# PSCC-Net: Progressive Spatio-Channel Correlation Network for Image Manipulation Detection and Localization

Xiaohong Liu[1*]    Yaojie Liu[2]    Jun Chen[1]    Xiaoming Liu[2]

[1]McMaster University    [2]Michigan State University

{liux173, chenjun}@mcmaster.ca    {liuyaoj1, liuxm}@msu.edu

## Abstract

*To defend against manipulation of image content, such as splicing, copy-move, and removal, we develop a Progressive Spatio-Channel Correlation Network (PSCC-Net) to detect and localize image manipulations. PSCC-Net processes the image in a two-path procedure: a top-down path that extracts local and global features and a bottom-up path that detects whether the input image is manipulated, and estimates its manipulation masks at 4 scales, where each mask is conditioned on the previous one. Different from the conventional encoder-decoder and no-pooling structures, PSCC-Net leverages features at different scales with dense cross-connections to produce manipulation masks in a coarse-to-fine fashion. Moreover, a Spatio-Channel Correlation Module (SCCM) captures both spatial and channel-wise correlations in the bottom-up path, which endows features with holistic cues, enabling the network to cope with a wide range of manipulation attacks. Thanks to the light-weight backbone and progressive mechanism, PSCC-Net can process $1,080P$ images at $50+$ FPS. Extensive experiments demonstrate the superiority of PSCC-Net over the state-of-the-art methods on both detection and localization.*

## 1. Introduction

<div align="center">

*Seeing is believing?*

</div>

Not anymore. Recent advances on image manipulation techniques [13, 32, 33, 37] enable easy editing of raw images, such as removing unwanted objects [31, 34, 35, 67], face swapping [32], attribute changing [50], *etc*. Although such techniques are neutral, malicious attackers may utilize them to create deceitful content to propagate false information, *e.g.*, fake news [23], insurance fraud [65], and Deepfake [12, 55]. Thus, concerns of the adverse impact on social media and even real-world systems have been raised [54, 60]. To
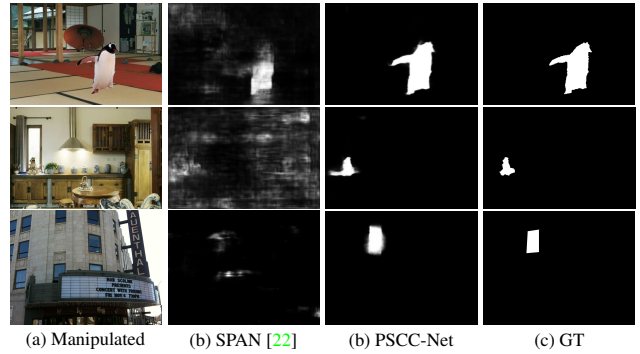


Figure 1: Examples of image manipulation localization. Three examples are splicing, copy-move, and removal manipulations respectively. With novel designs of progressive mechanism and correlation module, our method demonstrates robust and accurate estimation at different scales and types.

alleviate the concerns, it is crucial to develop reliable models to expose the manipulated images. While being used in machine and systems, the model is required to, at a minimal, distinguish manipulated images from pristine ones, where the objective is to **detect**. While being used for human's viewing, the model is further required to estimate tampered areas in forged images, where the objective is to **localize**.

Generally, image manipulation consists of the content-dependent process and content-independent process. The former includes splicing, copy-move, and removal, as shown in Fig. 1. Both splicing and copy-move are content-copying forgeries, where the splicing content is from a different donor image while the copy-move content is from the target image *per se*. Removal takes out certain objects from the target image and performs refilling via inpainting. Often, the content-dependent process follows the semantic arrangement in the target image, *e.g.*, placing a car on the road and replacing one face with another, which makes the resulting image visually "authentic" and indistinguishable from the pristine one. However, based on image/camera trace analysis [7, 11], subtle patterns can still be revealed to indicate the manipulation.

---

*This work was conducted when Xiaohong Liu was a visiting scholar at Michigan State University.

On the other hand, the content-independent process includes universal modifications such as brightness/contrast change, blurring, noising and image compression. They barely create any disinformation, but their resultant noise may undermine the analysis of image/camera traces and potentially hide the discrepancy between the manipulated and pristine areas.

To defend against manipulations, many image manipulation detection and localization (IMDL) methods have been proposed in the past. In the early stages, methods are designed to handle a single type of manipulation. In recent years, works [3, 4, 12, 22, 49, 65, 69] are proposed to build generic IMDL models for *multiple* manipulation types. However, there are still 3 major unsolved problems for IMDL:

1. **Scale variation** The forged area varies in sizes. Most prior works neglect the importance of scale variations and encounter difficulty when detecting forged areas of different sizes. Both the conventional encoder-decoder [4, 69] and no-pooling [22, 65] structures have difficulties in leveraging local and global features jointly, thus can only handle a limited scale variation.

2. **Image correlation** Manipulated regions can best be determined while comparing to pristine regions, especially for splicing attacks. A naive learning of mapping from the manipulated image to manipulation mask may lead to an overfitting to the specific attack type in training. In contrast, considering the image spatial correlation can lead to a more generalized localization solution. Yet, such correlation is mostly neglected in prior works.

3. **Detection** In principle, manipulation detection and localization are highly relevant tasks, where the detection score can be simply derived from the response of the predicted manipulation mask, *i.e.*, at least one part of the forged image has high response while no part of the pristine one does. However, most prior works assume the *existence of manipulation* in all input images. As a result, this could cause many false alarms on pristine images and make the detection unreliable.

To address the above issues, we propose a novel Progressive Spatio-Channel Correlation Network (PSCC-Net), as in Fig. 2. PSCC-Net consists of a top-down path and a bottom-up path. In the top-down path, a backbone encoder first extracts the local and global features from an input image. We adopt the network structure of [57] as our encoder, whose dense connections among different scales facilitate information exchange. In the bottom-up path, we leverage the learned features to estimate 4 manipulation masks from small scales to large ones, where each mask serves as a prior in the next-scale estimation. Thanks to such a design, the final mask is estimated in a coarse-to-fine fashion, harvesting both the local and global information. Moreover, this design enables a potential speed-up by terminating the bottom-up

mask estimation, if the intermediate mask is satisfactory. Moreover, rather than investigating the response of predicted manipulation masks, we feed the learned features into a detection head to produce the score for binary classification.

To cope with the image correlation, we propose a Spatio-Channel Correlation Module (SCCM) that grasps both spatial and channel-wise correlations at each bottom-up step. The spatial correlation aggregates the global context among local features. The channel-wise correlation computes the similarity among feature maps to enhance the representation in interest areas. Given the light-weight design of the encoder, PSCC-Net can process $1,080P$ at $50+$ FPS. Our proposed approach demonstrates a superior manipulation localization on several benchmarks. In addition, we show that the recent IMDL methods encounter difficulty in distinguishing manipulated images from pristine ones. By explicitly introducing a detection head, our method achieves the state of the art (SOTA) on manipulation detection.

We summarize the contributions of this work as follows:
⋄ We propose a new PSCC-Net that performs favorably on manipulation detection and enables progressive improvement of manipulation localization in a coarse-to-fine fashion;
⋄ We design a novel SCCM module to capture the spatial and channel-wise correlations for better generalization. SCCM avoids the use of massive annotated data to pre-train our feature extractor;
⋄ We achieve the SOTA results for both image manipulation detection and localization.

