# IMAGE FORGERY LOCALIZATION BASED ON MULTI-SCALE CONVOLUTIONAL NEURAL NETWORKS

*Yaqi Liu, Qingxiao Guan, Xianfeng Zhao, and Yun Cao*

1. State Key Laboratory of Information Security, Institute of Information Engineering,
Chinese Academy of Sciences, Beijing 100093, China
2. School of Cyber Security, University of Chinese Academy of Sciences, Beijing 100093, China

## ABSTRACT

In this paper, we propose to utilize Convolutional Networks (CNNs) and the segmentation-based multi-scale analysis to locate tampered areas in digital images. First, to deal with color input sliding windows of different scales, a unified CNN architecture is designed. Then, we elaborately design the training procedures of CNNs on sampled training patches. With a set of robust multi-scale tampering detectors based on CNNs, complementary tampering possibility maps can be generated. Last but not least, a segmentation-based method is proposed to fuse the maps and generate the final decision map. By exploiting the benefits of both the small-scale and large-scale analyses, the segmentation-based multi-scale analysis can lead to a performance leap in forgery localization of CNNs. Numerous experiments are conducted to demonstrate the effectiveness and efficiency of our method.

***Index Terms***— Image forensics, forgery localization, multi-scale analysis, Convolutional Neural Networks.

## 1. INTRODUCTION

Image forgery localization is one of the most challenging tasks in digital image forensics [1]. Different from forgery detection which simply discriminates whether a given image is pristine or fake, image forgery localization attempts to detect the accurate tampered areas [2]. Since forgery localization needs to conduct pixel-level analyses, it is more difficult than the conventional forgery detection task.

Different clues are investigated to locate the tampered areas, e.g., the photo-response nonuniformity noise (PRNU) [3], the artifacts of color filter array [4], the traces left by JPEG coding [5], the near-duplicate image analysis [6], and copy-move forgery detection [7], etc. The tampering operations inevitably distort some inherent relationships among the adjacent pixels, features motivated by steganalysis [8] are frequently adopted to localize tampered areas [9, 1]. In 2013,
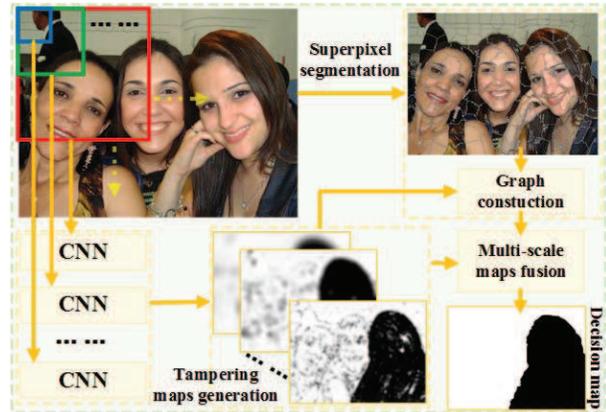
**Fig. 1**. The framework of MSCNNs. Note that the sliding windows (blue, green, red squares) and superpixels do not indicate their real sizes.

IEEE Information Forensics and Security Technical Committee (IFS-TC) established the First IFS-TC Image Forensics Challenge [10]. In the second phase, a complicated and practical situation for evaluating the performance of forgery localization was set up. The winner [11] and successors [6, 1] combined different clues to achieve high scores. As far as we know, the former best F1-score using a single clue from statistical features was achieved in [1] which was based on color rich models [12] and the ensemble classifier [13] (SCRM+LDA).

We focus on forgery localization utilizing statistical features extracted by Convolutional Neural Networks (CNNs) [14]. Booming in computer vision tasks, CNNs are also applied in image forensics. In [15], CNNs are applied in median filtering image forensics. In [16], a novel constrained convolutional layer is utilized to suppress the content of the image, and CNNs are adopted to detect multiple manipulations. In [17], a CNN with SRM kernels [8] for the first layer initialization is adopted for forgery detection. In [18], they show that residual-based descriptors can be regarded as a simple constrained CNN which can conduct forgery detection and localization. Numerous meaningful works have been

done to improve the performance of image forensics by adopting CNNs, they try to construct different CNN architectures. While in computer vision tasks, typical CNNs, e.g. AlexNet [14], VGG [19], ResNet [20], etc., are directly adopted for different purposes [21], and they mainly focus on the preprocessing and postprocessing. This kind of adoptions accelerate the development of many computer vision tasks. Thus, instead of designing a totally novel CNN, we adopt and modulate the state-of-the-art CNNs [22] to construct our framework for image forgery localization. More powerful CNNs can also be adopted in the proposed framework in the future.

