

# Zooming into Face Forensics: A Pixel-level Analysis

Jia Li<sup>1</sup>, Tong Shen<sup>2</sup>, Wei Zhang<sup>2</sup>, Hui Ren<sup>1</sup>, Dan Zeng<sup>3</sup>, Tao Mei<sup>2</sup>

<sup>1</sup>Communication University of China

<sup>2</sup>JD AI Research

<sup>3</sup>Shanghai University

## Abstract

*The stunning progress in face manipulation methods has made it possible to synthesize realistic fake face images, which poses potential threats to our society. It is urgent to have face forensics techniques to distinguish those tampered images. A large scale dataset FaceForensics++ has provided enormous training data generated from prominent face manipulation methods to facilitate anti-fake research. However, previous works focus more on casting it as a classification problem by only considering a global prediction. Through investigation to the problem, we find that training a classification network often fails to capture high quality features, which might lead to sub-optimal solutions. In this paper, we zoom in on the problem by conducting a pixel-level analysis, i.e. formulating it as a pixel-level segmentation task. By evaluating multiple architectures on both segmentation and classification tasks, We show the superiority of viewing the problem from a segmentation perspective. Different ablation studies are also performed to investigate what makes an effective and efficient anti-fake model. Strong baselines are also established, which, we hope, could shed some light on the field of face forensics.*

## 1. Introduction

Human faces play an important role in human communication, as a face is associated with the identity of a person. The unique face information, working as fingerprints, has been used in many applications such as phone unlocking, payment, etc., thanks to remarkable progress in face detection and recognition systems [45, 41, 37]. However, we have also seen stunning progress in image and video manipulation methods, which enable editing the images or videos in a visually plausible way. Some face specific manipulation methods [1, 3, 46, 47] are able to manipulate the face image of person and create an indistinguishable fake image. Current face manipulation methods can be roughly divided into two categories, facial reenactment and identity swap. Facial reenactment tries to transfer the facial expres-

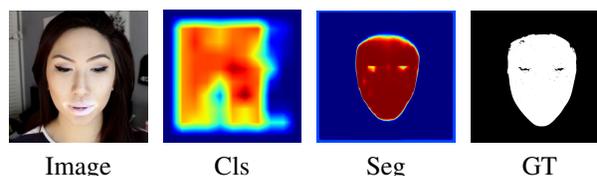


Figure 1: Predictions of a classification network and a segmentation network. The second image is the activation map of the classification network showing the high-response area. The third is the heatmap given by the segmentation network. Compared with the ground-truth on the right, the segmentation network localizes the tampered pixels on a far accurate level.

sions of one person to another person and synthesize realistic details. Face2Face [47] and NeuralTextures [46] are two representative works. Identity swap is a technique that enables replacing the face of a person with another person’s face. Deepfakes [1] and FaceSwap [3] are two of the most prominent methods. These methods enable effortless creation of fake face images and videos, which poses potential threats to our society. For example, fake news can be easily created by synthesizing a speech video of a politician [43].

To alleviate the potential issues caused by the fake face videos and images, great efforts have been dedicated to the field of face forensics, which aims to determine authenticity of a face photo. General image forensics techniques, relying on hand-crafted cues [18, 33, 8, 30], might not be suitable for face specific forensics tasks since faces are highly structured data. Recent works take advantage of great representation power of CNNs (Convolutional Neural Networks) and train a network using a large dataset containing authentic and manipulated face images [36, 4, 49, 9]. In [36], a large scale dataset called “FaceForensics++” is released to address the problem of face forensics. The dataset contains 5,000 videos generated from 4 popular face manipulation methods, Deefakes, FaceSwap, Face2Face and NeuralTextures, which provides rich data to train models as well as a standard benchmark for evaluation.

Most methods for face forensics cast the problem as a classification problem, in which given an image the model is expected to determine whether it is a real face or a manipulated face. Using deep networks has been proved effective in dealing with such a classification problem [36, 49]. In [36], a modified Xception network [11] is trained on “Face-Forensics++” dataset and achieves remarkable results, accuracy of 99.26 on the raw data. In [4], a compact network also achieves comparable performance. However, one question is raised: “*Is the problem well-defined?*” or “*Is it a good definition of the problem?*” In Figure 2, the second image shows the activation map of a classification model revealing the high response area for the fake face on the left. It is obviously that the activation map is not actually consistent with the ground-truth, which suggests that the features used to distinguish the fake images might have weak correlation to the real manipulated regions.

The example implies one of the limitations of a classification network that it can only produce a global scalar value representing the confidence of being fake but can not reflect the degree of how the image is manipulated. It would be more beneficial to have a pixel-level output that accurately reflects the manipulated pixels, as shown in the third image of Figure 2. Therefore, It would be more natural to formulate the problem of face forensics as a semantic segmentation task so that the model is forced to learn discriminative features to localize manipulated regions.

In this paper, we analyze the problem of face forensics from a pixel-level perspective using segmentation methods to complement the existing classification methods for face forensics. There are some questions that are still under investigated such as: 1) *By nature, whether face forensics is a classification or segmentation problem?* 2) *What is the most suitable network architecture for this problem?* 3) *Should we adopt shallow or deep networks?* 4) *Should we train the model from scratch or initialize it using general vision features.* We conduct experiments to try to answer these questions. By evaluating various architectures, we compare the performance of the segmentation networks and their counterpart classification networks from different aspects. We hope to provide more insight to the problem and establish a new baseline for the benchmark.

Our contributions are three folds:

- We conduct a pixel-level analysis to the problem of face forensics by using segmentation methods to be complementary to the existing classification methods.
- By redefining the problem to be a pixel-level task, we evaluate various architectures and create a strong new baseline for the problem.
- By performing different ablation studies, we analyze what makes an effective and efficient anti-fake model,

which, we hope, can shed some light on the field of research.