## 2. Related Work

**Image manipulation detection** Image manipulation detection aims to distinguish manipulated images from pristine ones via image-level binary classification. There are two major approach for this detection: the implicit manner [23, 63] and the explicit manner [24]. The former obtains the detection score by the statistics (*e.g.*, average [23] or maximum [63] value) of the predicted manipulation mask, and the latter explicitly outputs the score from a dedicated classification module. Recent works [22, 65] focus on pixel-level manipulation localization but neglect the importance of image-level detection. Instead, this work leverages both manipulated and pristine images in training and jointly considers detection and localization of image manipulation.

**Image manipulation localization** Early works propose to localize the manipulation of one specific type, *e.g.,* splicing [2, 5, 10, 11, 23, 29, 41, 62], copy-move [9, 24, 59, 63, 64], removal [71], and the content-preserved process [4, 27]. Although most methods perform well on detecting that specific forgery type, they fall short in handling real-world cases, where usually the forgery type is unknown in advance and various types of forgery might be utilized in manipulation. In the related problem of face anti-spoofing, researchers also study how to localize the facial pixels covered with various
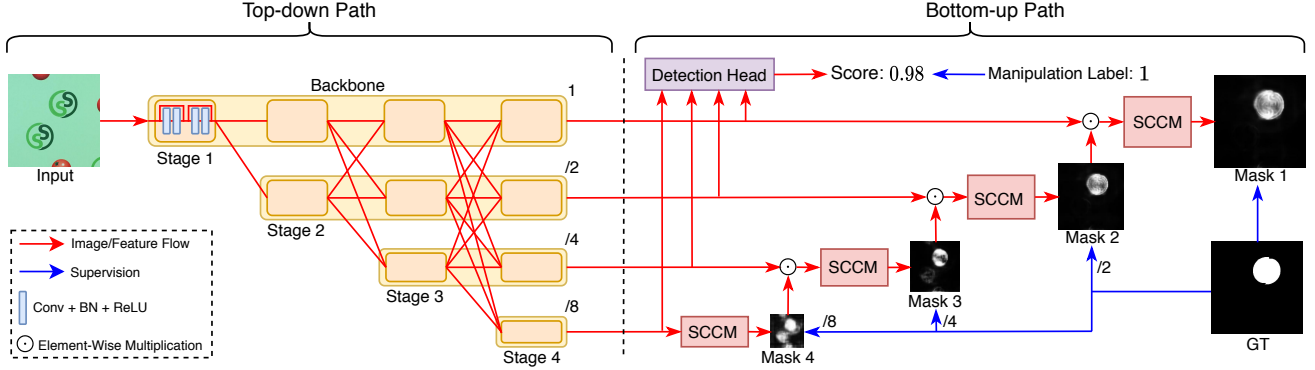
Figure 2: The architecture of the proposed PSCC-Net. The detection score predicted by the detection head indicates if the input is manipulated or not. The accuracy of manipulation localization from *Mask 4* to *Mask 1* is gradually improved, *e.g.*, the prediction of *Mask 4* confuses the pasted (forged) region with the pristine (copied) one, while *Mask 1* effectively fixes it.

spoof mediums [39].

Recent works attempt to tackle multiple forgeries in one model. J-LSTM [3] and H-LSTM [4] integrate the LSTM and CNN to capture the boundary-discriminative features. However, due to the patch-based design, both methods are time-consuming, and the size of detectable regions is limited by the preset patch size. RGB-N [69] adopts the steganalysis rich model [16] and Faster R-CNN [47], but it can only provide bounding boxes instead of segmentation masks. Later, ManTra-Net [65] learns features to distinguish 385 known manipulation types and treats the problem as anomaly detection. To learn the distinguishable features, auxiliary labeled data, such as camera sensors, are used. SPAN [22] extends ManTra-Net to further model the spatial correlation via local self-attention blocks and pyramid propagation. However, as the correlation is only considered in the local region, ManTra-Net and SPAN fail to take full advantage of the spatial correlation and consequently have limited generalizability. In this work, our PSCC-Net utilizes a progressive mechanism to improve the multi-scale feature representation and SSCM modules to better explore spatial and channel-wise correlations.

**Progressive mechanism** Progressive mechanism tackles a challenging task in a coarse-to-fine fashion. It has been widely adopted in many low-level and high-level vision tasks, such as denoising [34, 48], inpainting [66], super-resolution [8, 26], and object detection [6, 53, 68, 70]. The pyramid structure is commonly utilized to build multi-scale features. In this work, we propose a densely connected pyramid structure that progressively refines the manipulation mask from small scales to large ones, where each predicted mask serves as a prior for the next-scale estimation.

**Attention mechanism** The pioneer work [56] proposes an attention mechanism to improve the feature representation with relatively low cost, which has been widely employed in various vision tasks [12, 19, 21, 24, 25, 38, 58]. According

to the applied domain, the attention mechanism can be divided into two types: spatial attention [58] and channel-wise attention [21]. Recent works [17, 46, 61] take the benefit of both types to further improve the representation capability of DNN. These methods adopt separate schemes to explore the spatial and channel-wise attentions and thus require additional efforts to fuse them. In this work, a uniform SCCM jointly explores the image correlation and discrepancy in both spatial domain and feature channels, leading to better information sharing and faster inference.

## 3. PSCC-Net

Our PSCC-Net enables the detection and localization of various types of manipulations. As compared to the image-level detection, the pixel-level localization is more difficult. Therefore, PSCC-Net pays special attention to tackling the localization problem. Indeed, since the features for detection and localization are jointly learned, improving the localization performance will naturally benefit detection.

### 3.1. Network Architecture

#### 3.1.1 Top-Down Path

Most prior works use the conventional encoder-decoder [4, 69] and no-pooling structures [22, 65] to extract features. Since forged areas have various sizes, it is important to fuse local and global features to handle the scale variation. However, both structures extract features in a sequential pipeline and neglect feature fusion among different scales, and thus can only handle a limited scale variation. To address this issue, we adopt a light-weight backbone in [57], named HRNetV2p-W18. Following its default setting, the stage down-scaling ratio $s$ is set to 2, and there are totally 4 stages.

Compared to encoder-decoder and no-pooling structures, the benefits of our backbone are two-fold. First, features from different scales are computed in parallel. Hence, dense connections among different scales enable effective information exchange, which is beneficial for handling scale varia-

tions. Second, since the local and global feature fusion is performed for every scale, each feature contains sufficient information to predict a manipulation mask at the corresponding scale. Therefore, this backbone is in line with our progressive mechanism, where the prediction of each mask should rely on all local and global features to improve its accuracy. Indeed, except the predicted mask on the last scale, the others serve as a prior for the next-scale mask prediction. After the top-down path, the manipulated features on 4 scales are extracted. Then, we use the bottom-up path to perform manipulation detection and localization.

### 3.1.2 Bottom-Up Path

The bottom-up path in PSCC-Net estimates the detection score and the manipulation mask. Specifically, the detection score is predicted based on the extracted features from the top-down-path via a detection head [57], then the manipulation mask is generated through a progressive mechanism with full supervision. In particular, the coarse-to-fine progressive mechanism mimics how human tackles complicated problems in daily life.

