

This CVPR 2021 paper is the Open Access version, provided by the Computer Vision Foundation. Except for this watermark, it is identical to the accepted version; the final published version of the proceedings is available on IEEE Xplore.

# ForgeryNet: A Versatile Benchmark for Comprehensive Forgery Analysis

Yinan He<sup>1,2\*</sup> Bei Gan<sup>1,3\*</sup> Siyu Chen<sup>1,3\*</sup> Yichun Zhou<sup>1,4\*</sup> Guojun Yin<sup>1,3</sup> Luchuan Song<sup>5†</sup> Lu Sheng<sup>4</sup> Jing Shao<sup>1,3‡</sup> Ziwei Liu<sup>6</sup> <sup>1</sup>SenseTime Research <sup>2</sup>Beijing University of Posts and Telecommunications <sup>3</sup>Shanghai AI Laboratory <sup>4</sup>College of Software, Beihang University <sup>5</sup>University of Science and Technology of China <sup>6</sup>S-Lab, Nanyang Technological University {heyinan, ganbei, chensiyu, yinguojun, shaojing}@sensetime.com

{buaazyc, lsheng}@buaa.edu.cn slc0826@mail.ustc.edu.cn ziwei.liu@ntu.edu.sg



Figure 1: ForgeryNet is a new mega-scale face forgery dataset with comprehensive annotations and four forgery analysis tasks. It contains thousands of subjects, various manipulation methods and diverse re-rendering processes. In (a), can you distinguish which images are forged?

### Abstract

The rapid progress of photorealistic synthesis techniques have reached at a critical point where the boundary between real and manipulated images starts to blur. Thus, benchmarking and advancing digital forgery analysis have become a pressing issue. However, existing face forgery datasets either have limited diversity or only support coarse-grained analysis.

To counter this emerging threat, we construct the ForgeryNet dataset, an extremely large face forgery dataset

with unified annotations in image- and video-level data across four tasks: 1) Image Forgery Classification, including two-way (real / fake), three-way (real / fake with identity-replaced forgery approaches / fake with identityremained forgery approaches), and n-way (real and 15) respective forgery approaches) classification. 2) Spatial Forgery Localization, which segments the manipulated area of fake images compared to their corresponding real images. 3) Video Forgery Classification, which re-defines the video-level forgery classification with manipulated frames in random positions. This task is important because attackers in real world are free to manipulate any target frame. and 4) Temporal Forgery Localization, to localize the temporal segments which are manipulated. ForgeryNet is by far the largest publicly available deep face forgery dataset in terms of data-scale (2.9 million

<sup>\*</sup>Equal contribution.

<sup>&</sup>lt;sup>†</sup>Work done during an internship at SenseTime Research.

<sup>&</sup>lt;sup>‡</sup>Corresponding author.

<sup>\$</sup>https://yinanhe.github.io/projects/forgerynet. html

The label of images in Fig. 1(s)(from left to r-tght); fake, take, fake, take, fake, take, fake, take, fake, take, fake, fake

images, 221,247 videos), manipulations (7 image-level approaches, 8 video-level approaches), perturbations (36 independent and more mixed perturbations) and annotations (6.3 million classification labels, 2.9 million manipulated area annotations and 221,247 temporal forgery segment labels). We perform extensive benchmarking and studies of existing face forensics methods and obtain several valuable observations. We hope that the scale, quality, and variety of our ForgeryNet dataset will foster further research and innovation in the area of face forgery classification, as well as spatial and temporal forgery localization etc.

# 1. Introduction

Photorealistic facial forgery technologies, especially recent deep learning driven approaches [17, 26, 35], give rise to widespread social concerns on potential malicious abuse of these techniques to eye-cheatingly forge media (*i.e.*, images and videos, *etc.*) of human faces. Therefore, it is of vital importance to develop reliable methods for face forgery analysis<sup>1</sup>, so as to distinguish *whether* and *where* an image or video is manipulated.

Most recent progress about face forgery analysis are sparked by gathering of face forgery detection datasets [12, 38] and early attempts of profiling intrinsic characteristics within the forgery images. However, performances on most datasets have already saturated (*i.e.* over 99% accuracy [19,23,33,44]) due to their limited scales (*e.g.* number of images/videos and subject identities) and limited diversity (*e.g.* forgery approaches, scenarios, realistic perturbations, *etc.*). Moreover, in practical applications, it is often required to detect forged faces by locating tampered areas in an image and/or manipulated segments in an untrimmed video, rather than merely providing a binary label.

In this paper, we construct a new mega-scale dataset named ForgeryNet with comprehensive annotations, consisting of two groups (*i.e.* image- and video-level) and four tasks for real-world digital forgery analysis. We carefully benchmark existing forensics methods on ForgeryNet. Extensive experiments and in-depth analysis show that this larger and richer annotated dataset can boost the development of next-generation algorithms for forgery analysis. Specifically, ForgeryNet brings several unique advantages over existing datasets.

(1) Wild Original Data. Most current datasets are captured under controlled conditions (*e.g.* environment, angles and lighting). We collect original data with diversified dimensions of angle, expression, identity, lighting, scenario and

*etc.* from four datasets [6, 10, 14, 32]. Note that all the original data have a *Creative Commons Attribution* license that allows to share and adapt the material.

(2) Various Forgery Approaches. There are at most 8 forgery approaches in all current datasets, while ForgeryNet is manipulated by 15 approaches, including face transfer, face swap, face reenactment and face editing. We choose approaches that span a variety of learning-based models, including encoder-decoder structure, generative adversarial network, graphics formation and RNN/LSTM (Fig. 4).

(3) Diverse Re-rendering Process. In the process of transmission and re-rendering, media data (image/video) always undergo compression, blurring and other operations, which may smooth the traces of forgery and bring more challenge for forgery detection. The ForgeryNet dataset posts 36 perturbations, such as optical distortion, multiplicative noise, random compression, blur, and *etc*. As shown in Fig. 1(c), circle sizes refer to the number of forgery approaches with re-rendering process operations.