## 2. Related Work

We cover the most important related papers in the following paragraphs.

### 2.1. Digital Face Manipulation

A comprehensive state-of-the-art report of digital face manipulation can be found in [52]. Current facial manipulation methods can be separated into four categories: image-based approach, Audio-based approach, computer-graphics-based approach as well as learning based approach.

State of the Arts image-based approaches such as Video Rewrite [10], Video Face Replacement [5], Bringing Portraits to Life [7] and Deep Video Portraits [24]. These methods employ 2D warps to deform the image to match the expressions of a source actor. Synthesizing Obama [42] learned the mapping between audio and lip motions.

State-of-the-arts computer-graphics-based approaches such as Video Face Replacement [16], VDub [19] and Face2Face [47]. These methods usually reconstruct 3D models using blendshapes or other mesh editing process, based on high-quality 3D face capturing techniques as well as precise and rapid tracking techniques.

Recently, generative adversarial networks (GANs) are used to apply different facial attributes such as Aging [6], viewpoints [21], skin color [28], smiling [48], or other essential computer graphic renderings [24], which are implemented as an image-to-image translation, applying a patch-based GAN-loss.

### 2.2. Face forensics

Face forensics aims to ensure authenticity and origin of the face. Face forensics identify computer generated characters from computer graphics faces [4], print-scanned morphed faces [34], face splicing [15, 22], face swapping [51, 4], and face reenactment [4, 17]. Specific artifacts arising from the synthesis process such as color, texture [15] or eye blinking [26] can also be exploited. Learning-based approach propose a deep network trained to capture the subtle inconsistencies arising from low-level and/or high level features[4, 51]. Particularly, [20] uses a convolutional neural network to extract frame-level features, which are then used to train a recurrent neural network (RNN) that learns to classify if a video has been subject to manipulation or not. These approaches show impressive results, but can not precisely locate the manipulated area.

### 2.3. Pixel-level task

Instead of a rough prediction in global image-level view, there are many works towards to provide a local or pixel-

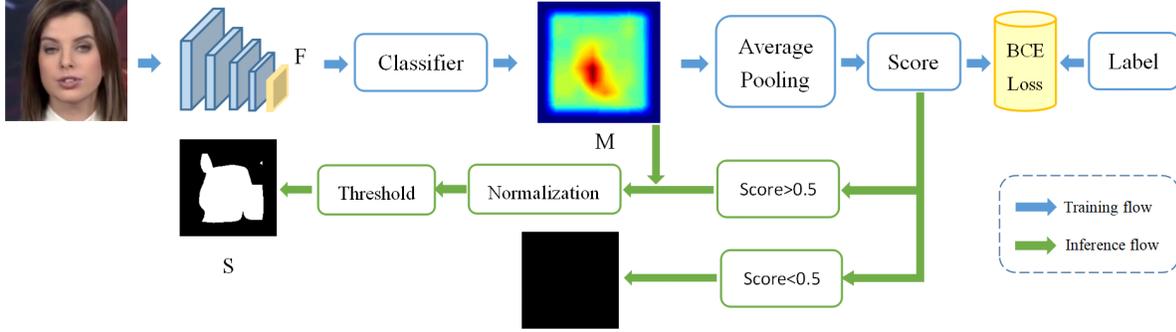


Figure 2: Pipeline of classification task. Different colors of arrows indicate different stages. Blue is for the training stage, Green is for the inference stage. When the classification score is above 0.5, it is classified as a fake image and is further processed to get the manipulated regions. When the score is below 0.5, indicating a real image, an all-zero mask is produced.

level prediction, such as Unet [35], fully convolutional network (FCN)[27], Deeplab.for semantic segmentation. As for image generation, pix2pix [23] realize the pixel-level transformation between different domains. There are lots of application concerning face parsing[29], pose parsing or scene segmentation.

As for face forensics, the mainstream methods are based on global classification at present, we drive the segmentation motivation of face manipulation to predict the region of local manipulation area. The face is often occluded by objects, but the face in the database [36] is generally unobstructed, so it can be trained directly.

### 3. Problem Setting

In this section, we first introduce the problem settings and methodologies for both the classification task and the segmentation task. Then we present an overview of the architectures used for evaluation.

#### 3.1. Classification Task

We first revisit the classification task. Formally Let  $\mathbf{x} \in \mathbb{R}^{H \times W \times 3}$  represent an image containing either an real or a tampered face, and  $l \in \{0, 1\}$  represent the label associated to it. We learn a mapping function  $f(\cdot) : \mathbb{R}^{H \times W \times 3} \rightarrow \{0, 1\}$  to predict the authenticity of a face image. Given a dataset  $\{(\mathbf{x}^t, l^t)\}_{t=1}^T$  containing T images, the network is trained by the following BCE (Binary Cross Entropy) loss:

$$\mathcal{L}_{cls} = -\frac{1}{T} \sum_{t=1}^T (1 - l^t) \log(p^t) + l^t \log(p^t) \quad (1)$$

where  $p^t$  is the output of the network for  $t$ th sample.

Since a classification network can only map an image to a scalar indicating the probability of an image being tampered, It is unclear whether the model has learned useful

features to localize the manipulated regions. There are some interpretation and visualization works trying to reveal more information from a classification network by investigating the activated regions on featuremaps. [13, 31, 50, 38] We adopt the most representative method, CAM (Class Activation Map), to help visualize what the model has learned.

CAM requires the network has an average pooling layer before the classifier, which collapses the output of the last convolution layer to a single vector. Suppose the featuremaps from the last convolution layer is  $\mathbf{F} \in \mathbb{R}^{H_f \times W_f \times K}$ ; the classifier has weight  $\mathbf{w} \in \mathbb{R}^{K \times 1}$ ; the activation map  $\mathbf{M} \in \mathbb{R}^{H_f \times W_f}$  of a tampered face is calculated as:

$$M_{ij} = \sum_k F_{ijk} \cdot w_k \quad (2)$$

where  $M_{ij}$ ,  $F_{ijk}$  and  $w_k$  are entries of  $\mathbf{M}$ ,  $\mathbf{F}$  and  $\mathbf{w}$  respectively.

