume that a person in the foreground is illuminated by flash-light, and a person in the background is illuminated by ambient light. Then, a difference in the color of the illuminants is expected. Such differences are hard to distinguish in a fully auto-mated manner,

but can be easily excluded in manual annotation. In this project the human face is detected based on the skin color segmentation techniques. This provides better result in less time.

Detection using skin color cue is fast and robust and can minimize the processing time. In addition, human skin color has its own feature color and can easily be distinguished from other objects. Therefore, in this application, skin color segmentation approaches are used as the detection instrument. The first thing to consider is the type of color space that is used and how to model it. Skin color segmentation can be defined as the process of discrimination between skin and non-skin pixels. However, there are some difficulties in robustly detecting the skin color. The ambient of the light and shadows can affect the appearance of the skin-tone color. Moreover, different cameras produce different color values even from the same person and moving object can cause blurring of colors. Finally, people have varied skin color-tones individually such as Asians skin gives big different with Caucasians skin type. We illustrate this setup in Fig. 2. The faces in Fig. 2(a) can be assumed to be exposed to the same illuminant. As Fig. 2(b) shows, the corresponding gray world illuminant map for these two faces also has similar values.

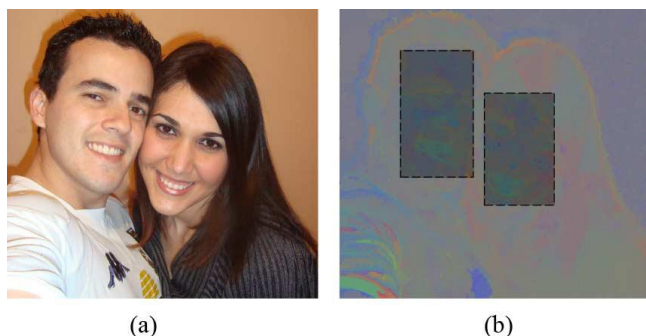


Fig. 2. Original image and its gray world map. Highlighted regions in the gray world map show a similar appearance. (a) Original. (b) Gray world with high-lighted similar parts.

Computation of Illuminant Features

For all face regions, texture-based and gradient-based features are computed on the IM values. Each one of them encodes complementary information for classification. One autocorrelation is computed using a specific fixed orientation, scale, and shift. Computing the mean and standard deviation of all such pixel values yields two feature dimensions. Repeating this computation for varying orientations, scales and shifts yields a 128-dimensional feature vector. As a final step, this vector is normalized by subtracting its mean value, and dividing it by its standard deviation.

Differing illuminant estimates in neighboring segments can lead to discontinuities in the illuminant map. Dissimilar illuminant estimates can occur for a number of reasons: changing geometry, changing material, noise, retouching or changes in the incident light. Thus, one can interpret an illuminant estimate as a low-level descriptor of the underlying image statistics. We observed that the edges, e.g., computed by a Canny edge detector, detect in several cases a combination of the segment borders and isophotes (i.e., areas of similar incident light in the image). When an image is spliced, the statistics of these edges is likely to differ from original images.

To characterize such edge discontinuities, we propose a new feature descriptor called *HOGedge*. It is based on the well-known HOG-descriptor, and computes visual dictionaries of gradient intensities in edge points. The full algorithm is described in the remainder of this section. Fig. 3 shows an algorithmic overview of the method. We first extract approximately equally distributed candidate points on the edges of illuminant maps. At these points, HOG descriptors are computed. These descriptors are summarized in a visual words dictionary. Each of these steps is presented in greater detail in the next subsections.

Extraction of Edge Points: Given a face region from an illuminant map, we first extract edge points using the Canny

edge detector [9]. This yields a large number of spatially close edge points. To reduce the number of points, we filter the Canny output using the following rule: starting from a seed point, we eliminate all other edge pixels in a region of interest (ROI) centered around the seed point. The edge points that are closest to the ROI (but outside of it) are chosen as seed points for the next iteration. By iterating this process over the entire image, we reduce the number of points but still ensure that every face has a comparable density of points. Fig. 4 depicts an example of the resulting points.

Point Description: We compute Histograms of Oriented

Gradients (HOG) to describe the distribution of the selected edge points. HOG is based on normalized local histograms of image gradient orientations in a dense grid. The HOG descriptor is constructed around each of the edge points. The neighborhood of such an edge point is called a cell. Each cell provides a local 1-D histogram of quantized gradient directions using all cell pixels.

To construct the feature vector, the histograms of all cells within a spatially larger region are combined and contrast-normalized. We

use the HOG output as a feature vector for the subsequent steps.

