

Camera-based Image Forgery Localization using Convolutional Neural Networks

Davide Cozzolino, Luisa Verdoliva

Department of Electrical Engineering and Information Technology

University Federico II of Naples, Naples, Italy

name.surname@unina.it

Abstract—Camera fingerprints are precious tools for a number of image forensics tasks. A well-known example is the photo response non-uniformity (PRNU) noise pattern, a powerful *device* fingerprint. Here, to address the image forgery localization problem, we rely on noiseprint, a recently proposed CNN-based camera *model* fingerprint. The CNN is trained to minimize the distance between same-model patches, and maximize the distance otherwise. As a result, the noiseprint accounts for model-related artifacts just like the PRNU accounts for device-related non-uniformities. However, unlike the PRNU, it is only mildly affected by residuals of high-level scene content. The experiments show that the proposed noiseprint-based forgery localization method improves over the PRNU-based reference.

Index Terms—Image forensics, PRNU, convolutional neural networks.

I. INTRODUCTION

With the widespread diffusion of powerful media editing tools, falsifying images and videos has become easier and easier in the last few years. Manipulated visual content, often used to support fake news, represents a growing menace in many fields of life, from politics, to journalism, to the judiciary. In response to this threat, a large number of methods have been recently proposed for image forgery detection and localization [1].

Supervised methods which exploit the presence of camera artifacts in the image under analysis [2], [3], especially those relying on the PRNU pattern [4], have shown great potential for many forensic tasks. The PRNU, caused by sensor defects arising in the manufacturing process, links univocally each photo to the device that acquired it, and can be therefore regarded as a sort of device fingerprint. PRNU-based methods have shown a very good performance for both source identification and image forgery detection [4]–[9]. The camera PRNU pattern must be accurately estimated in advance, which requires the availability of the camera itself or of a certain number of photos (typically, 100-200) taken from it. At testing time, this reference pattern is compared with the single-image

PRNU estimate extracted from the image under analysis. For camera identification, the comparison takes place on the whole image. Instead, for forgery detection and localization, a sliding-window correlation-based procedure is used. In the presence of a forgery, the reference PRNU is missing, and a low correlation is observed.

Ideally, this procedure allows one to detect and localize accurately all attacks to the image under test. In practice, the PRNU traces found in any individual image represent a very weak signal in strong noise (scene residuals, camera artifacts) which makes the whole process quite unreliable. To improve the signal-to-noise ratio, the high-level scene content (seen as noise in this context) is removed by means of suitable denoising filters, obtaining a noise residual which represents the desired single-image PRNU estimate. However, due to imperfections of denoising filters, some scene contents leak in the noise residual, inducing false alarms in dark, uniform, or very textured areas. Besides using state-of-the-art denoising filters [10], [11], a number of strategies have been proposed to face this problem. In [5] a predictor is used to identify potentially troubling regions and adapt the statistical test locally, while in [12] the interference of scene details is reduced by selective attenuation of wavelet coefficients.

A further source of noise, besides scene residuals, is represented by the so-called non-unique artifacts [13], traces of in-camera processes which are specific of a camera model but not of the individual device. Such artifacts are caused, for example, by JPEG compression or CFA interpolation, and are characterized by spatially periodic patterns. Again, various strategies have been proposed to remove them [5], [14].

However, one should remember that estimating the PRNU is not a goal in itself, but a means towards the completion of forensic tasks. Do model-related artifacts really hamper such tasks? Or else can they be exploited to improve performance? Under an information-theoretic point of view, the answer is obvious: all available information should be taken into account. Based on this line of thought, in [15] we proposed a new approach to identify the camera model traces and exploit them for multimedia forensics. Unlike in the prior literature, which focuses on compact features, we extract a PRNU-like camera model fingerprint, called *noiseprint*, which displays non-unique artifacts in the form of an image-size pattern (see Fig.1). The scene content is effectively removed by means of a siamese residual-based convolutional neural network,

This material is based on research sponsored by the Air Force Research Laboratory and the Defense Advanced Research Projects Agency under agreement number FA8750-16-2-0204. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the Air Force Research Laboratory and the Defense Advanced Research Projects Agency or the U.S. Government.

obtaining eventually a strong pattern affected by weak noise, which allows the easy (even visual) detection and localization of anomalies due to local image manipulations.