We denote the input image as $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$. The extracted features at 4 scales are $\mathbf{F}_1 \in \mathbb{R}^{H \times W \times C}$, $\mathbf{F}_2 \in \mathbb{R}^{H/s \times W/s \times sC}$, $\mathbf{F}_3 \in \mathbb{R}^{H/s^2 \times W/s^2 \times s^2 C}$ and $\mathbf{F}_4 \in \mathbb{R}^{H/s^3 \times W/s^3 \times s^3 C}$, and their corresponding masks are denoted as $\mathbf{M}_1 \in \mathbb{R}^{H \times W}$, $\mathbf{M}_2 \in \mathbb{R}^{H/s \times W/s}$, $\mathbf{M}_3 \in \mathbb{R}^{H/s^2 \times W/s^2}$ and $\mathbf{M}_4 \in \mathbb{R}^{H/s^3 \times W/s^3}$. Here $H$, $W$, and $C$ are the height, width, and channel number of the image/feature respectively. Formally, we have

$$\mathbf{M}_{n-1} = f_{n-1}(\tau(\mathbf{M}_n) \cdot \mathbf{F}_{n-1}), \quad n = 2, 3, 4, \quad (1)$$

where $f_n$ denotes the SCCM on the $n$th scale, and $\tau$ is the upsampling operation (*e.g.*, the bilinear interpolation). Since $\mathbf{M}_4$ is the mask on the last scale, it can be directly expressed as $\mathbf{M}_4 = f_4(\mathbf{F}_4)$. For Scales 1-3, the feature on the current scale is associated with the upsampled mask from the previous scale for feature modulation. Then, the modulated feature is fed into SCCM to produce a manipulation mask.

To reduce the prediction difficulty, the proposed progressive mechanism avoids generating the mask at the finest scale directly. Instead, the mask on the coarsest scale is first predicted to locate the regions that are potentially forged based on current available information. The subsequent prediction on the finer scale can leverage the previous mask and pay more attention to those selected regions. This process continues until the generation of the manipulation mask at the finest scale, which serves as the final prediction. However, without explicit supervision on each scale, the intermediate masks might not follow the coarse-to-fine order. Therefore, full supervisions are applied on all scales to guide the mask estimation.
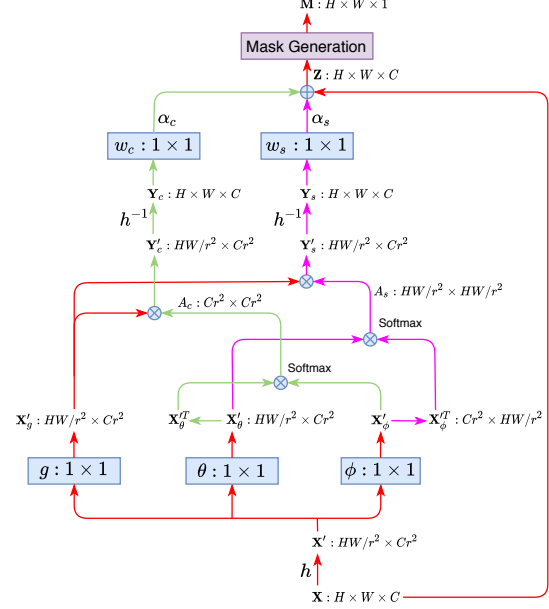


Figure 3: The structure of SCCM. Here $\otimes$ represents the matrix multiplication and $\oplus$ the element-wise addition; the red arrow shows the common feature flows; the pink and green arrows show the feature flows of spatial and channel-wise attentions respectively.

### 3.2. Spatio-Channel Correlation Module

Attention mechanisms are commonly used to modulate learned features according to their relative significance. As the final manipulation mask is binary, the localization can be considered as a pixel-level binary classification. Ideally, we expect the learned features on forged regions are similar to each other but distinct from those in pristine regions. In this case, a simple clustering method may suffice to produce an effective mask. Therefore, to better tackle manipulation localization, we propose a SCCM that employs the spatial attention to aggregate the pixel-level features based on their contextual correlations, and the channel-wise attention to consolidate the feature maps based on their channel correlations.

We illustrate the detailed structure of SCCM in Fig. 3, where the input feature $\mathbf{X}$ is of size $H \times W \times C$. Note that even though $\mathbf{X}$ is small ($256 \times 256$), the size of its spatial correlation can be enormous ($65,536 \times 65,536$), easily exceeding the memory limit. Therefore, we use function $h$ to reshape the input $\mathbf{X} \in \mathbb{R}^{H \times W \times C}$ to $\mathbf{X}' \in \mathbb{R}^{HW/r^2 \times Cr^2}$, where each feature map is flattened to form a vector based on SCCM down-scaling ratio $r$. This operation preserves all feature information and avoids modeling the spatial correlation of potentially large size $HW \times HW$.

To build the spatial and channel-wise correlations, one may directly leverage $\mathbf{X}'$. However, additional flexibility could be achieved by introducing the embedded Gaussian function [58]. Therefore, we use the $1 \times 1$ convolution

to build different functions $g$, $\theta$, and $\phi$ to transform $\mathbf{X}'$ into new linear embeddings as $\mathbf{X}'_g = g(\mathbf{X}')$, $\mathbf{X}'_\theta = \theta(\mathbf{X}')$, and $\mathbf{X}'_\phi = \phi(\mathbf{X}')$, all with the same size as $\mathbf{X}'$. Subsequently, the spatial and channel-wise correlations (denoted as $\mathbf{A}_s \in \mathbb{R}^{HW/r^2 \times HW/r^2}$ and $\mathbf{A}_c \in \mathbb{R}^{Cr^2 \times Cr^2}$) of embedded features $\mathbf{X}'_\theta$ and $\mathbf{X}'_\phi$ are computed, and the Gaussian operation is implemented by Softmax function. In the end, the spatial and channel-wise attentions are realized by performing matrix multiplications $\mathbf{A}_s \mathbf{X}'_g$ and $\mathbf{X}'_g \mathbf{A}_c$, respectively. Unlike prior methods [17, 46, 61] that employ two attentions on *different* features, we apply both on the *same* linear embedding for better information sharing and faster inference. Indeed, applying attentions in this way reduces the difficulty of subsequent fusion process, and also saves computational operations in SCCM. Specifically, the spatial attention can be formulated as:

$$\mathbf{Y}'_s = \mathbf{A}_s \mathbf{X}'_g = \text{softmax}(\mathbf{X}'_\theta \mathbf{X}'^T_\phi) \mathbf{X}'_g, \qquad (2)$$

where $\mathbf{Y}'_s \in \mathbb{R}^{HW/r^2 \times Cr^2}$ is the feature resulting from the application of spatial attention, and softmax$(\cdot)$ denotes the Softmax function. The element $(i, j)$ in $\mathbf{A}_s$ indicates the similarity between the feature vectors in the $i$th row of $\mathbf{X}'_\theta$ and $j$th row of $\mathbf{X}'_\phi$. The more similar they are, the higher correlation they have. This helps the network to learn feature representations for distinguishing forged regions from pristine ones and avoid overfitting to a specific attack type in training. Similarly, the channel-wise attention is expressed as:

$$\mathbf{Y}'_c = \mathbf{X}'_g \mathbf{A}_c = \mathbf{X}'_g \text{softmax}(\mathbf{X}'^T_\theta \mathbf{X}'_\phi), \qquad (3)$$

where $\mathbf{Y}'_c \in \mathbb{R}^{HW/r^2 \times Cr^2}$ is the feature resulting from the application of channel-wise attention. The element $(i, j)$ in $\mathbf{A}_c$ measures the similarity between the channel maps in the $i$th column of $\mathbf{X}'_\theta$ and $j$th column of $\mathbf{X}'_\phi$. Since the response from different channels might be associated with the same class, *e.g.*, manipulated or pristine, the channel-wise correlation aggregates feature maps based on their similarities to enhance the representation in forged regions.