What Equation 2 does is actually apply the classifier directly to the featuremaps  $\mathbf{F}$ , which performs classification on each spatial location. For simplicity, we modify the original CAM setting by switching the average pooling layer and the classifier. As shown in Figure 2, the activation map  $\mathbf{M}$  can be viewed as a dense prediction output for the image and the classification score is actually produced by averaging the activation map to a scalar.

In order to convert  $\mathbf{M}$  to a pixel-level mask, we need to further normalize it to the range of 0 to 1 and quantize it using a threshold. The normalization is operated as:

$$\tilde{\mathbf{M}} = \frac{\mathbf{M} - \min(\mathbf{M})}{\max[\mathbf{M} - \min(\mathbf{M})]} \quad (3)$$

The final pixel-level prediction is generated as:

$$\mathbf{S} = \mathbb{I}\{\tilde{\mathbf{M}} \geq \tau_1\} \quad (4)$$

where  $\mathbb{I}(\cdot)$  is a indicator function and  $\tau_1$  a threshold.

Now we have a pixel-level output that highlights the manipulated regions. Using these outputs makes it easier to investigate and analyze how well the classification model is able to learn discriminative and high-quality features on a pixel-level. Details and analysis are described in Section 4.



Figure 3: Pipeline of segmentation task. The network predicts a pixel-level output and is supervised directly by a pixel-level mask.

### 3.2. Segmentation Task

A classification network has limited capability to localize manipulated regions with a pixel-level manner because it is supervised only by a global label. Segmentation extends the task to a dense classification problem by assigning a label to each pixel of an image. The model is then forced to learn discriminative features to determine the authenticity of each pixel. Formally, the supervision for an image is defined as a mask  $y \in \{0, 1\}^{H \times W}$  instead of a single label and the loss is imposed on each pixel:

$$\mathcal{L}_{seg} = -\frac{1}{T} \sum_{t=1}^T \sum_{i=1}^H \sum_{j=1}^W (1 - y_{ij}^t) \log(p_{ij}^t) + y_{ij}^t \log(p_{ij}^t) \quad (5)$$

where  $y_{ij}^t$  and  $p_{ij}^t$  are the label and the prediction respectively for  $t$ th sample at position  $(i, j)$ .

Since a segmentation task requires pixel-level mask as supervision, annotation of the data is usually time-consuming. For example, as mentioned in [14], a high-resolution street view image for semantic segmentation requires around 1.5 hours for labelling. Fortunately for the face forensics task, the mask can be easily calculated by checking the pixel difference between the original image and the forged image without any extra annotation cost. Figure 4 shows some training images from “FaceForensics++” dataset as well as their corresponding mask indicating the manipulated area.

A classification network can be easily converted to a FCN (Fully Convolutional Network) [27] where the fully connected layers are replaced by convolutional layers. The pipeline for training a segmentation network is illustrated in Figure 3. Compared with the classification task in Figure 2, the main difference is that the average pooling is dropped and the BCE loss is directly applied to each pixel. The pixel-level prediction can be directly obtained from the trained model.

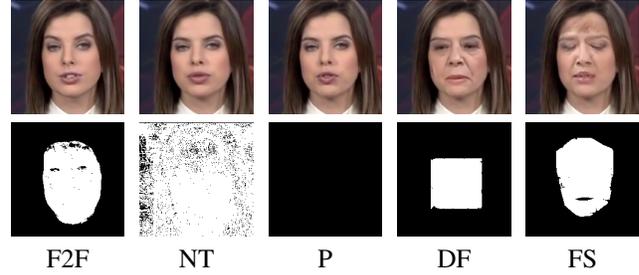


Figure 4: Illustration of example images and the corresponding masks for the “FaceForensics++” dataset. (P: Pristine, DF: DeepFakes, F2F: Face2Face, FS: FaceSwap, NT: NeuralTextures)

A segmentation model can be also evaluated from a global classification perspective by aggregating the dense prediction:

$$\hat{l} = \mathbb{I}\left\{\frac{1}{HW} \sum_{i=1}^H \sum_{j=1}^W \hat{y}_{ij} \geq \tau_2\right\} \quad (6)$$

where  $\hat{y}_{ij}$  represents the prediction at position  $(i, j)$  and  $\tau_2$  is the threshold.

In this way, we are able to make fair comparison between a segmentation network and its counterpart classification network under classification metrics. With extensive experiments in Section 4, we show the superiority of the segmentation networks for the face forensics task.

### 3.3. Architectures

In order to conduct deep analysis on the classification and segmentation task, we choose several representative architectures to evaluate the effectiveness on the problem of face forensics.

**Xception** [12] is a deep network architecture constructed by a series of modified inception modules [44] where the depthwise separable convolution is used. There are totally 36 convolutional layers involved to form the feature extraction base of the network. The architecture is adopted in [36] for the classification task of face forensics.

**MesoInception-4** [4] is a compact and light-weight network to address the problem of face forensics. It consists of two inception modules followed by two classic convolutional layers with maxpooling layers. We replace the all the operations after the last batchnorm layer with a single convolutional layer as the classifier.

**UNet** [35] is an effective and popular architecture for pixel-level tasks such as segmentation and pixel-to-pixel translation [23]. A Unet is basically defined by an encoder, consisting of convolutional layers and downsampling operations, and a decoder, consisting of convolutional layers

and upsampling operations. There are skip connections between the encoder and the decoder to enable passing information from low-level features. We choose two variants of UNet with different downsampling times in the encoder. **UNet8x** and **UNet4x** are downsampled by 8 and 4 times respectively.