(4) Rich Annotations and Comprehensive Tasks. According to the real application scenario, we propose four tasks, as shown in Fig. 1(b): 1) Image Forgery Classification, distinguishes whether an image is forgery or not and meanwhile tells its forgery type (*i.e.* manipulation approaches). We provide three types of annotations including two-way, three-way and n-way classification. Both intra- and crossforgery evaluations are set on three-way and *n*-way settings. 2) Spatial Forgery Localization, localizes manipulated areas of forgery images. Due to the fact that a forgery image may contain multiple faces and can be manipulated entirely or in part, it is more substantial to segment modified pixels in addition to only telling that it is forged. 3) Video Forgery Classification, similar to image-level classification, contains three types of annotations. Note that different from existing forgery video datasets, we construct our video dataset with untrimmed videos, each of which has part of the frames manipulated, considering the fact that forgery videos in real world are often manipulated on a certain subject and some key frames. 4) Temporal Forgery Localization, localizes the temporal segments which are manipulated. This is a new task for forgery analysis. Together with Video Forgery Classification and Spatial Forgery Localization, it provides comprehensive spatio-temporal forgery annotations.

# 2. Related Works

Due to the urgency in detecting face manipulation, many efforts have been devoted to creating face forgery detection datasets. Previous datasets can be grouped down into three generations. Their statistical information is listed in Tab. 1. **The first generation** consists of datasets such as DF-TIMIT [25], UADFV [43], SwapMe and FaceSwap [47]. DF-TIMIT manually selects 16 pairs of appearance-similar people from the publicly available VidTIMIT database, and

<sup>&</sup>lt;sup>1</sup>In this paper, the definition of the term "face forgery" refers to an image or a video containing modified identity, expressions or attribute(s) with a learning-based approach, distinguished with 1) a so-called "Cheap-Fakes" [34] that are created with off-the-shelf softwares without learnable components and 2) "DeepFakes" that only refer to manipulations with swapped identities [12].

Table 1: Comparison of various face forgery datasets. ForgeryNet surpasses any other dataset both in scale and diversity. It provides both video- and image-level data. The forgery data are constructed by 15 manipulation approaches within 4 categories. We also employ 36 types of perturbations from 4 kinds of distortions for post-processing.

| Datasat                  | Video  | o Clips | Still i   | mages     | Ammaaahaa  | Subjects | Uniq.    | Mix          | Annotationa |
|--------------------------|--------|---------|-----------|-----------|------------|----------|----------|--------------|-------------|
| Dataset                  | Real   | Fake    | Real      | Fake      | Approaches | Subjects | Perturb. | Perturb.     | Annotations |
| UADFV [43]               | 49     | 49      | 241       | 252       | 1          | 49       | -        | ×            | 591         |
| DF-TIMIT [25]            | 320    | 640     | -         | -         | 2          | 43       | -        | ×            | 1,600       |
| Deep Fake Detection [4]  | 363    | 3,068   | -         | -         | 5          | 28       | -        | ×            | 3,431       |
| Celeb-DF [27]            | 590    | 5,639   | -         | -         | 1          | 59       | -        | ×            | 6,229       |
| SwapMe and FaceSwap [47] | -      | -       | 4,600     | 2,010     | 2          | -        | -        | ×            | 6,610       |
| DFFD [11]                | 1,000  | 3,000   | 58,703    | 240,336   | 7          | -        | -        | ×            | 8,000       |
| FaceForensics++ [38]     | 1,000  | 5,000   | -         | -         | 5          | -        | 2        | ×            | 11,000      |
| DeeperForensics-1.0 [24] | 50,000 | 10,000  | -         | -         | 1          | 100      | 7        | $\checkmark$ | 60,000      |
| DFDC [12]                | 23,564 | 104,500 | -         | -         | 8          | 960      | 19       | ×            | 128,064     |
| ForgeryNet (Ours)        | 99,630 | 121,617 | 1,438,201 | 1,457,861 | 15         | 5400+    | 36       | $\checkmark$ | 9,393,574   |



Figure 2: Representative examples of original data collected from four face datasets respectively.

generates 640 videos with faces swapped. UADFV contains 98 videos, *i.e.* 49 real videos from YouTube and 49 fake ones generated by FakeAPP [3]. SwapMe and FaceSwap choose two face swapping Apps [1,2] to create 2010 forgery images in total on 1005 original real images.

The second generation includes Google DeepFake Detection dataset [4] with 3,068 forgery videos by five publicly available manipulation approaches, and Celeb-DF [27] containing 590 YouTube real videos mostly from celebrities and 5,639 manipulated video clips. FaceForensics++ [38] consists of 4000 fake videos manipulated by four approaches (*i.e.* DeepFakes, Face2Face, FaceSwap and NeuralTextures), and 1000 real videos from YouTube. The data scale and quality of the second generation have been improved. However, these datasets still lack diversity in forgery approaches and task annotations, and are not wellsuited for challenges encountered in real world.

The third generation datasets are the most recent face forgery datasets, *i.e.* DeeperForensics-1.0 [24], DFDC [12], and DFFD [11] which contains tens of thousands of videos and tens of millions of frames. DeeperForensics-1.0 consists of 60,000 videos for real-world face forgery detection. DFDC contains over 100,000 clips sourced from 960 paid actors, produced with several face replacement forgery approaches including learnable and non-learnable approaches. In a practical application, in addition to classification, it is necessary to locate the manipulated areas or segments in an image or an untrimmed video. A few datasets have taken these tasks into consideration. DFFD provides annotations of spatial forgery at the first time, yet it only presents binary masks without manipulation density.



Figure 3: Sampled forgeries in our ForgeryNet. (a) Identityremained forgery approaches: 1) *Face reenactment*, 2) *Face editing*. (b) Identity-replaced forgery approaches: 1) *Face transfer*, 2) *Face swap*, 3) *Face stacked manipulation*.