We use $h^{-1}$ to reshape $\mathbf{Y}'_s$ and $\mathbf{Y}'_c$ respectively back to $\mathbf{Y}_s$ and $\mathbf{Y}_c$ of size $H \times W \times C$. Further, two functions $\omega_s$ and $\omega_c$ are built by $1 \times 1$ convolution to improve their feature representations. The output features from $\omega_s$ and $\omega_c$ are complement to each other. As it is non-trivial to determine their relative significance, two learnable parameters $\alpha_s$ and $\alpha_c$, both initialized as 1, are used for trade-off. We also adopt the residual learning [20] to express the feature $\mathbf{Z}$ as:

$$\mathbf{Z} = \mathbf{X} + \alpha_s \cdot \omega_s(\mathbf{Y}_s) + \alpha_c \cdot \omega_c(\mathbf{Y}_c). \qquad (4)$$

The final output of SCCM is a predicted mask with only one channel. To reduce the channel number in $\mathbf{Z}$, we employ a mask generation block with the sequential order of *Conv-ReLU-Conv-Sigmoid*, where *Conv* is a $3 \times 3$ convolution.

## 3.3. Loss Function

To train the PSCC-Net, we adopt the binary cross-entropy loss ($L_{bce}$) for both detection and localization tasks. The predicted detection score ($s_d$) is supervised by the ground-truth (GT) label ($l_d$) with 0 standing for pristine image and 1 for forged image. Moreover, full supervisions are applied on each predicted mask by downsampling the GT mask $\mathbf{G}_1$ to $\mathbf{G}_2$, $\mathbf{G}_3$, and $\mathbf{G}_4$ according to their corresponding sizes, with 0 standing for pristine pixel and 1 for forged pixel. The masks predicted through the progressive mechanism at different scales are considered to be of equal importance. Therefore, our final loss function $\hat{L}$ can be expressed as:

$$\hat{L} = L_{bce}(s_d, l_d) + \frac{1}{4} \sum\nolimits_{m=1}^{4} L_{bce}(\mathbf{M}_m, \mathbf{G}_m). \qquad (5)$$

## 3.4. Training Data Synthesis

Since there is no standard IMDL dataset for training, a synthetic dataset is built to train and validate our PSCC-Net. This dataset includes four categories 1) splicing, 2) copy-move, 3) removal, and 4) pristine classes. For splicing, inspired by [40, 62], we use the MS COCO [36] to generate spliced images, where one annotated region is randomly selected per image, and pasted into a different image after several transformations. We adopt the same transformation as [62] including the scale, rotation, shift and luminance changes. Since the spliced region is not necessarily an object, we use the Bezier curve [43] to generate random contours, then fill them to produce splicing masks. We follow the same processes above but randomly select donor and target images in KCMI [52], VISION [51], and Dresden [18] that are commonly used to identify camera source [7], to generate additional spliced images as supplementary. For copy-move, the dataset from [63] is adopted. For removal, we adopt the SOTA inpainting method [35] to fill one annotated region that is randomly removed from each chosen MS COCO image. As to the pristine class, we simply select images from the original datasets mentioned above.

In summary, we have ∼100k images per class, thus 400k in total. As it is inefficient to train all manipulated images in one epoch, we uniformly sample 25k images per class to form a 100k dataset on-the-fly for training in each epoch. In addition, we also build a validation set that contains $4 \times 100$ images. The size of synthetic images are all set to $256 \times 256$.

## 4. Experiments

### 4.1. Experimental Setup

**Test data**   We evaluate the manipulation localization on 4 standard test datasets: Columbia [44], Coverage [59], CASIA [14] and NIST16 [1], and 1 real-world dataset: IMD20 [45]. To finetune PSCC-Net, we follow the same training/testing split on Coverage, CASIA, and NIST16 as

in [22, 69] for fair comparisons. Specifically, Columbia [44] is a splicing dataset of 180 images. Coverage [59] is a copy-move dataset of 100 images; for fine-tuning, it is split into 75/25 for training and testing. CASIA [14] (v1.0 + v2.0) includes both splicing and copy-move; for fine-tuning, 5, 123 images from v2.0 is adopted for training, and 921 images from v1.0 is for testing. NIST16 [1] has 564 images, involving all three manipulations; for fine-tuning, 404 images are used for training and 160 for testing. IMD20 [45] consists of 2, 010 real-life manipulated images collected from Internet.

As the manipulation detection is not considered by prior works, there is no standard dataset for benchmarking. To address this issue, we include both forged and pristine images in CASIA dataset and define a evaluation protocol for detection. This dataset is named CASIA-D and consists of 1, 842 images with 50% forged and 50% pristine.

**Metrics** To quantify the localization performance, following previous works [22, 65], we use pixel-level Area Under Curve (AUC) and F1 score on manipulation masks. To evaluate the detection performance, we use image-level AUC and F1 score, Equal Error Rate (EER), and True Positive Rate at 1% false positive rate (TPR$_{1\%}$). Since binary masks and detection scores are required to compute F1 scores, we adopt the EER threshold to binarize them.

**Implementation details** PSCC-Net is end-to-end trainable and light-weighted. Its top-down path and bottom-up path have 2.0 and 1.6 Million (M) parameters. In the bottom-up path, the detection head has 0.9 M and the rest part (for localization) has only 0.7 M parameters. In comparison, the ManTra-Net [65] and SPAN [22] have 3.8 and 3.7 M parameters, respectively. Implemented by PyTorch, our model is trained with GeForce GTX 1080Ti. We initialize our backbone with ImageNet pre-trained weights, and optimize the whole model by Adam [28] with a batch size of 10 and an initial learning rate of 2$e$-4. The learning rate is halved every 5 epochs and the total training period is 25 epochs.

Our network can take arbitrary-size images as input. To avoid performance degradation caused by size mismatch between training (*e.g.*, $256 \times 256$) and testing data (*e.g.*, $4, 000 \times 3, 000$), at the end of top-down path, we resample the extracted features from the first to the last scales respectively into fixed sizes $256 \times 256$, $128 \times 128$, $64 \times 64$, and $32 \times 32$, where the ratio $r$ in SCCM is set to 4, 2, 2, and 1 respectively to reduce the computational burden. The produced masks are resampled back to the same size as the input image for localization evaluation.