In [15] we developed a totally *blind* noiseprint-based forgery localization technique, working only on the image residual (single-image noiseprint estimate) with no side information. Here, we consider the same problem in a supervised setting. Given a suitable training set of images, a reliable estimate of the noiseprint is built, and used in a PRNU-like fashion to discover anomalies in the residual of the image under analysis. In fact, when a region of the image is tampered with, its noiseprint is perturbed, that is, replaced with a new one (splicing), strongly modified (rotation, resizing), or even deleted (inpainting), which allows one to detect and localize the attack. It is worth reminding that the noiseprint contains only weak traces of the PRNU itself, hence it is not device-specific. Therefore, it should not be regarded as a substitute of the PRNU. On the other hand, the noiseprint has a much higher signal-to-noise ratio than the PRNU pattern, a property which guarantees a better performance and allows its applications to more challenging situations.

In the reminder of the paper we give some more details on noiseprint, then describe the proposed noiseprint-based forgery localization procedure, finally we discuss experimental results and draw conclusions.

II. NOISEPRINT

In [15] we set the goal of extracting a noise-like image, called noiseprint, which works as a camera *model* fingerprint, much like the PRNU pattern can be regarded as a *device* fingerprint. In principle, this could be obtained by simply keeping the noise residual without removing non-unique artifacts. In practice, however, this residual does not possess the desired properties, because it suffers from the very same problems as the PRNU, caused by significant leakages from the high-level signal. Therefore, we designed a new system, based on deep learning, specifically dedicated to our goal.

To gain insight into our design choices, let us start from the ultimate goal. Our system must accept in input a generic image and provide in output a residual image, the noiseprint, having the same size as the input and containing only camera-model specific features with their natural spatial distribution. This task reminds of residual-based denoisers, which extract the additive Gaussian noise component of a given input image,

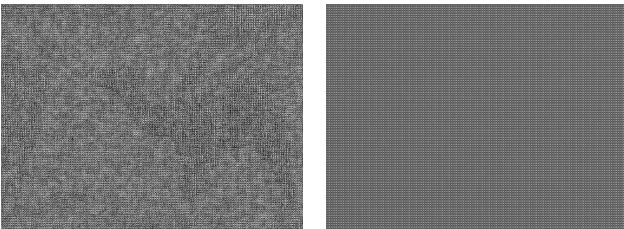


Fig. 1. Single-image (left) and 200-images (right) noiseprint estimates of the same camera. In the first case, some scene contents (camel) leaked in the estimate. In the second case, only periodic model-related artifacts are visible.

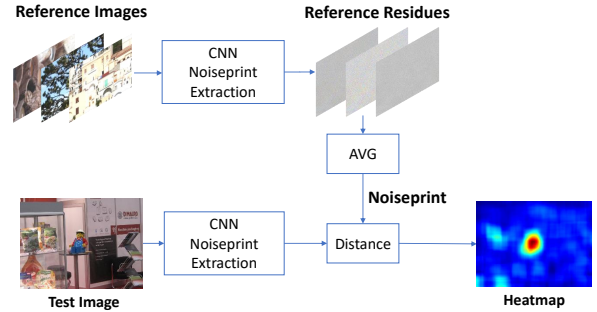


Fig. 2. Proposed supervised forgery localization scheme.

thus removing the high-level content. Such denoisers are trained by back-propagating the error between the the noise residual output by the CNN and the *true* noise pattern. In our problem, however, no true reference is available for training, as the ideal noiseprint of a camera model is unknown.

Nonetheless, we know that images acquired by the same camera model have the same noiseprint. We exploit this information in the training phase by using *pairs* of images, extracting their residuals in parallel, and computing their distance (mean squared error). When the pair comes from the same model, the error is back-propagated through the network to reduce the distance between residuals. On the contrary, when images come from different models, back-propagation is used to increase the residual distance. Therefore, it is like considering identical twin networks, working in parallel. Each one uses the output of the other twin, together with the same/different label, to adapt its weights towards convergence.