# 3. ForgeryNet Construction

Most of existing public face forgery datasets [4, 11, 12, 24, 25, 25, 27, 38, 43, 47] contain only single or no more than 10 specific manipulation approaches, and even the largest one [12] only operates 8 manipulations with 19 perturbations on 960 subjects. Moreover, these datasets take forgery analysis solely as a classification task. On the contrary, our proposed ForgeryNet dataset provides 15 manipulation approaches with more than 36 mix-perturbations on over  $5,400^2$  subjects, and defines four tasks (*i.e.* image and video classification, spatial and temporal localization) with a total of 9.4M annotations. Our whole dataset consists of two subsets: Image-forgery set provides over 2.9M still images and Video-forgery set has more than 220k video clips. These two subsets have their real data respectively randomly selected from the original data, and 15 forgery approaches are applied to image-forgery construction while 8 of them also generate the video-forgery data<sup>3</sup>. We compare our ForgeryNet with other publicly available datasets in Tab. 1.

<sup>&</sup>lt;sup>2</sup>Some original datasets do not provide the identity annotation.

 $<sup>^{3}\</sup>mbox{There}$  are 7 forgery approaches that are only suitable for generating images.



Figure 4: Pipeline of face forgery approaches. (a)-(c) Representation preparation: target image  $I_t$ , conditional source  $x_s$  and their intermediate representations. (d) Forgery models produce a forged target face  $\tilde{I}_t^f$  by processing the representations. (e)-(f) Re-render  $\tilde{I}_t^f$  to full image  $I_t$  and get the forgery image  $\tilde{I}_t$ . (g) Apply perturbations to  $\tilde{I}_t$  to obtain final forgery data.

Over all the comparison items listed in the table, our dataset surpasses the rest both in scale and diversity.

# 3.1. Original Data Collection

**Source of Original Data.** Four face datasets, CREMA-D [6], RAVDESS [32], VoxCeleb2 [10] and AVSpeech [14], are chosen as the original data to boost the diversity in dimensions of face identity, angle, expression, scenarios *etc*.

Note that CREMA-D is made available under the Open Database License, while others are released under a Creative Commons Attribution License. The resolutions of these original data range from 240p to 1080p, and face yaw angles ranging from -90 to 90 degrees are all covered. Representative examples are shown in Fig. 2.

**Preprocess Original Data.** For further manipulation, we crop original videos into a controllable set of source videos with reasonable lengths. Then we detect and select faces for manipulation and obtain their face attribute labels.

### **3.2. Forgery Approach**

To guarantee the diversity of forgery approaches in the proposed ForgeryNet, we introduce 15 face forgery approaches They are selected according to perspectives of modeling types, conditional sources, forgery effects and functions. We denote  $x_t$  as the *target* subject to be manipulated while the *source*  $x_s$  is regarded as the conditional media driving the *target* to change either identity or attributes, or even both.

#### 3.2.1 Forgery Category

According to the visual effects of facial manipulation, we divide the forgery approaches into two categories, *i.e. Identity-remained* and *Identity-replaced*. Sampled forgeries in Fig. 3 illustrate these categories and their sub-types.

**Identity-remained Forgery Approach** in Fig. 3(a) remains the identity of  $x_t$  and the identity-agnostic content like expression, mouth, hair and pose of  $x_t$  are changed, driven

by  $x_s$ . We adopt eight approaches and divide them into two sub-types: 1) *Face reenactment* on  $x_t(i, a)$  preserves its identity but has its *intrinsic* attributes like pose, mouth and expression manipulated by conditional source  $x_s$  and forms  $x_t(i, \tilde{a}^s)$ , where *i* refers to identity and *a* denotes attribute(s). Alternatively, with 2) *Face editing* on  $x_t(i, a)$  has its *external* attributes altered, such as facial hair, age, gender and ethnicity, to obtain  $x_t(i, \hat{a}^s)$ . We also include multiple attribute manipulation with two editing approaches, *e.g.* both hair and eyebrow are manipulated as shown with the first example in Fig. 3(a-2).

**Identity-replaced Forgery Approach** in Fig. 3(b) replaces the content of  $x_t$  with that of  $x_s$  preserving the identity of s. Seven approaches are divided into three sub-types as follows. 1) *Face transfer* transfers both identity-aware and identity-agnostic content (*e.g.* expression and pose) from  $x_s$ to  $x_t$ , resulting in  $x_t(\tilde{i}^s, \tilde{a}^s)$ . 2) *Face swap* which produces  $x_t(\tilde{i}^s, a)$  only swaps identity from the source  $x_s$  to the target  $x_t$ , and the identity-agnostic content a are preserved. 3) *Face stacked manipulation* refers to a combination of both *Identity-remained* and *Identity-replaced* approaches. We propose two assembles<sup>4</sup>, *i.e.*  $\langle editing \rightarrow transfer \rangle$  and  $\langle swap \rightarrow editing \rangle$ , where the former one transfers both the identity and attributes of the manipulated  $x_s(\hat{i}, \hat{a})$  to the target  $x_t$  to obtain  $x_t(\tilde{i}^s, \tilde{a}^s)$  and the latter alters the external attributes of the swapped target  $x_t(\tilde{i}^s, a)$  to get  $x_t(\tilde{i}^s, \hat{a}^s)$ .

### 3.2.2 Forgery Pipeline

Although there are a wild variety of architectures designed for the aforementioned approaches, most are created using variations or combinations of generative networks, encoderdecoder networks or graphics formation. We briefly summarize the forgery pipeline in Fig. 4.

The target is always an image marked as  $I_t$ , while there are various conditional source formats  $x_s$ , including image,

 $<sup>^{4}\</sup>mbox{StarGAN2-BlendFace-Stack}$  (SBS), DeepFakes-StarGAN2-Stack (DSS)



Figure 5: Annotations for Spatial Forgery Localization in ForgeryNet. Examples of (a) real image, (b) forgery image, (c) corresponding spatial annotations.

image sequence, sketch map, parsing mask, audio, label, or even noise. We first detect the *target* face  $\mathbf{I}_t^f$ , crop and align it, and then transform both the *target* face as well as *source* data to intermediate representations such as UV map, feature bank, 3DMM parameters and *etc*.