In this paper, an image forgery localization method based on Multi-Scale Convolutional Neural Networks (MSCNNs) is proposed, as shown in Figure 1. In our method, sliding windows of different scales are put into a set of CNNs to generate real-valued tampering possibility maps. Then, based on the graph constructed on superpixels [23], we can generate the final decision map by fusing those possibility maps. The contributions are two-fold: First, we propose to utilize multi-scale CNNs to detect forged regions. A unified CNN architecture is formulated for color patches, and multi-scale CNNs are treated as a set of "weak" classifiers to fully exploit the benefits of both the small-scale and large-scale analyses. Second, based on the fusion method in [2], the segmentation-based fusion method is proposed to efficiently process images of different sizes. Maps fusion based on conditional random fields is conducted on the superpixel-level graph, and two strategies for superpixel-level tampering possibility maps generation are proposed and compared.

On the IFS-TC dataset, MSCNNs can achieve the best performance among the forgery localization methods which merely utilize one kind of clue for splicing detection. To the best of our knowledge, only three methods, i.e., the winner [11] and successors [6, 1], can achieve higher scores than MSCNNs, but they all combine multiple different clues, e.g. statistical features, copy-move clues etc. The proposed MSCNNs only utilizes statistical features extracted by CNNs and can be further improved by adopting other clues. Besides, to demonstrate the robustness of the proposed framework, we also conduct experiments on another dataset, i.e. Realistic Tampering Dataset (RTD) [2, 24].

The rest of the paper is structured as follows. In Section 2, we elaborate the proposed method. In Section 3, experiments are conducted. In Section 4, we draw conclusions.

## 2. METHOD

### 2.1. CNNs architecture

Our motivation is that we want to replace the SCRM+LDA [1] with the end-to-end CNNs to estimate the tampering probability of a given patch. Adopting the sliding window manner, we can give the tampering possibility map of the investigated image. The CNNs proposed in [22] achieve the state-of-the-art

performance for steganalysis on gray-scale images. Considering the close relationship between image forensics and steganalysis, we adopt this kind of CNNs as the basic architecture in our work. In the first layer of their CNNs, a single high pass filter (we call it the base filter) is utilized to suppress the image content. In our work, to deal with color patches, two kinds of base filters are tested:

(1) Fixed SRM kernels: the base filters are fixed, and set as the SRM kernels [8]. In [17], 30 SRM kernels are adopted for the initialization of the first layer of their CNNs. We adopt all the SRM kernels as fixed base filters, and leave the task of validating their effectiveness to the backend network. Referring to [17], the 30 SRM kernels are formulated as $5 \times 5$ matrixes $\{\mathbf{F}_1, \cdots \mathbf{F}_{30}\}$ with zero-valued unused elements. The inputs are three-channel color patches, so we need $30 \times 3$ filters to generate 30 feature maps. For the $j$th feature map ($j \in \{1, 2, \cdots 30\}$), the corresponding filters are set as $\{\mathbf{F}_1^j, \mathbf{F}_2^j, \mathbf{F}_3^j\} = \{\mathbf{F}_{3k-2}, \mathbf{F}_{3k-1}, \mathbf{F}_{3k}\}$, where $k = ((j - 1) \bmod 10) + 1$.

(2) Constrained filters: in [16], a kind of constrained filter was proposed for manipulation detection. Here, we adopt it for forgery localization. The constraint means that the filter weight at the center $f(0, 0) = -1$, and $\sum_{r, c \neq 0} f(r, c) = 1$, $f(r, c)$ denotes the element in the base filter $\mathbf{F}$. For fair comparisons, 90 $5 \times 5$ constrained filters are adopted.

As we adopt 90 base filters, we modulate the parameters of CNNs in [22], and the unified CNN architecture can be depicted as Figure 2. For different scales of input patches, we only need to change $P$ in the last average pooling layer, ensuring that the input of the fully-connected layer is a 256-dimensional vector. Based on the CNN depicted in Figure 2, we can train a set of CNN detectors with input patches of different scales. The detailed training procedures are introduced in Section 3.