The idea of Siamese networks is not new in deep learning. Typically, each twin net extracts a low-dimensional vector (embedding) which summarizes the input keeping the relevant information for the specific task. The main goal is to perform a low-complexity comparison of the twin outputs. In our case, however, the aim is not reducing complexity, while it is important to preserve spatial information. Hence, our net must preserve image size and spatial relationships.

Turning to the practical implementation, we initialize the network (ideally, the Siamese nets) as the CNN denoiser proposed in [16] for AWGN (additive white Gaussian noise) image denoising. Removing the scene content, in fact, is a reasonable starting point for our goal. To limit training complexity, and also to ensure flexibility w.r.t. input size, we work on small image patches, rather than whole images. Hence, we feed the network with a large number of paired patches, with label +1 when they refer to the same camera model *and* the same spatial location, and label -1 otherwise. The constraint on the spatial location is necessary to preserve the precious spatial information. In fact, even a single-pixel shift impacts heavily on the local statistics of the residual. The distance between the CNN outputs of the two paired patches is then computed. In the loss function, pairs with negative label (different model and/or different position) are penalized. Therefore, at convergence the net should extract the same

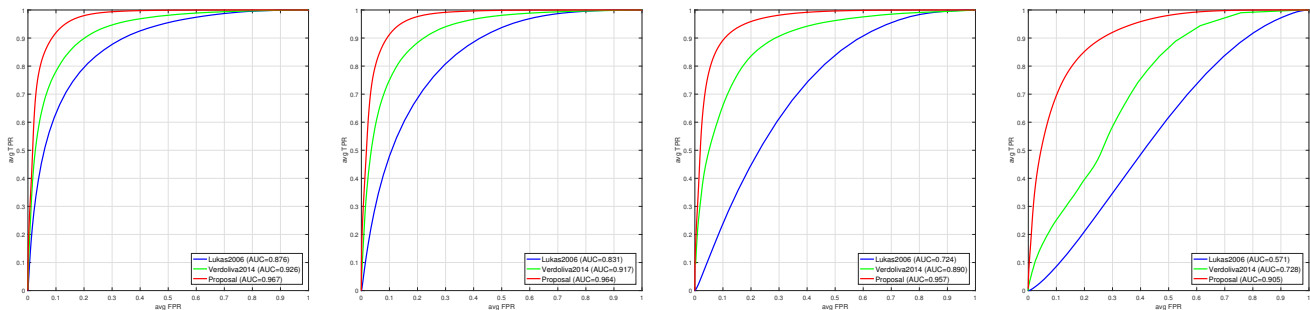


Fig. 3. ROCs at pixel level by varying the number of images for estimation (200,50,10,1).

residual for positive samples, and different residuals otherwise.

It is worth underlining that the network is trained once and for all, it can process any camera, both inside and outside the training set, and is not specific of a given experiment, or task. Therefore, once the training is over, the noiseprint is deterministically related to the original image.

We refer the interested reader to [15] for details on the training phase. Here, we mention only two key solutions adopted to improve the training efficiency, the method proposed in [17] to obtain $O(n^2)$ rather than $O(n)$ samples for each n -patch minibatch, and the distance based logistic (DBL) loss proposed in [18].

III. USING NOISEPRINT IN A SUPERVISED SETTING

In [15] we performed forgery localization with no side information, by looking for anomalies in the image residual. Here, we consider a supervised setting, assuming a reference noiseprint is available. The localization procedure is outlined in Fig.2 and follows the classic pipeline used with PRNU-based methods. We rely on a set of pristine images taken by the same camera of the image under analysis. Their residuals are averaged to obtain a clean reference, namely, a reliable estimate of the camera noiseprint where high-level scene leakages, as well as traces of the PRNU, are mostly removed. Fig.1 shows a single-image estimate, in which traces of the high-level content can be easily spotted (see the camel image in Fig.4), and a 200-image estimate which is virtually noise-free. This reference is then compared in sliding window modality with the residual of the image under test, using the Euclidean distance as a measure of similarity. As usual in these cases, the window size impacts on the trade-off between resolution and reliability. Here, a 64×64 window is used.