**Forgery Modeling.** These representations are forwarded to the forgery models to obtain a forged target face  $\tilde{\mathbf{I}}_t^f$ . We include five architecture variants as, 1) *Encoder-Decoder* [5], 2) *Vanilla GAN* [40], 3) *Pix2Pix* [26], 4) *RNN/LSTM* [7], and 5) *Graphics Formation* [13].

**Re-rendering Process.** To acquire the full forged target, the forged target face  $\tilde{\mathbf{I}}_t^f$  is re-rendered back to the target full image  $\mathbf{I}_t$  to obtain  $\tilde{\mathbf{I}}_t$ . In particular, according to different forgery procedures, 1)  $\tilde{\mathbf{I}}_t^f$  can be a *face mask*, shown in Fig. 4(e-1), which contains the area from the eyebrows to the face chin. 2)  $\tilde{\mathbf{I}}_t^f$  can also be a *face bounding-box*, illustrated in Fig. 4(e-2,3), which keeps the same bounding box as the original target face.

**Perturbation.** To better reflect real-world data distribution, we apply 36 types of perturbations to the forgery data  $\tilde{I}_t$ . We follow common practices in visual quality assessment [39] with distortions of compression, transmission, capture, color, *etc*.

### **3.3.** ForgeryNet Annotation

In contrast to most previous datasets, our ForgeryNet is annotated comprehensively both in image- and video-level across four tasks.

**Image Forgery Classification.** According to the forgery definition in Sec. 3.2.1, given a forgery image, we provide three types of forgery labels, *i.e.* labels for two-way (real / fake), three-way (real / fake with identity-replaced forgery approaches / fake with identity-remained forgery approaches), and *n*-way (n = 16, real and 15 respective forgery approaches) classification tasks respectively. These annotations make it possible to explore the correlation between different forgery meta-types or approaches.

**Spatial Forgery Localization.** As shown in Fig. 5, we take the forgery image  $\tilde{I}_t$  and the corresponding real image  $I_t$  to calculate their difference to obtain a *forgery distribution* 



Figure 6: Illustration of image- and video-level sets. From the inside to the outside are categories of *Identity-remained* and *Identity-replaced*, corresponding sub-types, specific forgery approaches and the situation of data split.

 $\tilde{\mathbf{I}}_{t}^{d}$ . In this paper, we define the Spatial Forgery Localization task as "localizing the face area manipulated by deep forgery approaches", and thus the forgery distribution before perturbation  $\tilde{\mathbf{I}}_{t}^{d}$  is taken as the ground-truth annotation. **Video Forgery Classification & Temporal Forgery Localization.** Note that in contrast to all the existing datasets, we construct our video forgery dataset with untrimmed forgery videos  $\tilde{\mathbf{V}}_{t}^{\prime}$ , each of which splices real and manipulated frames together. Same as image-forgery, *Video Forgery Classification* also contains three types of class annotations. We also provide the annotations on locations of manipulated segments in the untrimmed forgery video and propose a new task, *i.e. Temporal Forgery Localization*, to localize these forged segments.

# 4. ForgeryNet Settings

On ForgeryNet, we set up two benchmarks, image and video, with a series of tasks for face forgery analysis. **Dataset Preparation.** Both image- and video-level sets are split into training, validation and test subsets with a ratio close to 7:1:2. Forgery data distributions and catagories of the two sets are shown in Fig. **6**. Forgery data in each subset have identities matched with the corresponding real subset. The ratio of real to fake in each subset is close to 1:1.

### 4.1. Image Benchmark Settings

#### 4.1.1 Image Forgery Classification

In order to foster further researches on face forgery classification, we carefully design two protocols to evaluate forensics methods in this area.

**Protocol 1: Intra-forgery Evaluation.** In intra-forgery evaluation, all the real and fake data in the training set are used to train models, and the validation set is used for evaluation. This protocol has three variants, according to the definition in Sec. 3.3, *i.e.* two-/three-/*n*-way classification.

**Protocol 2: Cross-forgery Evaluation.** To further evaluate the generalization ability of training with our data, we conduct cross-forgery evaluation by training the evaluated

Table 2: **Image Forgery Classification (Protocol 1):** binary classification. We report accuracy and AUC scores of the compared forensics methods.

| Method                   | Param. | Acc   | AUC   |
|--------------------------|--------|-------|-------|
| MobileNetV3 Small [22]   | 1.7M   | 76.24 | 85.51 |
| MobileNetV3 Large [22]   | 4.2M   | 78.30 | 87.56 |
| EfficientNet-B0 [41]     | 4.0M   | 79.86 | 89.31 |
| ResNet-18 [21]           | 11.2M  | 78.31 | 87.75 |
| Xception [9]             | 20.8M  | 80.78 | 90.12 |
| ResNeSt-101 [45]         | 46.2M  | 82.06 | 91.02 |
| SAN19-patchwise [46]     | 18.5M  | 80.08 | 89.38 |
| ELA-Xception [20]        | 20.8M  | 73.77 | 82.69 |
| SNRFilters-Xception [8]  | 20.8M  | 81.09 | 90.52 |
| GramNet [31]             | 22.1M  | 80.89 | 90.20 |
| F <sup>3</sup> -Net [36] | 57.3M  | 80.86 | 90.15 |

forensics method with one certain type of manipulation and testing it with others. The manipulation type can either be general (*e.g. identity-replaced*), or specific (*e.g. ATVG-Net*). Note that this protocol only involves binary classification.

**Metrics.** For binary classification tasks, we evaluate with Accuracy (Acc) and the Area under ROC curve (AUC). For three- and *n*-way class settings, we use Accuracy (Acc) and mean Average Precision (mAP) as evaluation metrics.