### 2.2. Maps generation

For each input image, it is analysed by the sliding window of the scale as $s \times s$ with a stride of $st$ based on the CNN detectors described in Section 2.1. Then, we can get the tampering possibility map $\hat{\mathbf{M}}_s$ of size $h_s \times w_s$, where $h_s = \lfloor (h - s)/st \rfloor + 1$ and $w_s = \lfloor (w - s)/st \rfloor + 1$, $h$ and $w$ denote the height and width of the input image, and $\lfloor \cdot \rfloor$ denotes the floor function. The elements in $\hat{\mathbf{M}}_s$ denote the probabilities of the corresponding patches being fake. In order to get the possibility map $\mathbf{M}_s$ with the same size as the input image, the element $m_{i,j}^s$ in $\mathbf{M}_s$ is computed as:

$$m_{i,j}^s = \frac{1}{K} \sum_{k=1}^{K} \hat{m}_k^s \qquad (1)$$

where $K$ is the number of patches containing pixel $I_{i,j}$, and $\hat{m}_k^s$ denotes the corresponding value in $\hat{\mathbf{M}}_s$. Inevitably, for some pixels, $K$ is equal to 0, and the pixels always appear
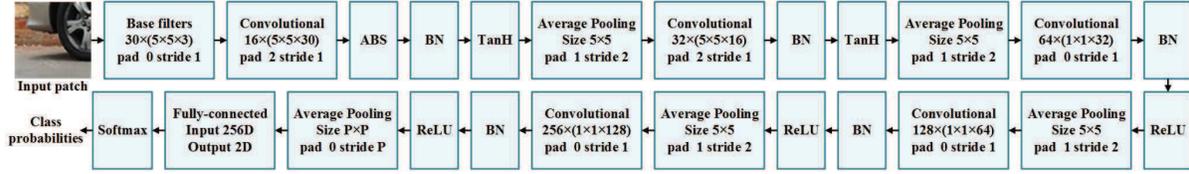
**Fig. 2**. The architecture and parameters of the unified CNNs.

around the edges of the image. We simply set the same probabilities as the nearest pixels whose $K \neq 0$. Since we have a large stride $st$, there are mosaic artifacts in the possibility map generated by formula (1). Naturally, it is expected that the map for an image tends to be smoother [1]. To smooth the possibility map, the mean filtering is applied as:

$$\bar{m}_{i,j}^{s} = \frac{1}{s \times s} \sum_{i'=-\frac{s}{2}}^{\frac{s}{2}-1} \sum_{j'=-\frac{s}{2}}^{\frac{s}{2}-1} m_{i+i',j+j'}^{s} \quad (2)$$

where $s$ is the size of corresponding patches. Thus, we can get the smoothed possibility map $\bar{\mathbf{M}}_s$ with elements as $\bar{m}_{i,j}^{s}$.

### 2.3. Maps fusion

With the analyses of multi-scale CNNs detectors, we can get a set of tampering possibility maps $\{\bar{\mathbf{M}}_s\}$ for each image, and $s$ denotes the scales of input patches. The final task is to fuse possibility maps to exploit the benefits of multi-scale analyses. In [2], the multi-scale analysis in PRNU-based tampering localization was proposed. By minimizing an energy function, possibility maps fusion is formulated as a random-field problem where decision fusion resolves to finding an optimal labeling of authentication units. The optimization problem is solved by the graph cut algorithm whose worst case running time complexity is $O(ev^2)$ [25], where $v$ is the number of nodes in the graph and $e$ is the edge number. They consider a 2nd-order neighborhood, which means that $e \approx 4v$, so the complexity of the method is $O(v^3)$. They adopt pixels as the nodes in the graph, thus the computing time of the large image is almost unacceptable. So we propose to construct graphs on superpixels, and find the optimal labels on the superpixel level.

Simple linear iterative clustering (SLIC) [23] is a commonly used efficient superpixel segmentation method, and we adopt SLIC to conduct oversegmentation on the investigated color images. The complexity of SLIC is linear, i.e. $O(v)$, and it is easy to generate superpixels by SLIC for large images. In the computer vision tasks, images are usually segmented into hundreds of superpixels. In the task of tampering possibility maps fusion, large superpixels can lead to information loss. Thus, thousands of superpixels must be generated in our task. Then, a graph on the superpixels is constructed, each superpixel is treated as a node in the graph and the adjacent superpixels are connected by an edge. The number of graph nodes is around several thousand, which is much easier to compute by the graph cut algorithm. Besides the efficiency of the superpixel-level computation, the segmentation-based method can also well adhere to the real boundaries, and avoid mislabeling of homogeneous pixels, resulting in the performance improvement.