Dense Local Illuminant Estimation (IE)

The input image is segmented into homogeneous regions. Per illuminant estimator, a new image is created where each region is colored with the extracted illuminant color. This resulting in-intermediate representation is called illuminant map (IM). we briefly examine the illuminant maps generated by the method of Riess and Angelopoulou [2]. To handle the illumination estimation the modified histogram based illumination estimation technique is used. Its local illuminant color estimate yields a so-called *illuminant map*. A human expert can then investigate the input image and the illuminant map to detect in-consistencies.

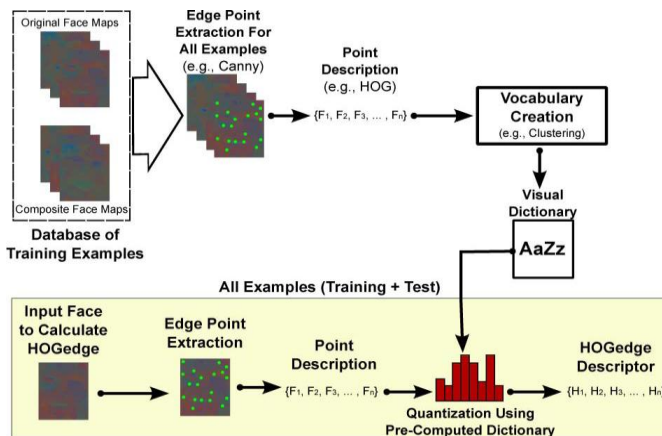


Fig. 3. Overview of the proposed HOGedge algorithm.

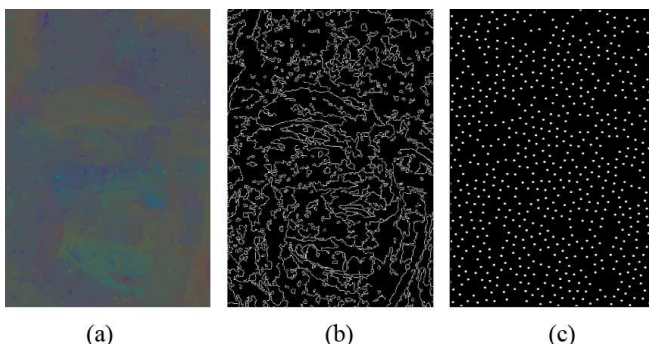


Fig. 4. (a) Gray world IM for the left face in Fig. 6(a). (b) Result of the Canny edge detector when applied on this IM. (c) Final edge points after filtering using a square region. (a) IM derived from gray world. (b) Canny Edges. (c) Filtered Points.

Visual Vocabulary: The number of extracted HOG vectors varies depending on the size and structure of the face under examination. We use visual dictionaries to obtain feature vectors of fixed length. Visual dictionaries constitute a robust

representation, where each face is treated as a set of region descriptors. The spatial location of each region is discarded. To construct our visual dictionary, we subdivide the training data into feature vectors from original and doctored images. Each group is clustered in clusters using the means algorithm. Then, a visual dictionary with visual words is constructed, where each word is represented by a cluster center. Algorithm 1 shows the pseudo code for the dictionary creation.

Quantization Using the Precomputed Visual Dictionary:

For evaluation, the HOG feature vectors are mapped to the visual dictionary. Each feature vector in an image is represented by the closest word in the dictionary (with respect to the Euclidean distance). A histogram of word counts represents the distribution of HOG feature vectors in a face. Algorithm 2 shows the pseudo code for the application of the visual dictionary on IMs.

Algorithm 1 HOGedge—Visual dictionary creation

Require: V_{TR} (training database examples) n (the number of visual words per class)

Ensure: V_D (visual dictionary containing $2n$ visual words)

$V_D \leftarrow \emptyset;$

$V_{NF} \leftarrow \emptyset;$

$V_{DF} \leftarrow \emptyset;$

for each face $IM\ i \in V_{TR}$ **do**

$V_{EP} \leftarrow$ edge points extracted from i ;

for each point $j \in V_{EP}$ **do**

$FV \leftarrow$ apply HOG in image i at position j ;

if i is a doctored face **then**

$V_{DF} \leftarrow \{V_{DF} \cup FV\};$

else

$V_{NF} \leftarrow \{V_{NF} \cup FV\};$

end if

end for

end for

Cluster V_{DF} using n centers;

Cluster V_{NF} using n centers; $V_D \leftarrow \{\text{centers of } V_{DF} \cup \text{centers of } V_{NF}\};$

return V_D ;

Face Pair

Our goal is to assess whether a pair of faces in an image is consistently illuminated. For an image with n_f faces, we construct n_f joint feature vectors, consisting of all possible pairs of faces.