### 4.2. Comparisons on Localization

Our baseline IMDL methods include ELA [30], NOI1 [42], CFA1 [15], J-LSTM [3], H-LSTM [4], RGB-N [69], ManTra-Net [65], and SPAN [22] where SPAN has reported the SOTA performance on localization. Following the evaluation protocol defined in SPAN [22], we compare

| Method | Columbia | Coverage | CASIA | NIST16 | IMD20 |
|---|---|---|---|---|---|
| ManTra-Net [65] | 82.4 | 81.9 | 81.7 | 79.5 | 74.8 |
| SPAN [22] | 93.6 | **92.2** | 79.7 | 84.0 | 75.0 |
| PSCC-Net | **98.2** | 84.7 | **82.9** | **85.5** | **80.6** |

Table 1: Localization AUC (%) of pre-trained models.

| Method | Type | Coverage | CASIA | NIST16 |
|---|---|---|---|---|
| ELA [30] | U | 58.3 / 22.2 | 61.3 / 21.4 | 42.9 / 23.6 |
| NOI1 [42] | U | 58.7 / 26.9 | 61.2 / 26.3 | 48.7 / 28.5 |
| CFA1 [15] | U | 48.5 / 19.0 | 52.2 / 20.7 | 50.1 / 17.4 |
| J-LSTM [3] | F | 61.4 / - | - / - | 76.4 / - |
| H-LSTM [4] | F | 71.2 / - | - / - | 79.4 / - |
| RGB-N [69] | F | 81.7 / 43.7 | 79.5 / 40.8 | 93.7 / 72.2 |
| SPAN [22] | F | 93.7 / 55.8 | 83.8 / 38.2 | 96.1 / 58.2 |
| PSCC-Net | F | **94.1 / 72.3** | **87.5 / 55.4** | **99.6 / 81.9** |

Table 2: Evaluation of the fine-tuned models. Localization AUC/F1s are reported (in %). Type U denotes an unsupervised model, and type F denotes a fine-tuned model. ManTra-Net is not shown here as it has only developed the pre-trained model.

the localization performance using two models: 1) the pre-trained model is trained on the synthetic dataset and evaluated on the *full* test datasets, and 2) the fine-tuned model is the pre-trained model further fine-tuned on the training split of test datasets and evaluated on their *test split*.

**Pre-trained model** We choose the best pre-trained model based on the performance on our validation set. Tab. 1 shows the localization performance of pre-trained models for different methods on 4 standard datasets and 1 real-world dataset under pixel-level AUC. The pre-trained PSCC-Net achieves the best localization performance on Columbia, CASIA, NIST16, and IMD20, and ranks the second on Coverage. The most significant performance gain is achieved while tackling real-life manipulated images (5.6% ↑). This validates that the PSCC-Net has the best generalization ability as compared to the others. We fail to achieve the best performance on Coverage, despite surpassing ManTra-Net 2.8% under AUC. The reason might be the imperfection of our training data for the case, where the copied object is intentionally moved to cover a pristine object with similar appearance. Indeed, by fine-tuning the pre-trained model on Coverage, PSCC-Net achieves the 0.4% gain over SPAN under AUC (Tab. 2).

**Fine-tuned model** We further fine-tune the pre-trained model on specific datasets using our training strategy. The cross validation on training data helps to select the best fine-tuned models on each test dataset. We compare the fine-tuned models in Tab. 2. For AUC, PSCC-Net surpasses baselines in all cases (over 2.5% to SPAN on average). As for F1 score, our model outperforms them with a large margin (over 19% to SPAN on average).
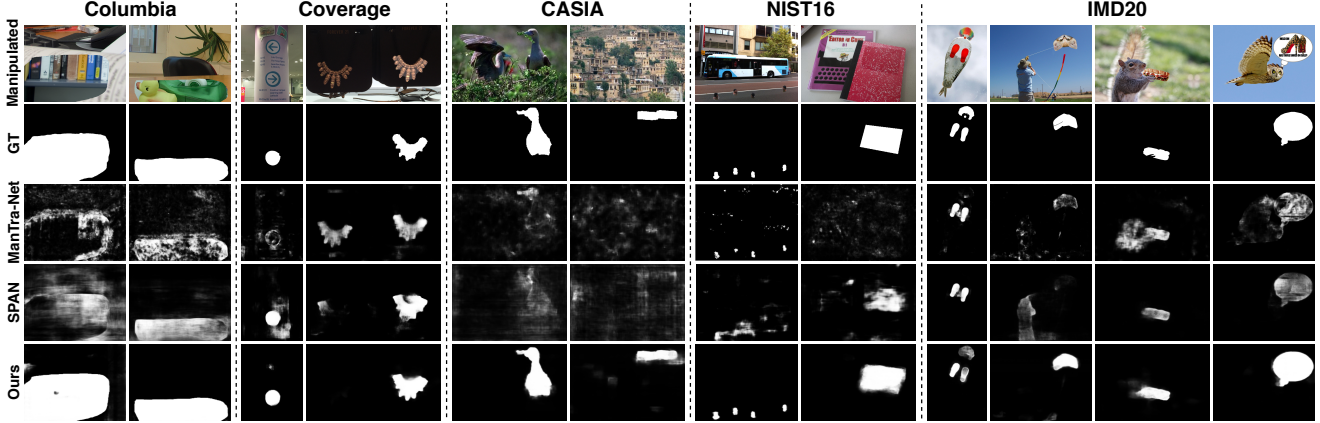
Figure 4: Qualitative localization evaluations on 5 datasets. From top to bottom, we show manipulated images, GT manipulation masks, predictions of ManTra-Net, SPAN, and ours. Best models are used to produce masks. Zoom in for details. See Suppl. for more results.
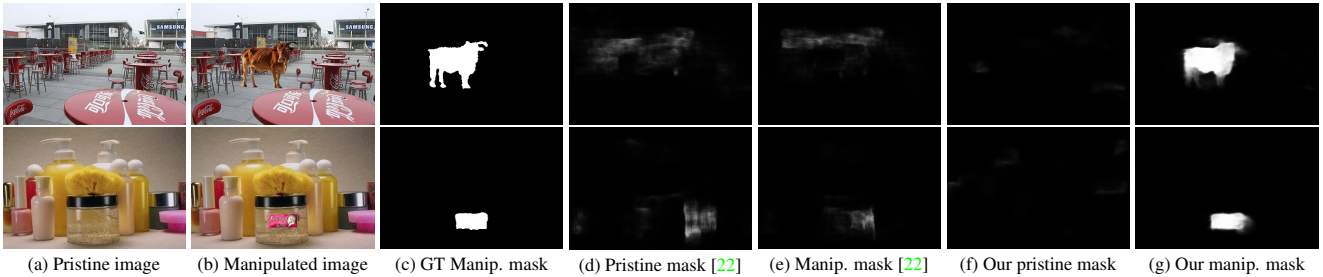


| (a) Pristine image | (b) Manipulated image | (c) GT Manip. mask | (d) Pristine mask [22] | (e) Manip. mask [22] | (f) Our pristine mask | (g) Our manip. mask |

Figure 5: Qualitative detection evaluations on CASIA-D. Since GT pristine masks are blank, they are not shown here to save space.

## 4.3. Comparisons on Detection

Since ManTra-Net and SPAN are the best performing baselines in the localization evaluation, and ManTra-Net does not develop the fine-tuned model, we choose to use the pre-trained model for detection evaluation, in order to make comparisons to both of them. Although these two baselines make no direct attempt to perform detection, their estimated manipulation masks can be leveraged for this purpose. As such, we simply regard the average of the mask as their scores. For fair comparisons, we build a variant that adopts the same averaging strategy to calculate this score, denoted as PSCC-Net$^{\dagger}$. In Tab. 3, owing to our well-predicted manipulation masks, the PSCC-Net$^{\dagger}$ achieves the best detection performance on all used metrics. It is evident that the detection performance can be dramatically improved by introducing a tailored head. With a favorable detection, the IMDL methods can be more efficient. That is, detection is performed before localization, and only the detected forgery is passed for localization. Our network design is compatible with this efficiency consideration as the detection head is placed at the beginning of the bottom-up path.