**VGG16** [39] is a classic deep network for recognition tasks, which consists of 16 convolutional layers. Since we found the full network fails to converge on face forensics tasks, we only use two shallow versions **VGG8** and **VGG5**, containing the first 7 and 4 feature layers of vgg respectively and a classifier.

**FN3** is a 3-layer network we design to explore the potential of shallow networks. This architecture only contains two ‘‘Conv-BN-ReLU’’ blocks, and a  $3 \times 3$  convolutional layer as the classifier. The first two convolutional layers are with kernel size 7 and stride 2. It is interesting that this minimum structure works surprisingly good, even outperforming most of those deep architectures. Please refer to Section 4 for more details.

## 4. Experiments

### 4.1. Experiment Setup

	DF	F2F	FS	NT	P	Avg
Xception-clc	<b>99.16</b>	<b>98.35</b>	<b>98.88</b>	<b>99.09</b>	<b>99.18</b>	<b>98.93</b>
Mesonet-clc	93.33	77.01	26.77	92.05	88.99	75.63
UNet8x-clc	56.57	33.4	22.96	47.21	92.55	50.5
UNet4x-clc	66.9	45.45	37.48	55.42	98.98	60.8
VGG7-clc	41.69	73.37	67.78	38.81	76.45	59.6
VGG4-clc	56.02	84.92	90.72	40.33	70.66	68.53
Conv3-clc	94.35	93.28	81.13	94.26	61.43	84.89
Xception-seg	96.45	97.98	99.02	98.39	99.92	98.35
Mesonet-seg	68.86	79.58	89.77	96.92	59.56	78.94
UNet8x-seg	<b>99.08</b>	<b>98.74</b>	97.17	<b>99.42</b>	66.65	92.21
UNet4x-seg	98.61	97.32	<b>99.05</b>	96.53	97.01	97.70
VGG7-seg	98.41	98.34	99.05	99.01	99.33	98.83
VGG4-seg	98.24	98.29	99.03	99.01	99.99	98.91
Conv3-seg	98.16	98.32	99.03	99.06	<b>99.99</b>	<b>98.91</b>

Table 1: Classification accuracy on different manipulation methods. (P: Pristine, DF: DeepFakes, F2F: Face2Face, FS: FaceSwap, NT: NeuralTextures)

**Dataset:** FaceForensics++ [36] is a large scale face forensics dataset consisting of 5,000 video clips in total. Video sequences are crawled from the internet and a manual screening is adopted to ensure high quality and avoid face

occlusion, resulting in 1,000 original videos. Four manipulation methods, Deepfakes, Face2Face, FaceSwap and NeuralTextures, are applied to create forged videos, resulting in 4,000 fake clips. The dataset also provides data with three different compression levels, raw, HQ and LQ. We only focus on the raw quality task because low quality videos usually suffer from strong loss of visual and identity information, which might not cause abuse as those clear ones. [36] also suggests a split of 720 videos for training, 140 for validation as well as testing. We follow the same setting.

**Evaluation protocol and metrics:** In [36], there are two types of training protocols involved, *method specific training* and *mixed training*. The former involves forged data from only one of the manipulation methods. The latter requires training a model with all the real and forged data and the performance is evaluated on each specific method. We only adopt *mixed training* as it poses a more challenging task and real scenario. The evaluation is frame-based, therefore we extract all frames for the training set and partial frames for validation and testing (every 10 frames).

In terms of evaluation metrics, we use classification accuracy for the classification tasks, which represents how many test images are correctly classified. For segmentation tasks, IoU (Intersection over Union) is used to represent the ratio of  $\frac{TP}{TP+FP+TN}$ , where TP (True Positive), FP (False Positive) and TN (True Negative) are calculated based on pixels. The IoU is calculated for both foreground and background, denoted as Fg-IoU and Bg-IoU. The two IoUs are averaged to get mIoU, the mean IoU.

**Implementation details:** In face forensics, faces are the most important regions. As shown in [36], the model trained with the whole images performs poorly. Therefore, instead of using the whole image, we extract the faces as a pre-processing step using a public face detection tool [2] and only use the face regions to train the models. In order to include more background information, we enlarge the bounding box to the scale of 2. The segmentation masks are calculated by checking the difference between a manipulated face image and its corresponding original image. For segmentation tasks, the images are randomly cropped by size 256x256 and the same operation is applied to the corresponding mask to get the cropped mask. For classification tasks, it is necessary to include most face regions in the crop. Therefore, the shorter dimension of the image is first resized to 256, then a patch of 256x256 is cropped from the resized image.

The implementation is based on PyTorch [32]. All the models are trained using the Adam [25] optimizer with parameters  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . Since the Adam optimizer can adjust the learning rate dynamically, we only set the default learning rate to  $10^{-3}$  and do not use any learning rate decay policies. The batchsize is set to 64.