### 4.1.2 Spatial Forgery Localization

Compared with classification tasks, spatial forgery localization aims to specify manipulated regions. Images along with forgery masks are used to train the localization model. **Metrics.** We utilize three metrics for evaluation: two variants of Intersection over Union (IoU) and L1 distance.

# 4.2. Video Benchmark Settings

**Video Forgery Classification.** Evaluation protocols for video forgery classification are generally similar to the ones designed for the image set, except that n=9 for n-class setting. Metrics are the same as those for image classification. **Temporal Forgery Localization.** For each video, forensics methods to be evaluated are expected to provide temporal boundaries of forgery segments and the corresponding confidence values. We follow metrics used in ActivityNet [18] evaluation, and employ Interpolated Average Precision (AP) as well as Average Recall@K (AR@K) for evaluating predicted segments with respect to the groundtruth ones.

# 5. Image Forgery Analysis Benchmark

# 5.1. Image Forgery Classification

**Protocol 1: Intra-forgery Evaluation.** For comprehensive evaluation, we provide results of two-way class classification with several representative models of different sizes. Considering the trade-off between performance and efficiency, we use Xception [9] as the baseline model. ELA-Xception [20] and SNRFilters-Xception [8]

Table 3: **Image Forgery Classification (Protocol 1):** multi-class settings and their mappings to binary classification. We report the accuracy, mAP and AUC scores.

|                     | 3-way class                     |                                  | 3→2-w                            | ay class                           |
|---------------------|---------------------------------|----------------------------------|----------------------------------|------------------------------------|
|                     | Acc.                            | mAP                              | Acc.                             | AUC                                |
| Xception            | 73.00                           | 89.90                            | 80.17                            | 89.92                              |
| GramNet             | 73.30                           | 90.00                            | 80.75                            | 90.13                              |
| F <sup>3</sup> -Net | 74.45                           | 90.41                            | 81.75                            | 90.63                              |
|                     |                                 |                                  |                                  |                                    |
|                     | 16-wa                           | y class                          | 16→2-w                           | vay class                          |
|                     | 16-wa<br>Acc.                   | y class<br>mAP                   | 16→2-w<br>Acc.                   | vay class<br>AUC                   |
| Xception            | 16-wa<br>Acc.<br>58.81          | y class<br>mAP<br>93.16          | 16→2-w<br>Acc.<br>81.00          | AUC 90.53                          |
| Xception<br>GramNet | 16-wa<br>Acc.<br>58.81<br>56.77 | y class<br>mAP<br>93.16<br>92.27 | 16→2-w<br>Acc.<br>81.00<br>80.83 | vay class<br>AUC<br>90.53<br>90.25 |

Table 4: **Image Forgery Classification (Protocol 2):** binary classification. We report the accuracy and AUC scores. Forensics methods trained with ID-replaced forgery approaches have significant performance drops when tested on unseen ID-remained forgery approaches, and *vice versa*.

|                     |             | ID-replaced ID-rem |       | nained |       |
|---------------------|-------------|--------------------|-------|--------|-------|
|                     |             | Acc.               | AUC   | Acc.   | AUC   |
| Vantion             | ID-replaced | 84.13              | 92.80 | 64.62  | 74.86 |
| лсерион             | ID-remained | 67.28              | 75.83 | 81.17  | 90.71 |
| CromNat             | ID-replaced | 82.82              | 92.54 | 62.72  | 74.28 |
| Grannvet            | ID-remained | 67.50              | 76.19 | 80.60  | 90.28 |
| $\mathbf{E}^3$ Not  | ID-replaced | 83.84              | 92.73 | 64.33  | 73.82 |
| F <sup>*</sup> -Net | ID-remained | 68.44              | 77.24 | 81.18  | 90.29 |

are two variants of Xception. Smaller models include MobileNetV3 [22], EfficientNet-B0 [41] and ResNet-18 [21]. We select ResNeSt-101 [45] as the large model. We also experiment with recent state-of-the-art methods for face forgery detection, *i.e.*  $F^3$ -Net [36] and GramNet [31], as well as a fully-attentional network SAN19 [46].

All experiments are conducted on face images cropped with face bounding boxes enlarged by  $1.3 \times$ . During training, we use several types of data augmentation to mimic distortions caused by compression and packet loss during transmission, so as to improve the generalization of developed models.

As presented in Tab. 2, we list binary classification metrics of all aforementioned forensics methods. We also show the corresponding ROC curves of these methods in Fig. 7(a). For three-way and 16-way classification experiments, as shown in Tab. 3, Acc scores show that classification becomes more difficult when the number of categories increases, yet the mAP metric indicates that the discrimination ability becomes higher instead. Moreover, after mapping back to binary classification, we can also observe slight performance boosts on F<sup>3</sup>-Net compared to training results with only binary labels. This suggests that more auxiliary information potentially makes the forensics model more discriminative.



Figure 7: **Image Forgery Classification (Protocol 1):** (a) We show the ROC curves of the compared methods under the setting of binary classification. (b)-(d) t-SNE feature visualization of the data manipulated by different forgery approaches, trained with binary, three-way and *n*-way classification respectively.



Figure 8: **Image Forgery Classification (Protocol 2):** (a) AUC score map, and (b) correlation map according to the AUC scores. X-axis denotes the tested forgery approach and Y-axis denotes the forgery approach for training.