The SASI and HOGedge descriptors capture two different properties of the face regions. From a signal processing point of view, both descriptors are *signatures* with different behavior. we computed the mean value and standard deviation per feature

dimension. This experiment empirically demonstrates two points. Firstly, SASI and HOGedge, in combination with the IIC-based and gray world illuminant maps create features that discriminate well between original and tampered images, in at least some dimensions. Secondly, the dimensions, where these features have distinct value, vary between the four combinations of the feature vectors. We exploit this property during classification by fusing the output of the classification on both feature sets, as described in the next section.

Classification

We use a machine learning approach to automatically classify the feature vectors. We consider an image as a forgery if at least one pair of faces in the image is classified as inconsistently illuminated.

Algorithm 2 HOGedge—Face characterization

Require: V_D (visual dictionary precomputed with $2n$ visual words) IM (illuminant map from a face)
Ensure: HFV (HOGedge feature vector)
 $HFV \leftarrow 2n$ -dimensional vector, initialized to all zeros;
 $V_{FV} \leftarrow \emptyset$;
 $V_{EP} \leftarrow$ edge points extracted from IM ;
for each point $i \in V_{EP}$ **do**
 $FV \leftarrow$ apply HOG in image IM at position j ;
 $V_{FV} \leftarrow \{V_{FV} \cup FV\}$;
end for
for each feature vector $i \in V_{FV}$ **do**
 $lower_distance \leftarrow +\infty$;
 $position \leftarrow -1$;
 for each visual word $j \in V_D$ **do**
 $distance \leftarrow$ Euclidean distance between i and j ;
 if $distance < lower_distance$ **then**
 $lower_distance \leftarrow distance$;
 $position \leftarrow$ position of j in V_D ;
 end if
 end for
 $HFV[position] \leftarrow HFV[position] + 1$;
end for
return HFV ;

We classify the illumination for each pair of faces in an image as either consistent or inconsistent. Assuming all selected faces are illuminated by the same light source, we tag an image as manipulated if one pair is classified as inconsistent. Individual feature vectors, i.e., SASI or HOG edge features on either gray world or IIC-based illuminant maps, are classified using a ELM (Extreme Learning Machine) method.

The information provided by the SASI features is complementary to the information from the HOGedge features. Thus, we use a machine learning-based fusion technique for improving the detection performance. We classify each

combination of illuminant map and feature type independently (i.e., SASI-Gray-World, SASI-IIC, HOGedge-Gray-World and HOGedge-IIC) using ELM classifier to obtain the distance between the image's feature vectors and the classifier decision boundary. ELM is used to increase the accuracy rate.

ELM possesses unique features to deal with regression and (multi-class) classification tasks. Consequently, ELM offers significant advantages such as fast learning speed, ease of implementation, and minimal human intervention. ELM has good potential as a viable alternative technique for large-scale computing and artificial intelligence.

III. EXPERIMENTAL RESULTS

To validate our approach, we performed six rounds of experiments using two different databases of images involving people. We show results using classical ROC curves where *sensitivity* represents the number of composite images correctly classified and *specificity* represents the number of original images (non-manipulated) correctly classified.

A. Evaluation Data

To quantitatively evaluate the proposed algorithm, and to compare it to related work, we considered two datasets. One consists of images that we captured ourselves, while the second one contains images collected from the internet. Additionally, we validated the quality of the forgeries using a human study on the first dataset. Human performance can be seen as a baseline for our experiments.

1) *DSO-I*: This is our first dataset and it was created by ourselves. It is composed of 200 indoor and outdoor images with an image resolution of 2,048X1,536. Out of this set of images, 100 are original, i.e., have no adjustments whatsoever, and 100 are forged. The forgeries were created by adding one or more individuals in a source image that already contained one or more persons. When necessary, we complemented an image splicing operation with post processing operations (such as color and brightness adjustments) in order to increase photorealism.

2) *DSI-I*: This is our second dataset and it is composed of 50 images (25 original and 25 doctored) downloaded from different websites in the Internet with different resolutions. Fig. 5 depicts some example images from our databases.