## 4.4. Visualization, Ablation and Analysis

**Qualitative results** We provide qualitative evaluations of manipulation localization and detection in Figs. 4, 5. PSCC-

| Method | AUC ↑ | F1 ↑ | EER ↓ | TPR$_{1\%}$ ↑ |
|---|---|---|---|---|
| ManTra-Net [65] | 59.94 | 56.69 | 43.21 | 5.43 |
| SPAN [22] | 67.33 | 63.48 | 36.47 | 5.54 |
| PSCC-Net$^{\dagger}$ | 74.40 | 66.88 | 33.21 | 28.37 |
| PSCC-Net | **99.65** | **97.12** | **2.83** | **95.65** |

Table 3: Detection evaluation on CASIA-D, all reported in %.

Net predicts more accurate and sharper manipulation masks while maintaining low false alarms on pristine regions, especially for *small* manipulations.

**Visualization of SCCM** To provide insights into SCCM, we visualize the spatial response map for forged and pristine pixels in $\mathbf{M}_3$, by examining its spatial correlation represented in $\mathbf{A}_s$. After interpolation, each row of $\mathbf{A}_s$ is associated with one pixel (*e.g.*, $P_1$) in the manipulated image, and its grayscale spatial response map can be obtained by reshaping this row vector from $1 \times HW$ to $H \times W$ (*e.g.*, $P_1$ response). In Fig. 6 (a), 2 examples of splicing and copy-move manipulations from CASIA are shown in the 1st and 2nd rows. We select 3 representative pixels for each image and annotate as $P_1$, $P_2$, and $P_3$, where $P_1$ and $P_2$ are from forged regions, and $P_3$ is from pristine regions. We project their grayscale spatial response maps into *Jet* color space and overlay them on the manipulated image as in Figs. 6 (c-e). It is evident

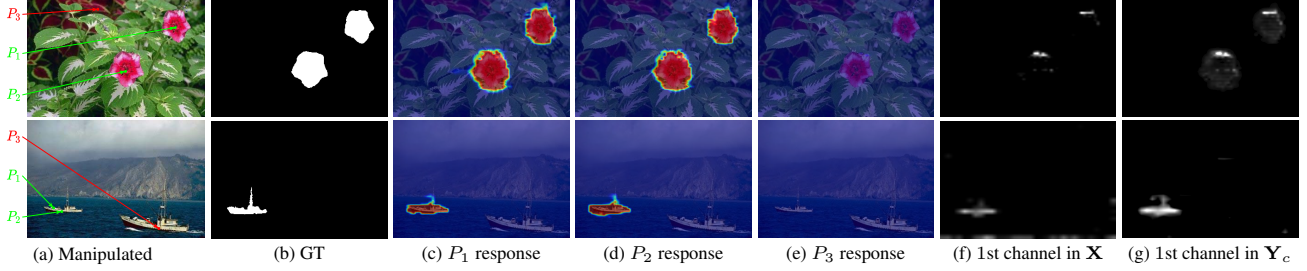|               | (a) Manipulated | (b) GT | (c) $P_1$ response | (d) $P_2$ response | (e) $P_3$ response | (f) 1st channel in $\mathbf{X}$ | (g) 1st channel in $\mathbf{Y}_c$ |

Figure 6: Visualization of spatial and channel-wise attentions in SCCM. For each row, we show a manipulated image, its GT mask, 3 spatial response maps (one for each selected pixel), and the 1st channel map in $\mathbf{X}$ and $\mathbf{Y}_c$. Zoom in for details.

| Variants | Columbia | Coverage | CASIA | NIST16 | Time |
|---|---|---|---|---|---|
| Mask 4 | 93.34 / 79.22 | 82.99 / 44.23 | 81.49 / 31.69 | 84.15 / 30.55 | 0.63 |
| Mask 3 | 98.08 / 92.41 | 83.48 / 47.29 | 82.55 / 34.64 | 85.25 / 33.55 | 0.75 |
| Mask 2 | 98.18 / 93.32 | 84.44 / 49.08 | 82.78 / 35.59 | 85.38 / 34.94 | 0.88 |
| w/o CA+SA | 85.78 / 70.32 | 79.95 / 43.27 | 79.26 / 31.06 | 79.58 / 31.73 | 0.84 |
| w/o SA | 90.70 / 75.68 | 80.56 / 43.50 | 79.51 / 31.08 | 83.49 / 32.34 | 0.92 |
| w/o CA | 94.50 / 85.34 | 82.16 / 45.04 | 82.63 / 35.97 | 84.65 / 33.42 | 0.92 |
| PSCC-Net | **98.19 / 93.45** | **84.65 / 49.78** | **82.93 / 36.27** | **85.47 / 35.73** | 1.00 |

Table 4: Ablation study of PSCC-Net. Average AUC/F1s are reported (in %). The run time (in proportion) is relative to that of PSCC-Net. Our full model takes $0.019s$ to process one $1,080$P image, whereas ManTra-Net and SPAN take $0.208s$ and $0.161s$, respectively. Terminating the prediction earlier on *Mask 4* can shorten the run time to $0.012s$, *i.e.*, $\sim 37\%$ additional saving.

| Distortion | Columbia | | | NIST16 | | |
|---|---|---|---|---|---|---|
| | [65] | [22] | Ours | [65] | [22] | Ours |
| Resize ($0.78\times$) | 71.66 | 89.99 | **93.40** | 77.43 | 83.24 | **85.29** |
| Resize ($0.25\times$) | 68.64 | 69.08 | **78.41** | 75.52 | 80.32 | **85.01** |
| GaussianBlur ($k = 3$) | 67.72 | 78.97 | **84.18** | 77.46 | 83.10 | **85.38** |
| GaussianBlur ($k = 15$) | 62.88 | 67.70 | **73.24** | 74.55 | 79.15 | **79.93** |
| GaussianNoise ($\sigma = 3$) | 68.22 | 75.11 | **82.64** | 67.41 | 75.17 | **78.42** |
| GaussianNoise ($\sigma = 15$) | 54.97 | 65.80 | **74.35** | 58.55 | 67.28 | **76.65** |
| JPEGCompress ($q = 100$) | 75.00 | 93.32 | **97.97** | 77.91 | 83.59 | **85.40** |
| JPEGCompress ($q = 50$) | 59.37 | 74.62 | **89.11** | 74.38 | 80.68 | **85.37** |
| w/o distortion | 77.95 | 93.60 | **98.19** | 78.05 | 83.95 | **85.47** |

Table 5: Robustness comparison with respect to various distortions. AUCs are reported (in %).

that the spatial response maps of $P_1$ and $P_2$ have high values in forged regions and low values in pristine regions, but the map of $P_3$ retains low values in all regions including the one providing the copied content (*e.g.,* the $P_3$ response in the 2nd row of Fig. 6 (e)). This visualization indicates that the features in forged regions are clustered together, thus justifies the effectiveness of spatial attention in SCCM.

For channel-wise correlation $\mathbf{A}_c$, it is hard to provide a comprehensible visualization. Instead, we choose to visualize one channel of $\mathbf{Y}_c$ and compare it to the same channel of $\mathbf{X}$ to see if any region is enhanced. We visualize the 1st channel of $\mathbf{X}$ and $\mathbf{Y}_c$ in Figs. 6 (f,g). Indeed, the forged region in $\mathbf{Y}_c$ is consolidated compared to the one in $\mathbf{X}$, which proves the effectiveness of channel-wise attention in SCCM.