	mIoU						Bg-IoU						Fg-IoU					
	DF	F2F	FS	NT	P	Avg	DF	F2F	FS	NT	P	Avg	DF	F2F	FS	NT	P	Avg
Xception-seg	89.32	88.18	87.7	62.81	<b>99.95</b>	85.59	95.95	93.79	94.19	41.94	<b>99.95</b>	85.16	82.7	82.56	81.21	83.67	-	82.54
Mesonet-seg	56.58	51.14	54.52	40.23	90.2	58.53	78.96	71.06	74.68	22.02	90.2	67.38	34.19	31.21	34.35	58.44	-	39.55
UNet8x-seg	87.83	86.97	85.02	50.51	86.02	79.27	94.7	92.32	91.82	28.27	86.02	78.63	80.96	81.62	78.22	72.75	-	78.39
UNet4x-seg	89.12	89.43	86.29	51.46	96.09	82.48	95.41	93.89	92.59	30.68	96.09	81.73	82.82	84.96	79.99	72.25	-	80.00
VGG8-seg	94.68	95.21	94.33	<b>76.04</b>	99.31	91.91	97.87	97.34	97.19	<b>59.42</b>	99.31	90.23	91.48	93.07	91.47	<b>92.67</b>	-	92.17
VGG5-seg	<b>95.78</b>	<b>96.21</b>	<b>94.81</b>	75.6	99.86	<b>92.45</b>	<b>98.36</b>	<b>97.94</b>	<b>97.51</b>	58.97	99.86	<b>90.53</b>	<b>93.21</b>	<b>94.48</b>	<b>92.11</b>	92.23	-	<b>93.01</b>
FN3-seg	92.68	93.05	89.01	64.42	99.72	87.78	97.16	96.24	94.81	43.89	99.72	86.36	88.21	89.86	83.21	84.95	-	86.56
Xception-clc	<b>47.6</b>	<b>58.98</b>	<b>56.21</b>	<b>58.83</b>	<b>99.23</b>	<b>64.17</b>	59.9	<b>71.9</b>	<b>75.62</b>	<b>23.27</b>	<b>99.23</b>	<b>65.98</b>	<b>35.3</b>	<b>46.06</b>	<b>36.8</b>	52.95	-	<b>42.78</b>
Mesonet-clc	45.96	37.14	37.48	24.78	95.46	48.16	<b>67.87</b>	55.75	65.08	13.93	95.46	59.62	24.05	18.53	9.88	35.63	-	22.02
UNet8x-clc	23	33.6	34.82	29.71	86.39	41.5	28.63	49.62	53.8	11.3	86.39	45.94	17.42	17.58	15.84	48.11	-	24.74
UNet4x-clc	22.3	32.95	34.38	35.14	97.59	44.47	26.25	46.11	51.92	13.99	97.59	47.17	18.35	19.79	16.83	56.29	-	27.82
VGG8-clc	28.73	23.45	26.12	29.84	63.66	34.36	40.91	21.65	30.72	12.1	63.66	33.81	16.54	25.25	21.51	47.57	-	27.72
VGG5-clc	39.18	37.92	38.85	15.73	80.39	42.41	66.81	63.83	69.56	13.51	80.39	58.82	11.56	12.01	8.14	17.96	-	12.42
FN3-clc	16.77	18.46	20.47	43.84	48.68	29.64	14.58	10.63	17.9	8.09	48.68	19.97	18.97	26.29	23.04	<b>79.59</b>	-	36.97

Table 2: Segmentation results for different architectures. (P: Pristine, DF: DeepFakes, F2F: Face2Face, FS: FaceSwap, NT: NeuralTextures)

## 4.2. Experimental Results

**Classification task:** Table 1 shows the classification accuracy of different architectures. The suffixes “-seg” and “-cls” represent a segmentation model and a classification model respectively. The pixel-wise output is aggregated to a global output according to Equation 6. From the scores of the classification models, Xception-clc reaches the best performance, which is consistent with [36]. It can be seen that UNet, as a popular segmentation model for various pixel-level prediction tasks, fails to perform well in the classification task. It is interesting to see that FN3-clc, a minimum structure with only 3 layer works surprisingly good. Although, the performance is lower than Xception, FN3-clc achieves far better performance than other models. For those segmentation models, it can be easily noticed that they obtain better classification results than their counterpart classification models, which shows the benefit of training models under pixel-level supervision.

**Segmentation task:** Table 2 shows the segmentation results of different architectures. The classification models are trained using a global image-level label and visualized by the CAM to get a pixel-level output, explained in Section 3.1. For segmentation models, VGG5-seg achieves the best performance in terms of both mIoU, Bg-IoU, and Fg-IoU. Mesonet-seg, as a compact and efficient architecture, does not achieve comparable results, outperformed by other methods by a large margin. We suspect it could be due to the limited capacity of the model. It is also worth noting that UNet still does not reach promising results as a popular segmentation architecture. On the contrary, the 3-

layer network FN3-seg shows better potential, even better than Xception-seg. For classification models, Xception-clc achieves the best results in most of the scores, which implies that Xception-clc can successfully learn high-quality features to locate manipulated regions even trained with a global image-level label. However, Xception-clc can be hardly compared with its segmentation counterpart that obtains far higher scores. The rest of the classification models all suffer from low scores. Even VGG5-clc, whose segmentation counterpart achieves the best results, is unable to produce plausible predictions without pixel-level supervision.

From the results above, obviously segmentation models show superiority over the classification models in terms of both pixel-level prediction and global-level prediction. Figure 5 shows outputs of different architectures, which further illustrates the benefit of analyzing fake faces on a pixel-level.

## 4.3. Analysis

### Deep vs Shallow

To explore the effect of model depth to the task of face forensics, we also take a closer look at the performance of models with different depth. In Table 3, we summarize the mIoUs of segmentation models with different depth. Apart from VGG8 and VGG5, we also include VGG3, which only uses the first two layers of VGG16 followed by a classifier. It is interesting to see that the deep model, Xception with 36 layers, does not reach to a high score, whereas the shallow models present better abilities. This reveals that face forensics is supposed to be defined as a low-level vision problem than a high-level perception problem.