The pixel-wise distance field can then be shown as a heatmap, which can be provided to the end user for visual inspection or subject to a suitable post-processing to extract a binary decision map. Large distances (red in the heatmaps) suggest that the original noiseprint has been corrupted, that is, deleted (inpainting), replaced with the noiseprint of another camera (splicing), or even of the same camera but after some geometrical distortion (resizing, rotation, even simple displacement). Therefore, even a rigid copy-move can be discovered, unless the displacement in both the vertical and horizontal directions is a multiple of the noiseprint fundamental period. It

is worth underlining, however, that this approach makes sense only if the test image is aligned with the pristine reference images. If the test image is geometrically distorted or subject to heavy compression, also the references should undergo the same processing chain to deliver a correct reference noiseprint.

IV. EXPERIMENTAL RESULTS

The network used to extract all noiseprints is trained on a large variety of models. To this end, we created a large dataset, including both cameras and smartphones, using various publicly available datasets, plus some other private cameras. In detail, we used 44 cameras from the Dresden dataset [19], 32 from Socrates dataset [20], 32 from VISION [21], 17 from our private dataset, totaling 125 individual cameras from 70 different models and 19 brands. In the experiments, this dataset is split in training-set and validation-set comprising 100 and 25 cameras, respectively. All images are originally in JPEG format with a quality factor in the range [96-99]. The network is initialized with the weights of the denoising network of [16]. During training, each minibatch contains 200 patches of 48×48 pixels extracted from 100 different images of 25 different cameras. In each batch, there are 50 sets, each one formed by 4 patches with same camera and position. Training is performed using ADAM optimizer; the hyperparameters (learning rate, number of iterations and weight of regularization term) are chosen using the validation set.

For image forgery localization, we compare results with the PRNU-based method proposed in [4] (Lukas2006) and with a feature-based approach relying on camera model artifacts [2] (Verdoliva2014). We use the same 300-image dataset with splicings already used in [6]. Images come from 4 camera models (CanonEOS 450D, Nikon D200, CanonIXUS_95IS, NikonCoolpix_S5100) none of which used for training the network.

In Fig.3(left) we show pixel-level localization results in terms of receiver operating curves (ROC) when 200 images are used to estimate the reference noiseprint, PRNU, or the statistic of [2]. The proposed method provides a large gain over the Lukas2006, and also over Verdoliva2014. Synthetic results reported in Tab.I, in terms of area under curve (AUC) and F-measure, fully confirm this analysis. For F-measure we used both the best threshold over all the dataset and (more

TABLE I
PIXEL-LEVEL LOCALIZATION PERFORMANCE.

	Lukas2006	Verdoliva2014	Proposed
AUC	0.876	0.926	0.967
F1	0.499	0.580	0.724
F1-oracle	0.572	0.707	0.850

favourable for performance) the best threshold for each image (F1-oracle). In all cases, the performance gain of the proposed method is clear. The very same conclusions can be drawn by visual inspection of the examples shown in Fig.4, where very different types of manipulation have been considered (splicing, rigid copy-move, inpainting). Note that all images have the native camera JPEG quality, and no further compression is carried out. The noiseprint-based method exhibits always a very good performance, highlighting clearly the manipulations without false alarms.

As we remarked several times, the noiseprint appears to be less noisy than the PRNU pattern. Therefore, unlike for PRNU-based methods, we may expect the performance to depend only weakly on the number of reference images. To investigate this point, we repeated the previous experiment using only 50, 10, and 1 reference images, respectively. The corresponding ROCs, also shown in Fig.3, confirm noiseprint robustness. While the performance of the other methods are largely impaired when less than 50 images are used, the proposed method works reasonably well even with a single reference image, and better than Lukas2006 with 200 images.