**Ablation study** To justify our network design, we test several variants of PSCC-Net to show the effectiveness of progressive mechanism and SCCM in Tab. 4, where all variants are pre-trained on our dataset. *Mask 4*, *Mask 3*, and *Mask 2* are the variants that truncate the original model after generating manipulation masks on the 4th, 3rd, and 2nd scales, respectively. The comparisons of *Mask 4*, *Mask 3*, *Mask 2*, and the original PSCC-Net demonstrate the gradual improvement in performance, which is a clear manifestation of our progressive mechanism. Since *Mask 4* performs well under AUC and F1 scores, the mask prediction can be terminated earlier to save time. The comparisons among the variants without spatial and channel-wise attentions (*w/o SA+CA*),

without spatial attention (*w/o SA*), without channel-wise attention (*w/o CA*), and original PSCC-Net illustrate that both SA and CA outperform the baseline (*w/o SA+CA*), where the performance gain acquired from SA is more than that from CA. Owing to SCCM, the original PSCC-Net achieves the best performance as compared to its attention variants.

**Robustness analysis** To analyze the robustness of PSCC-Net for localization, we follow the distortion settings in [22] to degrade the raw manipulated images from Columbia and NIST16. These distortions include resizing images to a different scale, applying Gaussian blur with kernel size $k$, adding Gaussian noise with standard deviation $\sigma$, and performing compression with quality factor $q$. Table 5 shows the robustness analysis under pixel-level AUC with pre-trained models. The PSCC-Net is more robust than ManTra-Net and SPAN under all distortions. It is worth noting that resizing is commonly performed when uploading images to social media. Indeed, benefiting from the operation that resamples the manipulation features into the fixed sizes, the impact of resizing to PSCC-Net is the least as compared to the others.

## 5. Conclusion

In this work, a novel PSCC-Net is proposed to meet the challenge of advanced image manipulation techniques. We employ a progressive mechanism to predict the manipulation mask on all backbone scales, where each mask serves as a prior to help predict the next-scale mask. Moreover, a SCCM is designed to perform spatial and channel-wise attentions

on extracted features, which provides holistic information to make our model more generalized to manipulation attacks. Extensive experiments demonstrate that our PSCC-Net outperforms the SOTA methods on both detection and localization. For future work, we will develop techniques for estimating the uncertainty of predicted manipulation masks.

## Acknowledgement

## References

[1] NIST: Nist nimble 2016 datasets. https://www.nist.gov/itl/iad/mig/, 2016. 5, 6

[2] Irene Amerini, Tiberio Uricchio, Lamberto Ballan, and Roberto Caldelli. Localization of JPEG double compression through multi-domain convolutional neural networks. In *CVPRW*, 2017. 2

[3] Jawadul H Bappy, Amit K Roy-Chowdhury, Jason Bunk, Lakshmanan Nataraj, and BS Manjunath. Exploiting spatial structure for localizing manipulated image regions. In *ICCV*, 2017. 2, 3, 6

[4] Jawadul H Bappy, Cody Simons, Lakshmanan Nataraj, BS Manjunath, and Amit K Roy-Chowdhury. Hybrid LSTM and encoder-decoder architecture for detection of image forgeries. *TIP*, 2019. 2, 3, 6

[5] Luca Bondi, Silvia Lameri, David Güera, Paolo Bestagini, Edward J Delp, and Stefano Tubaro. Tampering detection and localization through clustering of camera-based CNN features. In *CVPRW*, 2017. 2

[6] Garrick Brazil and Xiaoming Liu. Pedestrian detection with autoregressive network phases. In *CVPR*, 2019. 3

[7] Chang Chen, Zhiwei Xiong, Xiaoming Liu, and Feng Wu. Camera trace erasing. In *CVPR*, 2020. 1, 5

[8] Yu Chen, Ying Tai, Xiaoming Liu, Chunhua Shen, and Jian Yang. FSRNet: End-to-end learning face super-resolution with facial priors. In *CVPR*, 2018. 3

[9] Davide Cozzolino, Giovanni Poggi, and Luisa Verdoliva. Efficient dense-field copy–move forgery detection. *TIFS*, 2015. 2

[10] Davide Cozzolino, Giovanni Poggi, and Luisa Verdoliva. Splicebuster: A new blind image splicing detector. In *International Workshop on Information Forensics and Security (WIFS)*, 2015. 2

[11] Davide Cozzolino and Luisa Verdoliva. Noiseprint: a CNN-based camera model fingerprint. *TIFS*, 2019. 1, 2

[12] Hao Dang, Feng Liu, Joel Stehouwer, Xiaoming Liu, and Anil Jain. On the detection of digital face manipulation. In *CVPR*, 2020. 1, 2, 3

[13] Helisa Dhamo, Azade Farshad, Iro Laina, Nassir Navab, Gregory D Hager, Federico Tombari, and Christian Rupprecht. Semantic image manipulation using scene graphs. In *CVPR*, 2020. 1

[14] Jing Dong, Wei Wang, and Tieniu Tan. Casia image tampering detection evaluation database. In *China Summit and International Conference on Signal and Information Processing*, 2013. 5, 6

[15] Pasquale Ferrara, Tiziano Bianchi, Alessia De Rosa, and Alessandro Piva. Image forgery localization via fine-grained analysis of CFA artifacts. *TIFS*, 2012. 6

[16] Jessica Fridrich and Jan Kodovsky. Rich models for steganalysis of digital images. *TIFS*, 2012. 3

[17] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *CVPR*, 2019. 3, 5

[18] Thomas Gloe and Rainer Böhme. The 'dresden image database' for benchmarking digital image forensics. In *ACM Symposium on Applied Computing*, 2010. 5

[19] Sixue Gong, Xiaoming Liu, and Anil Jain. Mitigating face recognition bias via group adaptive classifier. In *CVPR*, 2021. 3

[20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5

[21] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, 2018. 3

[22] Xuefeng Hu, Zhihan Zhang, Zhenye Jiang, Syomantak Chaudhuri, Zhenheng Yang, and Ram Nevatia. SPAN: Spatial pyramid attention network for image manipulation localization. In *ECCV*, 2020. 1, 2, 3, 6, 7, 8

[23] Minyoung Huh, Andrew Liu, Andrew Owens, and Alexei A Efros. Fighting fake news: Image splice detection via learned self-consistency. In *ECCV*, 2018. 1, 2

[24] Ashraful Islam, Chengjiang Long, Arslan Basharat, and Anthony Hoogs. DOA-GAN: Dual-order attentive generative adversarial network for image copy-move forgery detection and localization. In *CVPR*, 2020. 2, 3

[25] Takashi Isobe, Songjiang Li, Xu Jia, Shanxin Yuan, Gregory Slabaugh, Chunjing Xu, Ya-Li Li, Shengjin Wang, and Qi Tian. Video super-resolution with temporal group attention. In *CVPR*, 2020. 3

[26] Kui Jiang, Zhongyuan Wang, Peng Yi, Chen Chen, Baojin Huang, Yimin Luo, Jiayi Ma, and Junjun Jiang. Multi-scale progressive fusion network for single image deraining. In *CVPR*, 2020. 3

[27] Rani Mariya Joseph and AS Chithra. Literature survey on image manipulation detection. *International Research Journal of Engineering and Technology (IRJET)*, 2015. 2