B. Human Performance in Spliced Image Detection

To demonstrate the quality of DSO-1 and the difficulty in dis-criminating original and tampered images, we performed an ex-periment where we asked humans to mark images as tampered or original. Note that on Mechanical Turk categorization experiments, each batch is evaluated only by experienced users which generally leads to a higher confidence in the outcome of the task. In our experiment, we setup five identical categorization experiments, where each one of them is called a batch. Within a batch, all DSO-1 images have been evaluated. For each image, two users were asked to tag the image as original or manipulated. Each image was assessed by ten different users, each user expended on average 47 seconds to tag an image. The final accuracy, averaged over all experiments, was 64.6%. However, for spliced images, the users achieved only an average accuracy of 38.3%, while human accuracy on the original images was 90.9%. According to the Landis and Koch [15] scale, suggests a slight degree of agreement between users, which further supports our conjecture about the difficulty of forgery detection in DSO-1 images.



Fig. 5. Original (left) and spliced images (right) from both databases.
(a) DSO-1 Original image. (b) DSO-1 Spliced image. (c) DSI-1 Original image. (d) DSI-1 Spliced image.

We compared the five variants SASI- IIC, SASI-Gray-World, HOGedge-IIC, HOGedge-Gray-World and Metafusion. Compound names, such as SASI-IIC, indicate the data source (in this case: IIC-based illuminant maps) and the subsequent feature ex-traction method (in this case: SASI). The single components are configured as follows:

- **IIC:** IIC-based illuminant maps are computed as described in [2].
- **Gray-World:** Gray world illuminant maps are computed by setting $n=1$, $p=1$, and $\sigma=3$ in

(2).

- **SASI:** The SASI descriptor is calculated over the Y channel from the $YCbCr$ color space. All remaining parameters are chosen as presented in .
- **HOGedge:** Edge detection is performed on the Y channel of the $YCbCr$ color space, with a Canny low threshold of 0 and a high threshold of 10. The square region for edge point filtering was set to 32×32 pixels. Furthermore, we used 8-pixel cells without normalization in HOG. If applied on IIC-based illuminant maps, we computed 100 visual words for both the original and the tampered images (i.e., the dictionary consisted of 200 visual words). On gray world illuminant maps, the size of the visual word dictionary was set to 75 for each class, leading to a dictionary of 150 visual words.
- **Metafusion:** We implemented a late fusion as explained in Section IV-F. As input, it uses SASI-IIC, SASI-Gray-World, and HOGedge-IIC. We excluded HOGedge-Gray-World from the input methods, as its weaker performance leads to a slightly worse combined classification rate.

Fig. 6 depicts a ROC curve of the performance of each method using the corner clicking face localization. The area under the curve (AUC) is computed to obtain a single numerical measure for each result.

From the evaluated variants, Metafusion performs best, re-sulting in an AUC of 86.3%. In particular for high specificity (i.e., few false alarms), the method has a much higher sensi-tivity compared to the other variants. Thus, when the detection threshold is set to a high specificity, and a photograph is classified as composite, Metafusion provides to an expert high confidence that the image is indeed manipulated.

Note also that Metafusion clearly outperforms human as-sessment in the baseline Mechanical Turk. Part of this improvement comes from the fact that Metafusion achieves, on spliced images alone, an average accuracy of 67%, while human performance was only 38.3%.

The second best variant is SASI-Gray-World, with an AUC of 84.0%. In particular for a specifi city below 80.0%, the sensitivity is comparable to Metafusion. SASI-IIC achieved an AUC of 79.4%, followed by HOGedge-IIC with an AUC of 69.9% and HOGedge- Gray-World with an AUC of 64.7%. The weak performance of HOGedge- Gray-World comes from the fact that illuminant color

estimates from the gray world algorithm vary more smoothly than IIC-based estimates. Thus, the differences in the illuminant map gradients (as extracted by the HOGedge descriptor) are generally smaller.

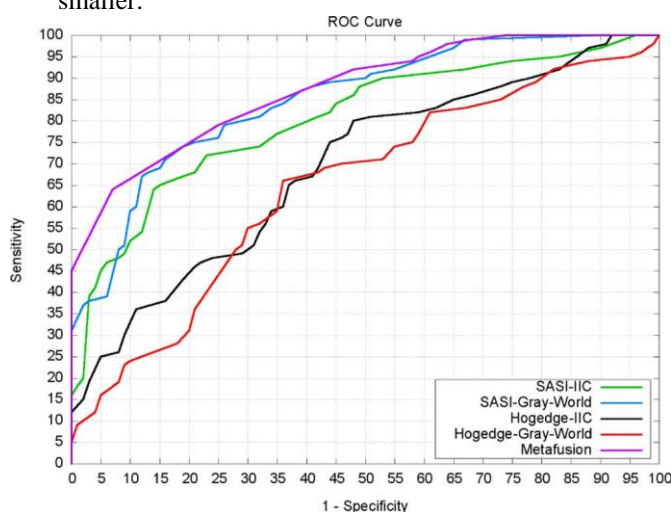


Fig. 6. Comparison of different variants of the algorithm using semiautomatic (corner clicking) annotated faces.