As for the superpixel-level tampering possibility maps $\mathbf{M}_s^{sup}$ generation, two strategies are proposed and compared. The one is "mean", and the tampering possibility $m_{sup_l}^{s}$ of superpixel $l$ under scale $s$ is computed as:

$$m_{sup_l}^{s} = \frac{1}{P_l} \sum_{p=1}^{P_l} \bar{m}_p^{s} \quad (3)$$

where $P_l$ denotes the number of pixels in superpixel $l$, and $m_{sup_l}^{s} \in \mathbf{M}_s^{sup}$. $\bar{m}_p^{s}$ is the element in $\bar{\mathbf{M}}_s$. The other strategy called "maxa" is:

$$m_{sup_l}^{s} = \bar{m}_{p_0}^{s}, p_0 = \arg \max_{p=1,\cdots,P_l} (\text{abs}(\bar{m}_p^{s} - \theta)) \quad (4)$$

where $\bar{m}_p^{s} \in [0,1]$, so we set $\theta = 0.5$. With the superpixel-level graph and superpixel-level maps at hands, it is easy to fuse the maps by minimizing the energy function in [2]:

$$\frac{1}{S} \sum_{i=1}^{N} \sum_{\{s\}} E_\tau(c_i^{(s)}, t_i) + \alpha \sum_{i=1}^{N} t_i + \sum_{i=1}^{N} \sum_{j \in \Xi_i} \beta_{ij}|t_i - t_j| \quad (5)$$

where $S$ is the number of candidate possibility maps. In our segmentation-based method, $N$ is the number of elements in $\mathbf{M}_s^{sup}$, $t_i = 1$ denotes tampered units, and $c_i^{(s)}$ denotes the element of the input candidate map with analysis windows of size $s$, i.e. $c_i^{(s)} = m_{sup_l}^{s}$. The three terms can penalize differences of different possibility maps, bias the decision towards the hypotheses and encode a preference towards piecewise-constant solutions. For space limitations, the detailed definitions and discussions of the terms are not provided here, readers can kindly refer to the seminal work [2] for details. In terms of the parameters in the energy function, we adopt the default settings of the codes provided by [2].

## 3. EXPERIMENTAL EVALUATION

Experiments are conducted on two publicly available datasets. In Section 3.1, we introduce the experimental results on the image corpus provided in the IFS-TC Image Forensics Challenge (IFS-TC) [10]. In Section 3.2, experiments are conducted on Realistic Tampering Dataset (RTD) [2, 24].

## 3.1. Experiments on IFS-TC

In the IFS-TC image dataset, there are two sets of images, i.e. 450 images in the training set with corresponding human-labeled ground truths, and 700 testing images without ground truths. The scores on the testing set have to be computed by the system provided by the IFS-TC challenge. Thus, in order to test the methods locally, we randomly select 368 images for training and 75 images for testing (7 images are deserted for imperfect ground truths) from the training set of IFS-TC. For the sake of clarity, the image set of 368 images is called *sub-training set*, the image set of 75 images is called *testing set-1*, and the testing set of IFS-TC with 700 images is called *testing set-2*.

During the patches generation, we also adopt the sliding window manner. The sliding window with a fixed scale slides across the full image. We set the stride $st$ as 8 to get plenty of sampled patches. In the training set, the tampered areas are marked as the ground truths, we can sample patches based on whether they contain tampered pixels. In [1], the patches tampered with 10% to 90% are regarded as fake patches for that discriminative features mostly appear around the contours of manipulated regions, we also adopt this strategy. The rates of the tampered areas in the full images differ greatly. In some images, more than ten thousand patches can be generated, while in some images, no patch can be generated. The imbalance of patches distribution can lead to overfitting, so we set an upper threshold $T$. While more than $T$ patches are generated, we randomly select $T$ patches, and we set $T = 500$ to make sure that we sample a similar number of patches on most images. With the sliding window sampling manner, no patch can be generated for some images. For those images, we resample patches which are centered at the tampered areas. If the tampered rates of patches are satisfied, the patches are selected. After the fake patches are generated, we sample the same number of pristine patches in the same images, and the pristine patches do not have any tampered pixels. With 5 groups of sampled patches of scales as $\{32, 48, 64, 96, 128\}$, 5 independent CNNs can be trained, and the CNNs are trained on the *sub-training set*.