[28] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6

[29] Vladimir V Kniaz, Vladimir Knyaz, and Fabio Remondino. The point where reality meets fantasy: Mixed adversarial generators for image splice detection. In *NeurIPS*, 2019. 2

[30] Neal Krawetz and Hacker Factor Solutions. A picture's worth. *Hacker Factor Solutions*, 2007. 6

[31] Avisek Lahiri, Arnav Kumar Jain, Sanskar Agrawal, Pabitra Mitra, and Prabir Kumar Biswas. Prior guided GAN based semantic inpainting. In *CVPR*, 2020. 1

[32] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. MaskGAN: Towards diverse and interactive facial image manipulation. In *CVPR*, 2020. 1

[33] Bowen Li, Xiaojuan Qi, Thomas Lukasiewicz, and Philip HS Torr. ManiGAN: Text-guided image manipulation. In *CVPR*, 2020. 1

[34] Jingyuan Li, Fengxiang He, Lefei Zhang, Bo Du, and Dacheng Tao. Progressive reconstruction of visual structure for image inpainting. In *ICCV*, 2019. 1, 3

[35] Jingyuan Li, Ning Wang, Lefei Zhang, Bo Du, and Dacheng Tao. Recurrent feature reasoning for image inpainting. In *CVPR*, 2020. 1, 5

[36] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014. 5

[37] Xihui Liu, Zhe Lin, Jianming Zhang, Handong Zhao, Quan Tran, Xiaogang Wang, and Hongsheng Li. Open-Edit: Open-domain image manipulation with open-vocabulary instructions. In *ECCV*, 2020. 1

[38] Xiaohong Liu, Yongrui Ma, Zhihao Shi, and Jun Chen. Grid-DehazeNet: Attention-based multi-scale network for image dehazing. In *ICCV*, 2019. 3

[39] Yaojie Liu, Joel Stehouwer, and Xiaoming Liu. On disentangling spoof traces for generic face anti-spoofing. In *ECCV*, 2020. 3

[40] Yaqi Liu, Xiaobin Zhu, Xianfeng Zhao, and Yun Cao. Adversarial learning for constrained image splicing detection and localization based on atrous convolution. *TIFS*, 2019. 5

[41] Siwei Lyu, Xunyu Pan, and Xing Zhang. Exposing region splicing forgeries with blind local noise estimation. *IJCV*, 2014. 2

[42] Babak Mahdian and Stanislav Saic. Using noise inconsistencies for blind image forensics. *Image and Vision Computing*, 2009. 6

[43] Michael E Mortenson. *Mathematics for computer graphics applications*. Industrial Press Inc., 1999. 5

[44] Tian-Tsong Ng, Jessie Hsu, and Shih-Fu Chang. Columbia image splicing detection evaluation dataset. *DVMM lab. Columbia Univ CalPhotos Digit Libr*, 2009. 5, 6

[45] Adam Novozamsky, Babak Mahdian, and Stanislav Saic. IMD2020: A large-scale annotated dataset tailored for detecting manipulated images. In *WACVW*, 2020. 5, 6

[46] Jongchan Park, Sanghyun Woo, Joon-Young Lee, and In So Kweon. BAM: Bottleneck attention module. In *BMVC*, 2018. 3, 5

[47] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 3

[48] Wenqi Ren, Lin Ma, Jiawei Zhang, Jinshan Pan, Xiaochun Cao, Wei Liu, and Ming-Hsuan Yang. Gated fusion network for single image dehazing. In *CVPR*, 2018. 3

[49] Ronald Salloum, Yuzhuo Ren, and C-C Jay Kuo. Image splicing localization using a multi-task fully convolutional network (MFCN). *Journal of Visual Communication and Image Representation*, 2018. 2

[50] Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of GANs for semantic face editing. In *CVPR*, 2020. 1

[51] Dasara Shullani, Marco Fontani, Massimo Iuliani, Omar Al Shaya, and Alessandro Piva. VISION: a video and image dataset for source identification. *EURASIP Journal on Information Security*, 2017. 5

[52] IEEE's Signal Processing Society. Camera model identification. https://www.kaggle.com/c/sp-society-camera-model-identification. 5

[53] Xiaolin Song, Kaili Zhao, Wen-Sheng Chu Honggang Zhang, and Jun Guo. Progressive refinement network for occluded pedestrian detection. In *ECCV*, 2020. 3

[54] T.J. Thomson, Daniel Angus, and Paula Dootson. Seeing no longer means believing. *In Daily*. 1

[55] Ruben Tolosana, Ruben Vera-Rodriguez, Julian Fierrez, Aythami Morales, and Javier Ortega-Garcia. Deepfakes and beyond: A survey of face manipulation and fake detection. *Information Fusion*, 2020. 1

[56] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 3

[57] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learning for visual recognition. *PAMI*, 2020. 2, 3, 4

[58] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, 2018. 3, 4

[59] Bihan Wen, Ye Zhu, Ramanathan Subramanian, Tian-Tsong Ng, Xuanjing Shen, and Stefan Winkler. COVERAGE–A novel database for copy-move forgery detection. In *ICIP*, 2016. 2, 5, 6

[60] AJ Willingham. Is that video real? *CNN*. 1

[61] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. CBAM: Convolutional block attention module. In *ECCV*, 2018. 3, 5

[62] Yue Wu, Wael Abd-Almageed, and Prem Natarajan. Deep matching and validation network: An end-to-end solution to constrained image splicing localization and detection. In *ACMMM*, 2017. 2, 5

[63] Yue Wu, Wael Abd-Almageed, and Prem Natarajan. BusterNet: Detecting copy-move image forgery with source/target localization. In *ECCV*, 2018. 2, 5

[64] Yue Wu, Wael Abd-Almageed, and Prem Natarajan. Image copy-move forgery detection via an end-to-end deep neural network. In *WACV*, 2018. 2

[65] Yue Wu, Wael AbdAlmageed, and Premkumar Natarajan. ManTra-Net: Manipulation tracing network for detection and localization of image forgeries with anomalous features. In *CVPR*, 2019. 1, 2, 3, 6, 7, 8

[66] Peng Yi, Zhongyuan Wang, Kui Jiang, Junjun Jiang, and Jiayi Ma. Progressive fusion video super-resolution network via exploiting non-local spatio-temporal correlations. In *ICCV*, 2019. 3

[67] Yu Zeng, Zhe Lin, Jimei Yang, Jianming Zhang, Eli Shechtman, and Huchuan Lu. High-resolution image inpainting

with iterative confidence feedback and guided upsampling. In *ECCV*, 2020. 1

[68] Xiaoning Zhang, Tiantian Wang, Jinqing Qi, Huchuan Lu, and Gang Wang. Progressive attention guided recurrent network for salient object detection. In *CVPR*, 2018. 3

[69] Peng Zhou, Xintong Han, Vlad I Morariu, and Larry S Davis. Learning rich features for image manipulation detection. In *CVPR*, 2018. 2, 3, 6

[70] Jiashu Zhu, Dong Li, Tiantian Han, Lu Tian, and Yi Shan. ProgressFace: Scale-aware progressive learning for face detection. In *ECCV*, 2020. 3

[71] Xinshan Zhu, Yongjun Qian, Xianfeng Zhao, Biao Sun, and Ya Sun. A deep learning approach to patch-based image inpainting forensics. *Signal Processing: Image Communication*, 2018. 2