In the last experiment we consider a more challenging situation, where the images have a different format, and hence are not well aligned with our dataset. We take some examples from the dataset used by Korus et al. [9]. The images are taken from 4 different camera models (Canon 60D, Sony A57, Nikon D7000 and Nikon D90) and attacked with different forms of manipulations. They are in raw format, hence no JPEG artifacts are present, a significant mis-alignment with respect to the trained CNN. We use the PRNU provided with this dataset, estimated over 200 natural images for Nikon and Canon, and over 20 flat images for Sony. For noiseprint and for [2] we used only the 53 available pristine images, taking care to avoid using the same background in training and test. As expected, Lukas2006 works better than in the previous case, since images are not compressed, while the performance of the proposed method is impaired due to the mis-alignment. Even in this critical situation, however, the noiseprint-based method keeps providing valuable hints for forgery localization.

V. CONCLUSIONS

We performed image forgery localization by a sliding-window comparison of image residual and camera fingerprint. Unlike in the recent literature, we do not use the PRNU fingerprint to this end, but a new camera *model* fingerprint, called noiseprint, extracted by means of a suitably trained Siamese CNN. Experiments show that the noiseprint-based

method outperforms largely the PRNU-based baseline, and keeps providing very good results even when the fingerprint is estimated on a very small number of images. Overall, noiseprint appears to have a great potential for multimedia forensic analyses, both in supervised and blind settings.

REFERENCES

- [1] P. Korus, "Digital image integrity a survey of protection and verification techniques," *Digital Signal Processing*, vol. 71, pp. 1–26, 2017.
- [2] L. Verdoliva, D. Cozzolino, and G. Poggi, "A feature-based approach for image tampering detection and localization," in *IEEE International Workshop on Information Forensics and Security*, 2014, pp. 149–154.
- [3] L. Bondi, S. Lameri, D. Güera, P. Bestagini, E. Delp, and S. Tubaro, "Tampering Detection and Localization through Clustering of Camera-Based CNN Features," in *IEEE Computer Vision and Pattern Recognition Workshops*, 2017.
- [4] Lukáš, J. Fridrich, and M. Goljan, "Detecting digital image forgeries using sensor pattern noise," in *Proc. SPIE*, vol. 6072-0Y, 2006, pp. 362–372.
- [5] M. Chen, J. Fridrich, M. Goljan, and J. Lukas, "Determining image origin and integrity using sensor noise," *IEEE Transactions on Information Forensics and Security*, vol. 3, no. 4, pp. 74–90, 2008.
- [6] G. Chierchia, G. Poggi, C. Sansone, and L. Verdoliva, "A Bayesian-MRF approach for PRNU-based image forgery detection," *IEEE Transactions on Information Forensics and Security*, vol. 9, no. 4, pp. 554–567, 2014.
- [7] P. Korus and J. Huang, "Evaluation of random field models in multi-modal unsupervised tampering localization," in *IEEE Workshop on Information Forensics and Security*, december 2016, pp. 1–6.
- [8] S. Chakraborty and M. Kirchner, "PRNU-based forgery detection with discriminative random fields," in *International Symposium on Electronic Imaging: Media Watermarking, Security, and Forensics*, February 2017.
- [9] P. Korus and J. Huang, "Multi-scale analysis strategies in PRNU-based tampering localization," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 4, pp. 809–824, april 2017.
- [10] G. Chierchia, S. Parrilli, G. Poggi, C. Sansone, and L. Verdoliva, "On the influence of denoising in PRNU based forgery detection," in *ACM Workshop on Multimedia in Forensics, Security and Intelligence*, october 2010, pp. 117–122.
- [11] M. Al-Ani and F. Khelifi, "On the SPN estimation in image forensics: a systematic empirical evaluation," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 5, pp. 1067–1081, May 2017.
- [12] C.-T. Li, "Source camera identification using enhanced sensor pattern noise," *IEEE Transactions on Information Forensics and Security*, vol. 5, no. 2, pp. 280–287, 2010.
- [13] J. Fridrich, "Sensor defects in digital image forensic," in *Digital Image Forensics*, H. Sencar and N. Memon, Eds. Springer-Verlag, 2012, pp. 179–218.
- [14] C. Li and Y. Li, "Color-Decoupled Photo Response Non-Uniformity for Digital Image Forensics," *IEEE Transactions on Information Forensics and Security*, vol. 22, no. 3, pp. 260–271, february 2012.
- [15] D. Cozzolino and L. Verdoliva, "Noiseprint: a CNN-based camera model fingerprint," *submitted*, 2018.
- [16] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a Gaussian Denoiser: Residual Learning of Deep CNN for Image Denoising," *IEEE Transactions on Image Processing*, vol. 26, no. 7, pp. 3142–3155, July 2017.
- [17] H. Song, Y. Xiang, S. Jegelka, and S. Savarese, "Deep metric learning via lifted structured feature embedding," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [18] N. Vo and J. Hays, "Localizing and orienting street views using overhead imagery," in *European Conference on Computer Vision*, 2016.
- [19] T. Gloe and R. Böhme, "The dresden image database for benchmarking digital image forensics," in *ACM Symposium on applied Computing*, 2010, pp. 1584–1590.
- [20] C. Galdi, F. Hartung, and J.-L. Dugelay, "Videos versus still images: asymmetric sensor pattern noise comparison on mobile phones," in *IS&T Electronic Imaging: Media Watermarking, Security and Forensics*, 2017.
- [21] D. Shullani, M. Fontani, M. Iuliani, O. A. Shaya, and A. Piva, "Vision: a video and image dataset for source identification," *EURASIP Journal on Information Security*, pp. 1–16, 2017.