Our method is implemented via Caffe and Matlab. Mini-batch gradient descent is adopted for training, the momentum is 0.99 and weight decay is 0.0005. The learning rate is initialized to 0.001 and scheduled to decrease 10% for every 8000 iterations. The convolution kernels are initialized by random numbers generated from zero-mean Gaussian distribution with standard deviation of 0.01, and bias learning is disabled. The parameters in the fully-connected layer are initialized using "Xavier". Note that the input patches for the CNNs should all subtract the mean values of each channel.

We summarize localization performance as an average F1-score [1]. As shown in Table 1, the comparisons between SCRM+LDA (codes provided by [12, 13]) and different variants of CNNs are conducted. "MF" denotes mean filtering,

**Table 1**. The comparisons on the IFS-TC *testing set-1*. Time-1 denotes the training time, and Time-2 denotes the average computing time.

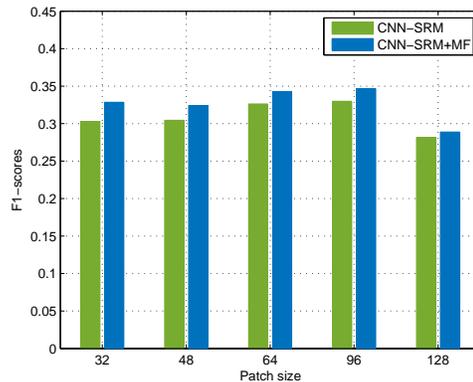| Method | Size | Stride | Time-1 (s) | Time-2 (s) | F1-score |
|---|---|---|---|---|---|
| SCRM+LDA | 64 | 16 | $3.20 \times 10^5$ | 2854.75 | 0.2847 |
| SCRM+LDA+MF | 64 | 16 | | 2855.05 | 0.3123 |
| CNN-SRM | 64 | 8 | 3376.07 | 17.11 | 0.3263 |
| CNN-SRM+MF | 64 | 8 | | 17.32 | **0.3423** |
| CNN-SRM+MF | 64 | 16 | 3376.07 | 8.38 | 0.3354 |
| CNN-C-SRM+MF | 64 | 8 | 3843.03 | 31.66 | 0.2816 |
| CNN-C-GAU+MF | 64 | 8 | 3849.09 | 31.71 | 0.2718 |



**Fig. 3**. The F1-scores of CNNs with input patches of different sizes on IFS-TC *testing set-1*.

and it can certainly improve the F1-scores based on the experimental observation. The results in Figure 3 also corroborate that, and the main reason of the improvement achieved by MF smoothing is that the map for an image tends to be smoother without mosaic artifacts caused by sliding-window operations. The training procedure of SCRM+LDA takes too much time. Although we have a powerful CPU, it takes almost 4 days. Furthermore, its average computing time on the images is also unacceptable. With the same patch size (64) and stride (16), the computing time of CNN is $1/340$ of SCRM+LDA. CNN-SRM denotes the CNN with fixed SRM base filters, CNN-C-SRM denotes constrained filters with SRM initialization and the base filters of CNN-C-GAU are constrained filters with Gaussian initialization. It can be seen that CNN-SRM can achieve higher F1-scores. Because there are many zero values in the SRM base filters, it is also more efficient than CNN-C-SRM and CNN-C-GAU.

For the good performance of CNN-SRM, we adopt this form of CNN for multi-scale analyses, and the stride is set as 8. As shown in Figure 3, CNNs with scales as 64 and 96 can achieve higher scores, and no single-scale CNN can achieve a score higher than 0.35. Nevertheless, as shown in Figure 4, the multi-scale analysis can improve the performance significantly. As an alternative, we resize the maps, and conduct
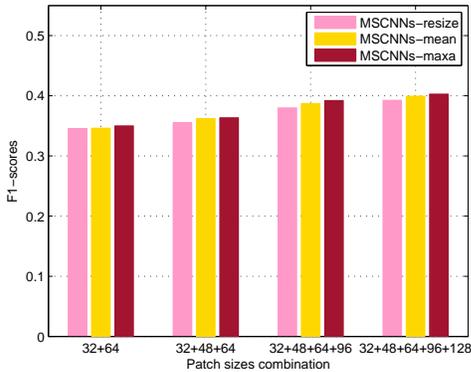
**Fig. 4**. The F1-scores of MSCNNs with different combinations on IFS-TC *testing set-1*.