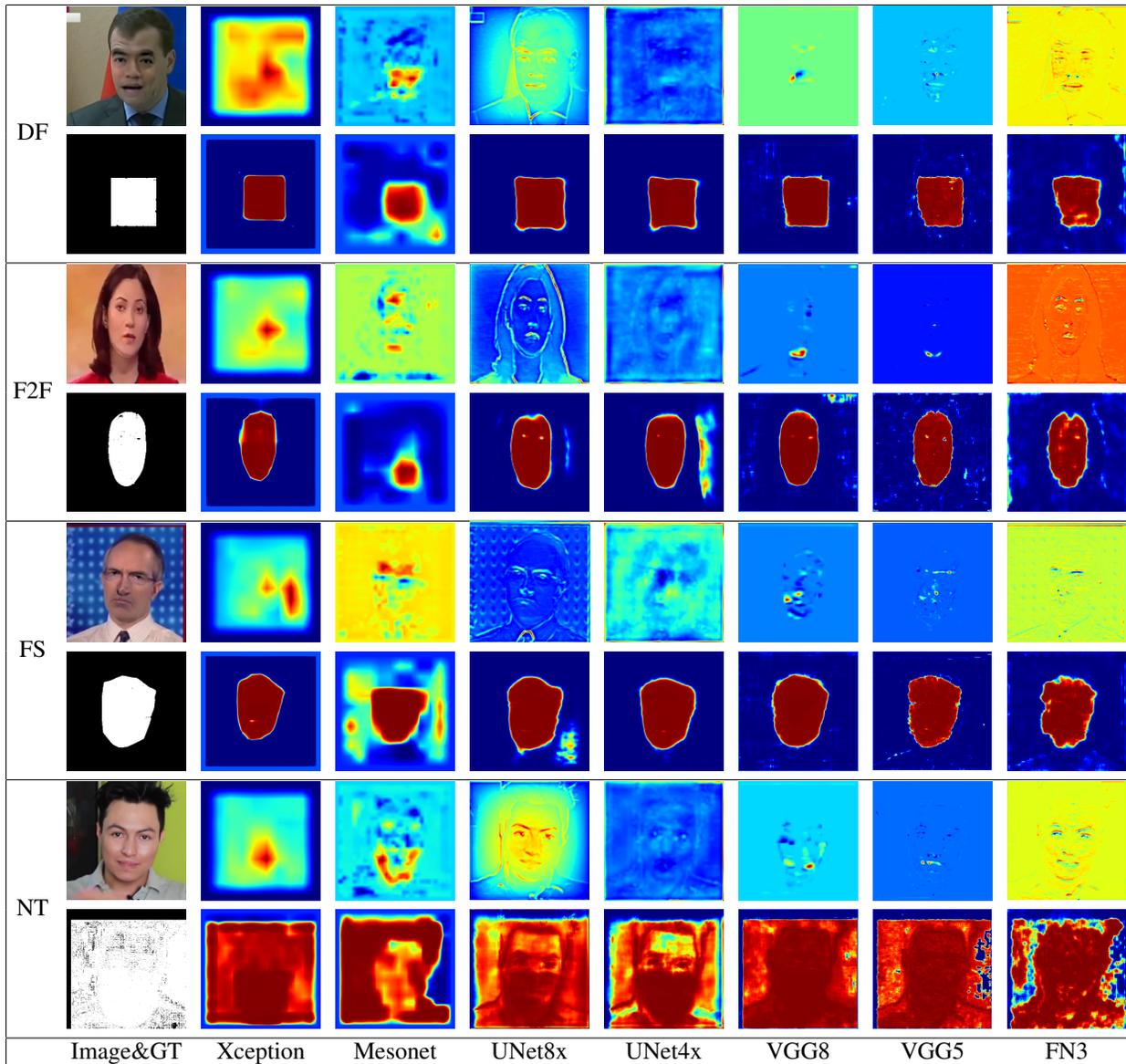


Figure 5: Qualitative results of the classification and segmentation models. Each of two rows relates to a specific manipulation method. For each method, on the left are the input fake image and the ground-truth indicating the manipulated area. The upper row shows the pixel-level results of all the classification models, and the lower row displays the predictions of the segmentation models. (DF: DeepFakes, F2F: Face2Face, FS: FaceSwap, NT: NeuralTextures)

### Pretrained or From Scratch

As implied by the analysis in the last section, face forensics is more like a low-level vision task. Another question is that “can the models benefit from the features used for general vision recognition tasks?” We conduct another ablation study where we compare the performance on the segmentation task using models with and without ImageNet-pretraining. The results are shown in Table 4. According to the numbers, there is little difference between the pre-trained model and the trained-from-scratch model. The fea-

tures learned in a general vision recognition task such as ImageNet did not help quickly find a better local optima.

### Kernel visualization

In order to have a better understanding of the features learned by the model, we analyze the kernels by visualizing them using the technique in [40]. In Figure 6, for each fake image, we visualize two kernels in each convolutional layer. Apart from the features in conv1, which are mostly low-level edges and corners, the kernels in following layers do not make much sense to us humans. Intuitively, the