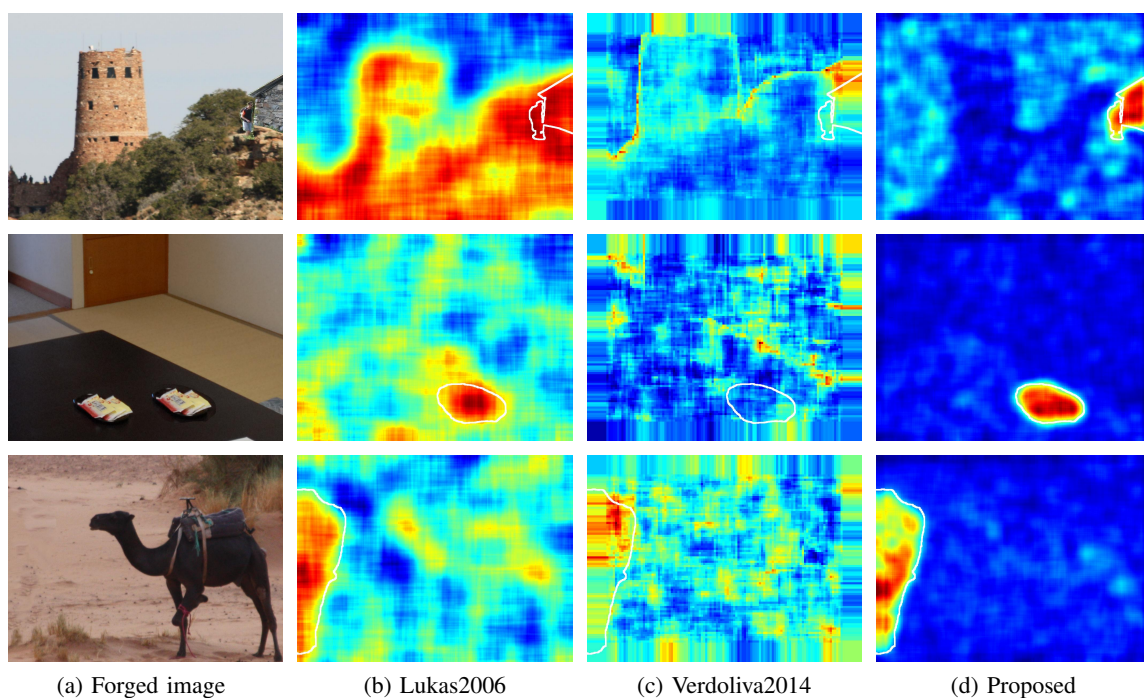


Fig. 4. Forgery localization results for some selected examples. From top to bottom: splicing, rigid copy-move, inpainting.