Protocol 2: Cross-forgery Evaluation. For this protocol, we show the generalization ability of forensics methods across forgery approaches. Tab. 4 lists the results of models trained on ID-replaced but evaluated on ID-remained, and vice versa. The more exhaustive cross-forgery setting with 15 specific forgery approaches is also evaluated and shown in Fig. 8. We observe from these results that intra-forgery testing naturally performs the best. From Fig. 8(a), we can also see that training on ATVG-Net, StyleGAN2 or Blend-*Face* gives the best generalization performance on average. On the other hand, *DiscoFaceGAN* is the most generalizable forgery approach, while SC-FEGAN is the most difficult approach to generalize to. There is another interesting finding that forgery approaches with stronger similarity tend to induce better cross-forgery performance. For example, DiscoFaceGAN is a StyleGAN-based approach, thus training on the latter approach produces favorable results on the former. Similarly, StarGAN2 and the two face stack manipulations which both involve StarGAN2 generalize well to each other. In addition, as shown in Fig. 8(b), forgery approaches belonging to the same meta-category usually have higher correlations mutually. For example, for meta-category Face reenactment, if a forensics method can obtain good perfor-

| Table 5:   | Spatial   | Forgery              | Locali | zation.             | We co  | mpare   | re- |
|------------|-----------|----------------------|--------|---------------------|--------|---------|-----|
| sults with | n three m | netrics, <i>i.e.</i> | , IOU, | $IOU_{\text{diff}}$ | and L1 | distand | ce. |

|                            | incures     | ,, 1     |                |                     | u El u    | iotunee.           |
|----------------------------|-------------|----------|----------------|---------------------|-----------|--------------------|
| Method                     | Ic          | U        | 0.01           | IoU <sub>diff</sub> | 0.1       | Loss <sub>11</sub> |
|                            | 0.1         | 0.2      | 0.01           | 0.05                | 0.1       |                    |
| Xception+Reg.              | 89.55       | 93.70    | 67.57          | 83.25               | 89.22     | 0.0131             |
| Xeption+Unet [37]          | 95.99       | 98.76    | 79.71          | 92.70               | 97.13     | 0.0134             |
| HRNet [42]                 | 96.27       | 98.78    | 88.73          | 92.99               | 96.27     | 0.0114             |
| Original Ta                | rget Before | Perturb. | After Perturb. | Ground              | Truth Pre | dicted Mask        |
| (a)<br>Face<br>Replacement |             |          |                | 10° - 00-           |           |                    |
| (b)<br>Face<br>Reenactment |             |          | 6              | 113 5 (11)          |           | 18 10              |
| (c)<br>Face<br>Editing     |             |          |                |                     |           |                    |
| (d)<br>Real<br>Face        |             |          |                |                     |           | 15                 |
| (e)<br>Real<br>Face        |             |          | 6              |                     |           |                    |

Figure 9: **Spatial Forgery Localization.** Examples of predicted manipulation masks by HRNet.

mance on ATVG-Net, it may also work for FirstOrderMotion and Talking-headVideo.

### 5.2. Spatial Forgery Localization

We evaluate pixel regression and other two segmentation methods for the spatial localization task. UNet [37] is a popular segmentation architecture, which has been widely used. For comparison, we also adopt HRNet [42] because of its superior performance on other datasets.

In Tab. 5, HRNet outperforms other methods. Especially in terms of  $IoU_{diff}$  with threshold 0.01, HRNet surpasses

Table 6: Video Forgery Classification (Protocol 1): binary classification. We report accuracy and AUC scores under two crop strategies. Video-level classification has better results than the image-level setting.

|                |            | Single-crop |       | Multi-crop |       |
|----------------|------------|-------------|-------|------------|-------|
| Method         | Parameters | Acc         | AUC   | Acc        | AUC   |
| X3D-M [15]     | 2.9M       | 87.93       | 93.75 | 88.97      | 96.99 |
| Slow-only [16] | 31.6M      | 86.76       | 92.64 | 87.37      | 95.96 |
| TSM [28]       | 23.5M      | 88.04       | 93.05 | 89.11      | 96.25 |
| SlowFast [16]  | 33.6M      | 88.78       | 93.88 | 89.92      | 97.28 |

Table 7: **Video Forgery Classification (Protocol 1):** multiclass settings and their mappings to binary classification. We report the accuracy, mAP and AUC scores.

| Mada al       | 3-way                  | / class                 | 3→2-w                  | $3 \rightarrow 2$ -way class |  |  |
|---------------|------------------------|-------------------------|------------------------|------------------------------|--|--|
| Method        | Acc.                   | mAP                     | Acc.                   | AUC                          |  |  |
| X3D-M [15]    | 84.00                  | 94.55                   | 87.69                  | 93.78                        |  |  |
| SlowFast [16] | 85.73                  | 94.89                   | 89.11                  | 94.37                        |  |  |
|               |                        |                         |                        |                              |  |  |
|               | 9-way                  | / class                 | 9→2-w                  | ay class                     |  |  |
|               | 9-way<br>Acc.          | class<br>mAP            | 9→2-w<br>Acc.          | ay class<br>AUC              |  |  |
| X3D-M [15]    | 9-way<br>Acc.<br>76.91 | v class<br>mAP<br>95.06 | 9→2-w<br>Acc.<br>87.51 | ay class<br>AUC<br>93.81     |  |  |

other methods by more than 10%. We also present predicted manipulation maps for several test samples. In Fig. 9(c), the slight beard change is hard to detect, while in Fig. 9(d), a real image is misjudged as manipulated.

# 6. Video Forgery Analysis Benchmark

### 6.1. Video Forgery Classification

In this section, we select several typical video backbones of different sizes: X3D-M [15], Slow-only R-50 [16], TSM [28], and SlowFast R-50 [16]. We sample 16 frames with temporal stride 4 as input to all models.

Binary classfication results of video-level forensics methods are listed in Tab. 6. Compared to image-level evaluation, video-level Acc and AUC are generally higher. SlowFast [16] obtains the best performance on video classification, while X3D-M [15], with only a very small number of parameters, also gives satisfying results. We select these two as representatives of large and small models respectively in subsequent experiments, as displayed in Tab. 7 and Tab. 8. Cross-forgery evaluation results are worse than their image counterparts, suggesting harder generalization with temporal information.