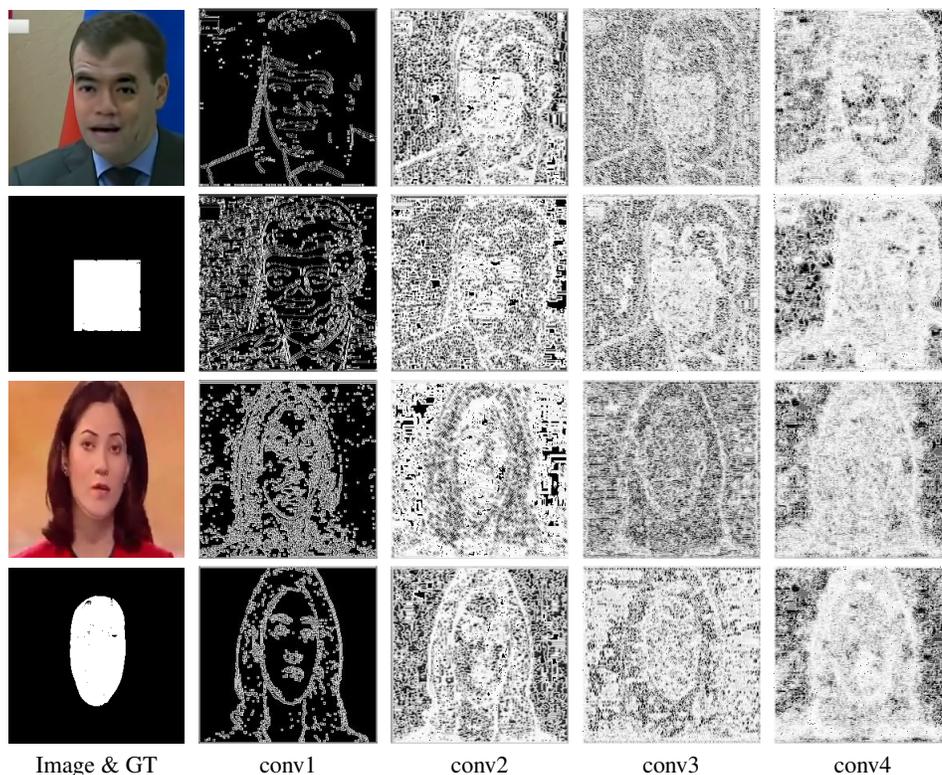


Figure 6: Kernel visualization of VGG5. The left column is the input fake image and the ground-truth. Each column on the right shows kernels from a specific convolutional layer.

	DF	F2F	FS	NT	P	Avg
Xception (36)	89.32	88.18	87.7	62.81	<b>99.95</b>	85.59
VGG8 (7)	94.68	95.21	94.33	<b>76.04</b>	99.31	91.91
VGG5 (4)	<b>95.78</b>	<b>96.21</b>	<b>94.81</b>	75.6	99.86	<b>92.45</b>
FN3 (3)	92.68	93.05	89.01	64.42	99.72	87.78
VGG3 (3)	88.79	89.92	79.65	57.93	96.64	82.58

Table 3: Comparison among models with different depth. The number in the parentheses indicates the depth of the model. The numbers are mIoU.

model tries to learn subtle features, to which humans are not sensitive to. Humans are good at recognizing things on a semantic level, but fake faces, generated by advanced manipulation methods, seem beyond humans ability. This further emphasizes the demand of a good face forensics model.

## 5. Conclusion

Face forensics has become increasingly important as face manipulation methods have made stunning progress to enable effortless generation of indistinguishable fake face images. Most previous works cast the problem as a clas-

	DF	F2F	FS	NT	P	Avg
Xception (pretrained)	89.32	88.18	87.7	62.81	99.95	<b>85.59</b>
Xception (non-pretrained)	88.72	87.88	88.70	62.84	99.74	85.57
VGG5 (pretrained)	95.78	96.21	94.81	75.6	99.86	<b>92.45</b>
VGG5 (non-pretrained)	95.69	96.2	94.75	75.35	99.86	92.37
VGG8 (pretrained)	94.68	95.21	94.33	76.04	99.31	91.91
VGG8 (non-pretrained)	95.67	95.93	95.06	75.18	99.83	<b>92.33</b>

Table 4: Comparison between a pretrained model and the model trained from scratch. The numbers are mIoU.

sification task, which suffers from limitations. In this paper, we analyze the problem from pixel-level perspective by using segmentation methods to complement the traditional classification methods. With comprehensive experiments, we show the superiority of formulating it as a segmentation problem instead of a classification problem. In addition, we also perform different ablation studies to analyze the important factors of being an effective face forensics model, which establishes a strong new baseline for the benchmark. We hope that our analysis can provide more insight to the field of face forensics.