V. CONCLUSION AND FUTURE WORK

In this project, a new method for detecting forged images of people using the illuminant color is proposed. The estimation of the illuminant color using a statistical gray edge method and a physics-based method which exploits the inverse intensity-chromaticity color space. These illuminant maps as texture maps. It also extract information on the distribution of edges on these maps. In order to describe the edge information, we propose a new algorithm based on edge-points and the HOG ROC curve provided by cross-database experiment. descriptor, called HOGedge. It combine these complementary cues (texture- and edge-based) using machine learning late fusion. These results are encouraging, yielding an AUC of over 86% correct classification. Good results are also achieved over internet images and under cross-database training/testing. Although the proposed method is custom-tailored to detect splicing on images containing faces, there is no principal hindrance in applying it to other, problem-specific materials in the scene. The proposed method requires only a minimum amount of human interaction and provides a crisp statement on the authenticity of the image. Additionally, it is a significant advancement in the exploitation of illuminant color as a forensic cue. Prior color-based work either assumes complex user interaction or imposes very limiting assumptions.

An incorporation of this method is subject of future work. Reasonably effective skin detection methods have been presented in the computer vision

literature in the past years. Incorporating such techniques can further expand the applicability of this method. Such an improvement could be employed, for instance, in detecting pornography compositions which, according to forensic practitioners, have become increasingly common nowadays.

REFERENCES

- [1] A. Rocha, W. Scheirer, T. E. Boult, and S. Goldenstein, "Vision of the unseen: Current trends and challenges in digital image and video forensics," *ACM Comput. Surveys*, vol. 43, pp. 1–42, 2011.
- [2] C. Riess and E. Angelopoulou, "Scene illumination as an indicator of image manipulation," *Inf. Hiding*, vol. 6387, pp. 66–80, 2010.
- [3] H. Farid and M. J. Bravo, "Image forensic analyses that elude the human visual system," in *Proc. Symp. Electron. Imaging (SPIE)*, 2010, pp. 1–10.
- [4] Y. Ostrovsky, P. Cavanagh, and P. Sinha, "Perceiving illumination inconsistencies in scenes," *Perception*, vol. 34, no. 11, pp. 1301–1314, 2005.
- [5] H. Farid, "A 3-D lighting and shadow analysis of the JFK Zapruder film (Frame 317)," *Dartmouth College, Tech. Rep. TR2010-677*, 2010.
- [6] S. Gholap and P. K. Bora, "Illuminant colour based image forensics," in *Proc. IEEE Region 10 Conf.*, 2008, pp. 1–5.
- [7] X. Wu and Z. Fang, "Image splicing detection using illuminant color inconsistency," in *Proc. IEEE Int. Conf. Multimedia Inform. Networking and Security*, Nov. 2011, pp. 600–603.
- [8] W. R. Schwartz, A. Kembhavi, D. Harwood, and L. S. Davis "Human detection using partial least squares analysis," in *Proc. IEEE Int. Conf. Comput. Vision (ICCV)*, 2009, pp. 24–31.
- [9] J. Canny, "A computational approach to edge detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 8, no. 6, pp. 679–698, Jun. 1986.
- [10] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Proc. IEEE Conf. Comput. Vision and Pattern Recognition*, 2005, pp. 886–893.
- [11] G. Csurka, C. R. Dance, L. Fan, J. Willamowski, and C. Bray, "Visual categorization with bags of keypoints," in *Proc. Workshop on Statistical Learning in Comput. Vision*, 2004, pp. 1–8.

- [12] J. Winn, A. Criminisi, and T. Minka, "Object categorization by learned universal visual dictionary," in Proc. IEEE Int. Conf. Comput. Vision (ICCV), 2005, pp. 1800–1807.
- [13] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc, 2006.
- [14] O. Ludwig, D. Delgado, V. Goncalves, and U. Nunes, "Trainable classifier-fusion schemes: An application to pedestrian detection," in Proc. IEEE Int. Conf. Intell. Transportation Syst., 2009, pp. 1–6.
- [15] J. R. Landis and G. G. Koch, "The measurement of observer agreement for categorical data," *Biometrics*, vol. 33, no. 1, pp. 159–174, 1977.