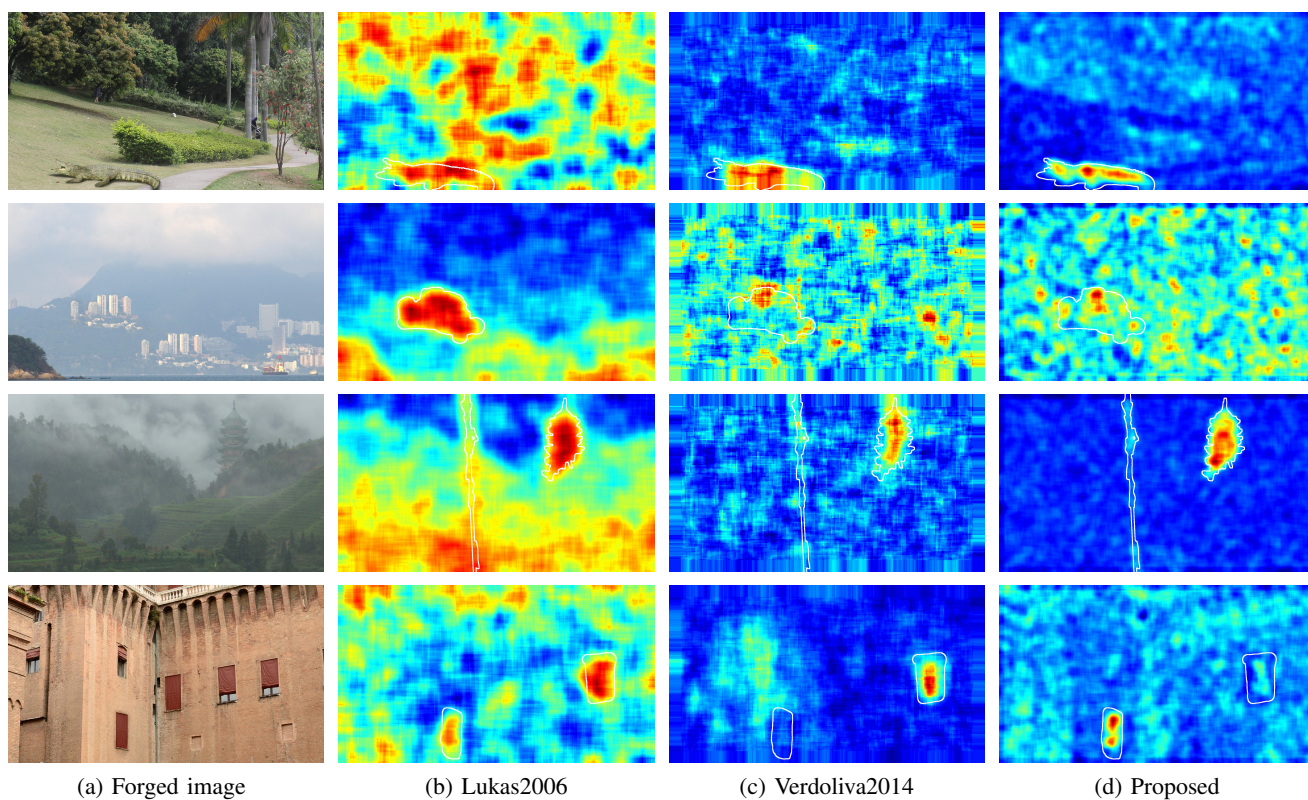


Fig. 5. Forgery localization results for some selected examples using the dataset proposed in [9]. Even if the images are in raw format and estimation is carried out on less data with respect to PRNU, results for the proposed approach are promising.