maps fusion on the resized pixel-level maps directly. Let $w$ and $h$ denote the width and height of the maps, if $w > 2000$ and $h > 2000$, the map is reduced to $1/10$ of the original map; if $w < 1000$ and $h < 1000$, the map is reduced to $1/2$; otherwise, it is reduced to $1/4$. We call this kind of method as "MSCNN-resize". "mean" and "maxa" represent the two strategies for superpixel-level tampering possibility maps generation. It can be seen that MSCNNs can achieve higher scores with more scales, and MSCNNs-maxa can achieve a higher score than MSCNNs-resize and MSCNNs-mean. In MSCNNs-maxa and MSCNNs-mean, all the images are empirically segmented into $4000$ superpixels, and adaptive segmentation strategies for maps fusion need further research in the future.

Subsequently, we adopt five single-scale CNNs and the 5-scale MSCNNs to test on the *testing set-2*. As shown in Table 2, the right side presents the results of different variants of our method, and the left side presents results of the state-of-the-art methods for splicing detection. In another word, the compared methods are not designed for some particular cases, e.g. copy-move forgery detection, and can be utilized to detect any splicing forgeries. Their results are borrowed from their papers [1, 6, 11]. SCRM+LDA adopts the sliding window manner with the scale of $64$, and our CNN with $s = 64$ can achieve better performance than SCRM+LDA. Multi-scale analyses can greatly improve the performance of CNNs, and MSCNNs-maxa can achieve a similar F1-score as the winner of IFS-TC challenge [11] (0.4063 vs. 0.4072). The winner makes use of three different clues, while MSCNNs-maxa only utilizes features extracted by CNNs and can be further improved by combining other clues.

We evaluate the computing time on the *testing set-2* in which the sizes of images vary from $922 \times 691$ to $4752 \times 3168$ (most images are around $1024 \times 768$). Experiments are conducted on a machine with Intel(R) Core(TM) i7-5930K CPU @ 3.50GHz, 64GB RAM and a single GPU (TITAN X). As

**Table 2**. Results on the IFS-TC *testing set-2*.

| Method | F1-score | Variant | F1-score |
|---|---|---|---|
| S3+SVM [11] | 0.1115 | CNN-SRM32MF | 0.3436 |
| S3+LDA [1] | 0.1737 | CNN-SRM48MF | 0.3526 |
| PRNU [6] | 0.2535 | CNN-SRM64MF | 0.3570 |
| SCRM+LDA [1] | 0.3458 | CNN-SRM96MF | 0.3423 |
| | | CNN-SRM128MF | 0.3135 |
| | | MSCNNs-resize | 0.4014 |
| | | MSCNNs-mean | 0.4025 |
| | | MSCNNs-maxa | **0.4063** |

**Table 3**. Computing time on IFS-TC *testing set-2*.

| | | 32 | 48 | 64 | 96 | 128 |
|---|---|---|---|---|---|---|
| CNNs | Average time (s) | 15.47 | 15.10 | 17.74 | 19.21 | 19.56 |
| | Median time (s) | 7.96 | 7.67 | 8.89 | 9.33 | 9.20 |
| MF | Average time (s) | 0.08 | 0.13 | 0.19 | 0.37 | 0.63 |
| | Median time (s) | 0.04 | 0.07 | 0.11 | 0.20 | 0.35 |
| | multi-scales fusion: 32+48+64+96+128 | | | | | |
| Fusion | Average time (s) | | | 20.88 | | |
| | Median time (s) | | | 11.75 | | |

shown in Table 3, the computing time of 5-scales MSCNNs is around $60$ s for most images. The MF and Fusion (including SLIC) procedures are implemented on CPU which can be further accelerated by implementing on GPU.

### 3.2. Experiments on RTD

The RTD dataset contains 220 realistic forgeries created by hand and covers various challenging tampering scenarios involving both object insertion and removal. The images were captured by four different cameras: Canon 60D (C60D), Nikon D90 (ND90), Nikon D7000 (ND7000), Sony $\alpha57$ (S57). All images are $1920 \times 1080$ px RGB uint8 bitmaps stored in the TIFF format [2, 24]. Each kind of camera contains 55 images, and we randomly select 27 as the *training set*, and the left 28 images compose the *testing set*. In another words, there are 108 images in the *training set* and 112 images in the *testing set*. We adopt the same manner to sample patches on RTD, readers can refer to Section 3.1 for details.