# 6.2. Temporal Forgery Localization

We experiment with both frame-based and video-based models for temporal localization. For frame-based model, after binarizing frame predictions with a fixed threshold (0.25), we select consecutive fake sequences, with different tolerance levels for real frames in the middle, as final proposals. The confidence of a proposal is simply the average of the original frame scores. We adopt Boundary-Sensitive Network (BSN) [30] and Boundary-Matching

Table 8: Video Forgery Classification (Protocol 2): binary classification. Forensics methods trained with IDreplaced forgery approaches have substantial performance drops (even more significant than their image-level counterparts) when tested on unseen ID-remained forgery approaches, and *vice versa*.

|          |             | ID-replaced |       | ID-rer | nained |
|----------|-------------|-------------|-------|--------|--------|
|          |             | Acc.        | AUC   | Acc.   | AUC    |
| V2D M    | ID-replaced | 87.92       | 92.91 | 55.25  | 65.59  |
| A3D-W    | ID-remained | 55.93       | 62.87 | 88.85  | 95.40  |
| ClawEast | ID-replaced | 88.26       | 92.88 | 52.64  | 64.83  |
| SlowFast | ID-remained | 52.70       | 61.50 | 87.96  | 95.47  |

Table 9: **Temporal Forgery Localization.** We show AP, AR and mAP scores of all compared methods.

|                   | AR    |       |       | avg.  |       |       |
|-------------------|-------|-------|-------|-------|-------|-------|
|                   | 2     | 5     | 0.5   | 0.75  | 0.9   | AP    |
| Xception [9]      | 25.83 | 73.95 | 68.29 | 62.84 | 58.30 | 62.83 |
| X3D-M+BSN [30]    | 81.33 | 86.88 | 80.46 | 77.24 | 55.09 | 70.29 |
| X3D-M+BMN [29]    | 88.44 | 91.99 | 90.65 | 88.12 | 74.95 | 83.47 |
| SlowFast+BSN [30] | 83.63 | 88.78 | 82.25 | 80.11 | 60.66 | 73.42 |
| SlowFast+BMN [29] | 90.64 | 93.49 | 92.76 | 91.00 | 80.02 | 86.85 |

Network (BMN) [29] on top of X3D-M and SlowFast features as the video-based models.

Tab. 9 compares these methods on the validation set. In particular, video-based methods perform significantly better than the frame-based method, demonstrating the importance of applying a boundary-aware network. Additionally, BMN outperforms BSN with large margins, and achieves  $\sim$ 87 average AP. This is of great significance since it shows our model is capable of effectively locating manipulated media in a large video database. We hope our results can inspire more future works on forgery localization.

# 7. Conclusion

In this paper, we present ForgeryNet, a new mega-scale benchmark for both image- and video-level face forgery analysis. Compared with existing datasets for face forgery, ForgeryNet possesses more variety and is more comprehensive in terms of wild sources, various manipulation approaches, diverse re-rendering process and richness of annotations. We further introduce four possible applications with ForgeryNet: image and video classification, spatial and temporal localization. The results obtained in these tasks help us better understand facial forgery towards realworld scenarios. For future works, we welcome interested researchers to contribute more novel facial forgery approaches. More forgery analysis can also be studied on our dataset to improve the defense capabilities.

Acknowledgments This work is supported by Key Research and Development Program of Guangdong Province, China, under Grant No. 2019B010154003, as well as NTU NAP and A\*STAR through the Industry Alignment Fund -Industry Collaboration Projects Grant, the National Natural Science Foundation of China under Grant No. 61906012.

# References

- Faceswap. https://github.com/ MarekKowalski/FaceSwap/.
- [2] Swapme. https://itunes.apple.com/us/app/ swapme-by-faciometrics/.
- [3] Fakeapp. https://www.fakeapp.com/, 2018.
- [4] Google ai blog. contributing data to deepfake detection research. https://ai.googleblog.com/ 2019/09/contributing-data-to-deepfakedetection.html, 2019.
- [5] faceswap. https://github.com/deepfakes/ faceswap, 2020.
- [6] Houwei Cao, David G Cooper, Michael K Keutmann, Ruben C Gur, Ani Nenkova, and Ragini Verma. Crema-d: Crowd-sourced emotional multimodal actors dataset. *IEEE Trans. Affect. Comput.*, 5(4):377–390, 2014.
- [7] Lele Chen, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In CVPR, 2019.
- [8] Mo Chen, Vahid Sedighi, Mehdi Boroumand, and Jessica Fridrich. Jpeg-phase-aware convolutional neural network for steganalysis of jpeg images. In *the 5th ACM Workshop*, 2017.
- [9] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *CVPR*, 2017.
- [10] Joon Son Chung, Arsha Nagrani, and Andrew Zisserman. Voxceleb2: Deep speaker recognition. arXiv preprint arXiv:1806.05622, 2018.
- [11] Hao Dang, Feng Liu, Joel Stehouwer, Xiaoming Liu, and Anil K Jain. On the detection of digital face manipulation. In *CVPR*, pages 5781–5790, 2020.
- [12] Brian Dolhansky, Joanna Bitton, Ben Pflaum, Jikuo Lu, Russ Howes, Menglin Wang, and Cristian Canton Ferrer. The deepfake detection challenge dataset. arXiv preprint arXiv:2006.07397, 2020.
- [13] Bernhard Egger, William AP Smith, Ayush Tewari, Stefanie Wuhrer, Michael Zollhoefer, Thabo Beeler, Florian Bernard, Timo Bolkart, Adam Kortylewski, Sami Romdhani, et al. 3d morphable face models—past, present, and future. *TOG*, 39(5):1–38, 2020.
- [14] A. Ephrat, I. Mosseri, O. Lang, T. Dekel, K Wilson, A. Hassidim, W. T. Freeman, and M. Rubinstein. Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation. arXiv preprint arXiv:1804.03619, 2018.
- [15] Christoph Feichtenhofer. X3d: Expanding architectures for efficient video recognition. In CVPR, 2020.
- [16] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *ICCV*, 2019.
- [17] Ohad Fried, Ayush Tewari, Michael Zollhöfer, Adam Finkelstein, Eli Shechtman, Dan B Goldman, Kyle Genova, Zeyu Jin, Christian Theobalt, and Maneesh Agrawala. Text-based editing of talking-head video. *TOG*, 38(4):1–14, 2019.
- [18] Bernard Ghanem, Juan Carlos Niebles, Cees Snoek, Fabian Caba Heilbron, Humam Alwassel, Victor Escorcia,

Ranjay Krishna, Shyamal Buch, and Cuong Duc Dao. The activitynet large-scale activity recognition challenge 2018 summary. *arXiv preprint arXiv:1808.03766*, 2018.