## References

- [1] Deepfakes [github](https://github.com/deepfakes/faceswap). <https://github.com/deepfakes/faceswap>.
- [2] facerecognition [github](https://github.com/ageitgey/face_recognition). [https://github.com/ageitgey/face\\_recognition](https://github.com/ageitgey/face_recognition).
- [3] Faceswap. <https://github.com/MarekKowalski/FaceSwap>.
- [4] D. Afchar, V. Nozick, J. Yamagishi, and I. Echizen. Mesonet: a compact facial video forgery detection network. In *2018 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–7, Dec 2018.
- [5] Mahmoud Afifi, Khaled Hussain, Hosny Ibrahim, and Nagwa Omar. Video face replacement system using a modified poisson blending technique. *2014 International Symposium on Intelligent Signal Processing and Communication Systems, ISPACS 2014*, 12 2014.
- [6] G. Antipov, M. Baccouche, and J. Dugelay. Face aging with conditional generative adversarial networks. In *2017 IEEE International Conference on Image Processing (ICIP)*, pages 2089–2093, Sep. 2017.
- [7] Hadar Averbuch-Elor, Daniel Cohen-Or, Johannes Kopf, and Michael F. Cohen. Bringing portraits to life. *ACM Transactions on Graphics (Proceeding of SIGGRAPH Asia 2017)*, 36(6):196, 2017.
- [8] Mauro Barni, Luca Bondi, Nicol Bonettini, Paolo Bestagini, Andrea Costanzo, Marco Maggini, Benedetta Tondi, and Stefano Tubaro. Aligned and non-aligned double jpeg detection using convolutional neural networks. *Journal of Visual Communication and Image Representation*, 49, 08 2017.
- [9] Belhassen Bayar and Matthew C. Stamm. A deep learning approach to universal image manipulation detection using a new convolutional layer. In *Proceedings of the 4th ACM Workshop on Information Hiding and Multimedia Security, IH&#38;MMSec '16*, pages 5–10, New York, NY, USA, 2016. ACM.
- [10] Christoph Bregler, Michele Covell, and Malcolm Slaney. Video rewrite: Driving visual speech with audio. In *Proceedings of the 24th Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH '97*, pages 353–360, New York, NY, USA, 1997. ACM Press/Addison-Wesley Publishing Co.
- [11] Francois Chollet. Xception: Deep learning with depthwise separable convolutions. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [12] Francois Chollet. Xception: Deep learning with depthwise separable convolutions. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [13] R. G. Cinbis, J. Verbeek, and C. Schmid. Weakly supervised object localization with multi-fold multiple instance learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(1):189–203, Jan 2017.
- [14] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [15] T. J. d. Carvalho, C. Riess, E. Angelopoulou, H. Pedrini, and A. d. R. Rocha. Exposing digital image forgeries by illumination color classification. *IEEE Transactions on Information Forensics and Security*, 8(7):1182–1194, July 2013.
- [16] Kevin Dale, Kalyan Sunkavalli, Micah K. Johnson, Daniel Vlastic, Wojciech Matusik, and Hanspeter Pfister. Video face replacement. In *Proceedings of the 2011 SIGGRAPH Asia Conference, SA '11*, pages 130:1–130:10, New York, NY, USA, 2011. ACM.
- [17] D. Dang-Nguyen, G. Boato, and F. G. B. De Natale. Identify computer generated characters by analysing facial expressions variation. In *2012 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 252–257, Dec 2012.
- [18] Hany Farid. *Photo Forensics*. The MIT Press, 2016.
- [19] P. Garrido, Levi Valgaerts, Hamid Sarmadi, I. Steiner, Kiran Varanasi, P. Prez, and C. Theobalt. Vdub: Modifying face video of actors for plausible visual alignment to a dubbed audio track. *Computer Graphics Forum*, 34, 05 2015.
- [20] D. Gera and E. J. Delp. Deepfake video detection using recurrent neural networks. In *2018 15th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, pages 1–6, Nov 2018.
- [21] Rui Huang, Shu Zhang, Tianyu Li, and Ran He. Beyond face rotation: Global and local perception gan for photorealistic and identity preserving frontal view synthesis. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [22] Minyoung Huh, Andrew Liu, Andrew Owens, and Alexei A. Efros. Fighting fake news: Image splice detection via learned self-consistency. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- [23] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. *arxiv*, 2016.
- [24] H. Kim, P. Garrido, A. Tewari, W. Xu, J. Thies, N. Nießner, P. Pérez, C. Richardt, M. Zollhöfer, and C. Theobalt. Deep Video Portraits. *ACM Transactions on Graphics 2018 (TOG)*, 2018.
- [25] Diederick P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.
- [26] Y. Li, M. Chang, and S. Lyu. In icu oculi: Exposing ai created fake videos by detecting eye blinking. In *2018 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–7, Dec 2018.
- [27] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3431–3440, June 2015.
- [28] Yongyi Lu, Yu-Wing Tai, and Chi-Keung Tang. Attribute-guided face generation using conditional cyclegan. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- [29] Ping Luo. Hierarchical face parsing via deep learning. In *Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, CVPR '12, pages

- 2480–2487, Washington, DC, USA, 2012. IEEE Computer Society.
- [30] James F. O’Brien and Hany Farid. Exposing photo manipulation with inconsistent reflections. *ACM Transactions on Graphics*, 31(1):4:1–11, Jan. 2012. Presented at SIGGRAPH 2012.
- [31] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1717–1724, June 2014.
- [32] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017.
- [33] A.C. Popescu and H. Farid. Exposing digital forgeries by detecting traces of resampling. *Trans. Sig. Proc.*, 53(2):758–767, Feb. 2005.
- [34] R. Raghavendra, K. B. Raja, S. Venkatesh, and C. Busch. Transferable deep-cnn features for detecting digital and print-scanned morphed face images. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 1822–1830, July 2017.
- [35] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, volume 9351 of *LNCS*, pages 234–241. Springer, 2015. (available on arXiv:1505.04597 [cs.CV]).
- [36] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Niessner. FaceForensics++: Learning to detect manipulated facial images. In *International Conference on Computer Vision (ICCV)*, 2019.
- [37] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, 2015.
- [38] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [39] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, May 2015.
- [40] J.T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller. Striving for simplicity: The all convolutional net. In *ICLR (workshop track)*, 2015.
- [41] Yi Sun, Yuheng Chen, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation by joint identification-verification. In *Advances in Neural Information Processing Systems*, pages 1988–1996, 2014.
- [42] Supasorn Suwajanakorn, Steven Seitz, and Ira Kemelmacher. Synthesizing obama: learning lip sync from audio. *ACM Transactions on Graphics*, 36:1–13, 07 2017.
- [43] Supasorn Suwajanakorn, Steven M. Seitz, and Ira Kemelmacher-Shlizerman. Synthesizing obama: Learning lip sync from audio. *ACM Trans. Graph.*, 36(4):95:1–95:13, July 2017.
- [44] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- [45] Yaniv Taigman, Ming Yang, Marc’Aurelio Ranzato, and Lior Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1701–1708, 2013.
- [46] Justus Thies, Michael Zollhöfer, and Matthias Niessner. Deferred neural rendering: Image synthesis using neural textures. *ACM Trans. Graph.*, 38(4):66:1–66:12, July 2019.
- [47] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Niessner. Face2face: Real-time face capture and reenactment of rgb videos. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [48] Paul Upchurch, Jacob Gardner, Geoff Pleiss, Robert Pless, Noah Snavely, Kavita Bala, and Kilian Q. Weinberger. Deep feature interpolation for image content changes. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [49] Sheng-Yu Wang, Oliver Wang, Andrew Owens, Richard Zhang, and Alexei A Efros. Detecting photoshopped faces by scripting photoshop. *arXiv preprint arXiv:1906.05856*, 2019.
- [50] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [51] P. Zhou, X. Han, V. I. Morariu, and L. S. Davis. Two-stream neural networks for tampered face detection. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 1831–1839, July 2017.
- [52] Michael Zollhofer, Justus Thies, Pablo Garrido, Derek Bradley, Thabo Beeler, Patrick Prez, Marc Stamminger, Matthias Niener, and Christian Theobalt. State of the Art on Monocular 3D Face Reconstruction, Tracking, and Applications. *Computer Graphics Forum*, 2018.