Firstly, we adopt the models trained on the *sub-training set* of IFS-TC to test on the RTD *testing set*. The CNN is the model based on CNN-SRM and mean filtering, and the results of SCRM+LDA are also processed by mean filtering. The size of the sliding window is $64 \times 64$, and the stride is set as 16 for fair comparison. The models based on MSCNNs are the 5-scale models as above mentioned. As shown in Figure 5, the performance of all the models decline than the performance on IFS-TC. It proves that both CNN and SCRM+LDA tend to be sensitive to the training sets for that the images may be captured from different cameras and the quality of manipulations may be different. In a different dataset, MSCNNs can still achieve better performance.

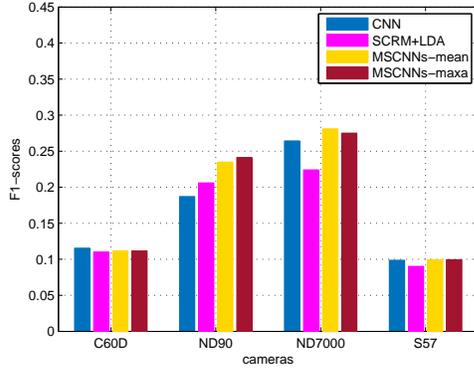Then, models trained on the RTD *training set* are com-

**Fig. 5**. The F1-scores on RTD *testing set*. All the models are trained on the *sub-training set* of IFS-TC.
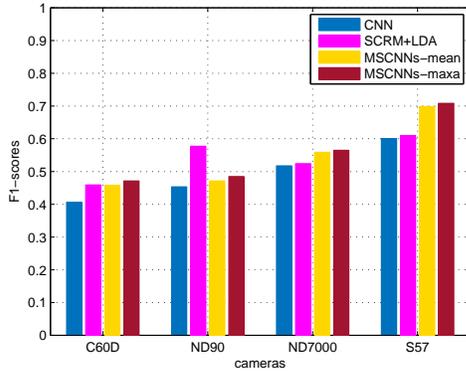


**Fig. 6**. The F1-scores on RTD *testing set*. All the models are trained on the *training set* of RTD.

pared. As shown in Figure 6, it can be seen that the performance of CNN is worse than SCRM+LDA. However, with the help of multi-scale analyses, MSCNNs can achieve better performance than SCRM+LDA except for results on ND90. Furthermore, the CNN and MSCNNs are very efficient, the average computing time of CNN is 6.58 s, and the computing time of MSCNNs is 34.62 s (5 CNNs on GPU) +30.36 s (the fusion procedure on CPU), while SCRM+LDA takes 2220.45 s per image. Thus, MSCNNs is a better alternative of SCRM+LDA in the image forgery localization tasks.

## 4. CONCLUSIONS

In this paper, a novel forgery localization method based on Multi-Scale Convolutional Neural Networks is proposed. CNNs for color patches of different scales are well designed and trained as a set of forgery detectors. Then, segmentation-based multi-scale analysis is utilized to dig out the information given by the different-scale analyses. Full experiments on the publicly available datasets demonstrate the effectiveness

and efficiency of the proposed method named MSCNNs. Although the proposed method can achieve the state-of-the-art performance, it still has a long way to go for real applications. The robustness of existing works against post compression, manipulation qualities and camera models still needs to be further studied. In the future, MSCNNs can also be improved by adopting more powerful CNNs.

## 5. REFERENCES

[1] Haodong Li, Weiqi Luo, Xiaoqing Qiu, and Jiwu Huang, "Image forgery localization via integrating tampering possibility maps," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 5, pp. 1240–1252, 2017.

[2] Paweł Korus and Jiwu Huang, "Multi-scale analysis strategies in prnu-based tampering localization," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 4, pp. 809–824, 2017.

[3] Mo Chen, Jessica Fridrich, Miroslav Goljan, and Jan Lukás, "Determining image origin and integrity using sensor noise," *IEEE Transactions on Information Forensics and Security*, vol. 3, no. 1, pp. 74–90, 2008.

[4] Pasquale Ferrara, Tiziano Bianchi, Alessia De Rosa, and Alessandro Piva, "Image forgery localization via fine-grained analysis of cfa artifacts," *IEEE Transactions on Information Forensics and Security*, vol. 7, no. 5, pp. 1566–1577, 2012.