- [19] Luca Guarnera, Oliver Giudice, and Sebastiano Battiato. Deepfake detection by analyzing convolutional traces. In *CVPRW*, 2020.
- [20] Teddy Surya Gunawan, Siti Amalina Mohammad Hanafiah, Mira Kartiwi, Nanang Ismail, Nor Farahidah Za'bah, and Anis Nurashikin Nordin. Development of photo forensics algorithm by detecting photoshop manipulation using error level analysis. *IJEECS*, 7(1):131–137, 2017.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In CVPR, 2016.
- [22] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *ICCV*, 2019.
- [23] Nils Hulzebosch, Sarah Ibrahimi, and Marcel Worring. Detecting cnn-generated facial images in real-world scenarios. In CVPRW, 2020.
- [24] Liming Jiang, Ren Li, Wayne Wu, Chen Qian, and Chen Change Loy. Deeperforensics-1.0: A large-scale dataset for real-world face forgery detection. In CVPR, 2020.
- [25] Pavel Korshunov and Sébastien Marcel. Deepfakes: a new threat to face recognition? assessment and detection. arXiv preprint arXiv:1812.08685, 2018.
- [26] Lingzhi Li, Jianmin Bao, Hao Yang, Dong Chen, and Fang Wen. Faceshifter: Towards high fidelity and occlusion aware face swapping. *arXiv preprint arXiv:1912.13457*, 2019.
- [27] Yuezun Li, Xin Yang, Pu Sun, Honggang Qi, and Siwei Lyu. Celeb-df: A large-scale challenging dataset for deepfake forensics. In *CVPR*, 2020.
- [28] Ji Lin, Chuang Gan, and Song Han. Tsm: Temporal shift module for efficient video understanding. In *ICCV*, 2019.
- [29] Tianwei Lin, Xiao Liu, Xin Li, Errui Ding, and Shilei Wen. Bmn: Boundary-matching network for temporal action proposal generation. In *ICCV*, 2019.
- [30] Tianwei Lin, Xu Zhao, Haisheng Su, Chongjing Wang, and Ming Yang. Bsn: Boundary sensitive network for temporal action proposal generation. In *ECCV*, 2018.
- [31] Zhengzhe Liu, Xiaojuan Qi, and Philip HS Torr. Global texture enhancement for fake face detection in the wild. In *CVPR*, 2020.
- [32] Steven R Livingstone and Frank A Russo. The ryerson audio-visual database of emotional speech and song (ravdess): A dynamic, multimodal set of facial and vocal expressions in north american english. *PLOS ONE*, 13(5):e0196391, 2018.
- [33] Francesco Marra, Cristiano Saltori, Giulia Boato, and Luisa Verdoliva. Incremental learning for the detection and classification of gan-generated images. In *WIFS*, 2019.
- [34] Britt Paris and Joan Donovan. Deepfakes and cheap fakes. United States of America: Data & Society, 2019.
- [35] Ivan Petrov, Daiheng Gao, Nikolay Chervoniy, Kunlin Liu, Sugasa Marangonda, Chris Umé, Jian Jiang, Luis RP, Sheng

Zhang, Pingyu Wu, et al. Deepfacelab: A simple, flexible and extensible face swapping framework. *arXiv preprint arXiv:2005.05535*, 2020.

- [36] Yuyang Qian, Guojun Yin, Lu Sheng, Zixuan Chen, and Jing Shao. Thinking in frequency: Face forgery detection by mining frequency-aware clues. In *ECCV*, 2020.
- [37] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015.
- [38] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. Faceforensics++: Learning to detect manipulated facial images. In *CVPR*, pages 1–11, 2019.
- [39] Muhammad Shahid, Andreas Rossholm, Benny Lövström, and Hans-Jürgen Zepernick. No-reference image and video quality assessment: a classification and review of recent approaches. *EURASIP J IMAGE VIDE*, 2014(1):40, 2014.
- [40] Yujun Shen, Bolei Zhou, Ping Luo, and Xiaoou Tang. Facefeat-gan: a two-stage approach for identity-preserving face synthesis. arXiv preprint arXiv:1812.01288, 2018.
- [41] Mingxing Tan and Quoc V Le. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv* preprint arXiv:1905.11946, 2019.
- [42] 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. *TPAMI*, 2020.
- [43] Xin Yang, Yuezun Li, and Siwei Lyu. Exposing deep fakes using inconsistent head poses. In *ICASSP*, 2019.
- [44] Ning Yu, Larry Davis, and Mario Fritz. Attributing fake images to gans: Analyzing fingerprints in generated images. *arXiv preprint arXiv:1811.08180*, 2018.
- [45] Hang Zhang, Chongruo Wu, Zhongyue Zhang, Yi Zhu, Zhi Zhang, Haibin Lin, Yue Sun, Tong He, Jonas Mueller, R Manmatha, et al. Resnest: Split-attention networks. arXiv preprint arXiv:2004.08955, 2020.
- [46] Hengshuang Zhao, Jiaya Jia, and Vladlen Koltun. Exploring self-attention for image recognition. In *CVPR*, 2020.
- [47] Peng Zhou, Xintong Han, Vlad I Morariu, and Larry S Davis. Two-stream neural networks for tampered face detection. In *CVPRW*, 2017.