[5] Tiziano Bianchi and Alessandro Piva, "Image forgery localization via block-grained analysis of jpeg artifacts," *IEEE Transactions on Information Forensics and Security*, vol. 7, no. 3, pp. 1003–1017, 2012.

[6] Lorenzo Gaborini, Paolo Bestagini, Simone Milani, Marco Tagliasacchi, and Stefano Tubaro, "Multi-clue image tampering localization," in *IEEE International Workshop on Information Forensics and Security (WIFS)*. IEEE, 2014, pp. 125–130.

[7] Davide Cozzolino, Giovanni Poggi, and Luisa Verdoliva, "Efficient dense-field copy–move forgery detection," *IEEE Transactions on Information Forensics and Security*, vol. 10, no. 11, pp. 2284–2297, 2015.

[8] Jessica Fridrich and Jan Kodovsky, "Rich models for steganalysis of digital images," *IEEE Transactions on Information Forensics and Security*, vol. 7, no. 3, pp. 868–882, 2012.

[9] Davide Cozzolino and Luisa Verdoliva, "Single-image splicing localization through autoencoder-based anomaly detection," in *IEEE International Workshop on Information Forensics and Security (WIFS)*. IEEE, 2016, pp. 1–6.

[10] IFS-TC, "The 1st ieee ifs-tc image forensics challenge," http://ifc.recod.ic.unicamp.br/fc.website/index.py, 2013.

[11] Davide Cozzolino, Diego Gragnaniello, and Luisa Verdoliva, "Image forgery localization through the fusion of camera-based, feature-based and pixel-based techniques," in *IEEE International Conference on Image Processing (ICIP)*. IEEE, 2014, pp. 5302–5306.

[12] Miroslav Goljan, Jessica Fridrich, and Rémi Cogranne, "Rich model for steganalysis of color images," in *IEEE International Workshop on Information Forensics and Security (WIFS)*. IEEE, 2014, pp. 185–190.

[13] Jan Kodovsky, Jessica Fridrich, and Vojtěch Holub, "Ensemble classifiers for steganalysis of digital media," *IEEE Transactions on Information Forensics and Security*, vol. 7, no. 2, pp. 432–444, 2012.

[14] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[15] Jiansheng Chen, Xiangui Kang, Ye Liu, and Z Jane Wang, "Median filtering forensics based on convolutional neural networks," *IEEE Signal Processing Letters*, vol. 22, no. 11, pp. 1849–1853, 2015.

[16] Belhassen Bayar and Matthew C Stamm, "A deep learning approach to universal image manipulation detection using a new convolutional layer," in *The 4th ACM Workshop on Information Hiding and Multimedia Security*. ACM, 2016, pp. 5–10.

[17] Yuan Rao and Jiangqun Ni, "A deep learning approach to detection of splicing and copy-move forgeries in images," in *IEEE International Workshop on Information Forensics and Security (WIFS)*. IEEE, 2016, pp. 1–6.

[18] Davide Cozzolino, Giovanni Poggi, and Luisa Verdoliva, "Recasting residual-based local descriptors as convolutional neural networks: an application to image forgery detection," in *The 5th ACM Workshop on Information Hiding and Multimedia Security*. ACM, 2017, pp. 159–164.

[19] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Identity mappings in deep residual networks," in *European Conference on Computer Vision*. Springer, 2016, pp. 630–645.

[21] Yaqi Liu, Xiaoyu Zhang, Xiaobin Zhu, Qingxiao Guan, and Xianfeng Zhao, "Listnet-based object proposals ranking," *Neurocomputing*, vol. 267, pp. 182–194, 2017.

[22] Guanshuo Xu, Han-Zhou Wu, and Yun-Qing Shi, "Structural design of convolutional neural networks for steganalysis," *IEEE Signal Processing Letters*, vol. 23, no. 5, pp. 708–712, 2016.

[23] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Süsstrunk, "Slic superpixels compared to state-of-the-art superpixel methods," *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 11, pp. 2274–2282, 2012.

[24] Paweł Korus and Jiwu Huang, "Evaluation of random field models in multi-modal unsupervised tampering localization," in *Proc. of IEEE Int. Workshop on Inf. Forensics and Security*, 2016.

[25] Yuri Boykov and Vladimir Kolmogorov, "An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision," *IEEE transactions on pattern analysis and machine intelligence*, vol. 26, no. 9, pp. 1124–1137, 2004.