# Detecting and Recovering
# Sequential DeepFake Manipulation

Rui Shao, Tianxing Wu, and Ziwei Liu⋆

S-Lab, Nanyang Technological University
{rui.shao, twu012, ziwei.liu}@ntu.edu.sg
https://rshaojimmy.github.io/Projects/SeqDeepFake

**Abstract.** Since photorealistic faces can be readily generated by facial manipulation technologies nowadays, potential malicious abuse of these technologies has drawn great concerns. Numerous deepfake detection methods are thus proposed. However, existing methods only focus on detecting *one-step* facial manipulation. As the emergence of easy-accessible facial editing applications, people can easily manipulate facial components using *multi-step* operations in a sequential manner. This new threat requires us to detect a sequence of facial manipulations, which is vital for both detecting deepfake media and recovering original faces afterwards. Motivated by this observation, we emphasize the need and propose a novel research problem called Detecting Sequential DeepFake Manipulation (**Seq-DeepFake**). Unlike the existing deepfake detection task only demanding a binary label prediction, detecting Seq-DeepFake manipulation requires correctly predicting a sequential vector of facial manipulation operations. To support a large-scale investigation, we construct the first Seq-DeepFake dataset, where face images are manipulated sequentially with corresponding annotations of sequential facial manipulation vectors. Based on this new dataset, we cast detecting Seq-DeepFake manipulation as a specific image-to-sequence (*e.g.* image captioning) task and propose a concise yet effective Seq-DeepFake Transformer (**SeqFakeFormer**). Moreover, we build a comprehensive benchmark and set up rigorous evaluation protocols and metrics for this new research problem. Extensive experiments demonstrate the effectiveness of SeqFakeFormer. Several valuable observations are also revealed to facilitate future research in broader deepfake detection problems.
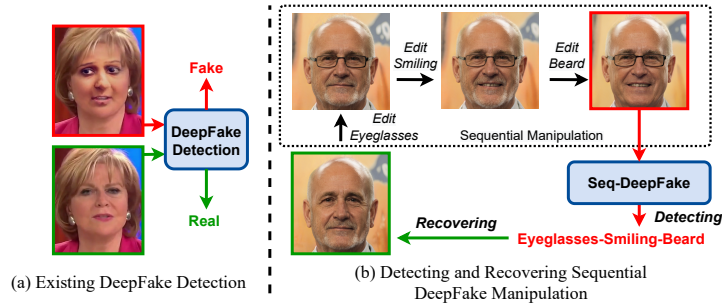
**Keywords:** DeepFake Detection, Sequential Facial Manipulation

## 1 Introduction

In recent years, hyper-realistic face images can be generated by deep generative models which are visually extremely indistinguishable from real images. Meanwhile, the significant improvement for image synthesis brings security

---

⋆ Corresponding author

(a) Existing DeepFake Detection

(b) Detecting and Recovering Sequential
DeepFake Manipulation

**Fig. 1:** Comparison between (a) existing deepfake detection and (b) proposed detecting and recovering sequential deepfake manipulation.

concerns on potential malicious abuse of these techniques that produce misinformation and fabrication, which is known as *deepfake*. To address this security issue, various deepfake detection methods have been proposed to detect such forged faces. As illustrated in Fig. 1 (a), given the manipulated face image generated by face swap algorithm [27] and the original face image, the existing deepfake detection task requires the model to predict the correct binary labels (Real/Fake).
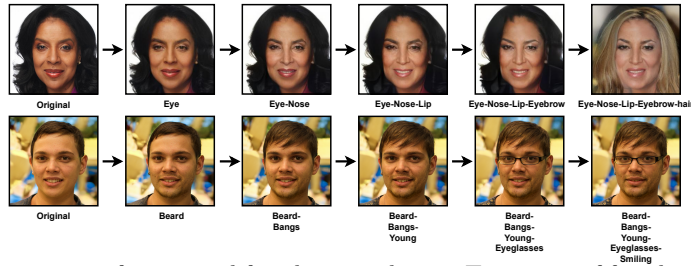
With the increasing popularity of easy-accessible facial editing applications, such as YouCam Makeup[1], FaceTune2[2], and YouCam Perfect[3], it is convenient for people to edit face images in daily life. Compared to existing deepfake techniques mainly carrying out *one-step* facial manipulation [27,11], we can now easily manipulate face images using *multi-step* operations in a *sequential* manner. As shown in Fig. 1 (b), the original image can be manipulated by adding eyeglasses, making a bigger smile and removing beard sequentially. This expands the scope of existing deepfake problem by adding sequential manipulation information and poses a new challenge for current *one-step* deepfake detection methods. This observation motivates us to introduce a new research problem — Detecting Sequential Deepfake Manipulation (**Seq-Deepfake**). We summarize several key differences between detecting Seq-Deepfake and the existing deepfake detection: 1) rather than only predicting binary labels (Real/Fake), detecting Seq-Deepfake aims to detect sequences of facial manipulations with diverse sequence lengths. For example, the model is required to predict a 3-length sequence as 'Eyeglasses-Smiling-Beard' for the manipulated image as shown in Fig. 1 (b). 2) As illustrated in Fig. 1 (b), beyond pure forgery detection, we can further **recover** the original faces (refer to Section 5.4 of Experiments) based on the detected sequences of facial manipulation in Seq-Deepfake. This greatly enriches the benefits of detecting Seq-Deepfake manipulation.

To facilitate the study of detecting Seq-Deepfake, this paper contributes the first Seq-Deepfake dataset. Fig. 2 shows some samples in Seq-Deepfake dataset. From Fig. 2, it can be seen that one face image can be sequentially

---

[1] https://apps.apple.com/us/app/youcam-makeup-selfie-editor/id863844475

[2] https://apps.apple.com/us/app/facetune2-editor-by-lightricks/id1149994032

[3] https://apps.apple.com/us/app/youcam-perfect-photo-editor/id768469908

**Fig. 2:** Illustration of sequential facial manipulation. Two types of facial manipulation approaches are considered, *i.e.* facial components manipulation [16] in the first row and facial attributes manipulation [13] in the second row.

manipulated with different number of steps (from minimum 1 step to maximum 5 steps), which leads to facial manipulation sequences with diverse lengths. It is extremely difficult to distinguish the original and manipulated face images, and even harder to figure out the exact manipulation sequences. To make our study more comprehensive, we consider two different facial manipulation techniques, facial components manipulation [16] and facial attributes manipulation [13], which are displayed in the first and second row, respectively in Fig. 2.

Most current facial manipulation applications are built based on Generative Adversarial Network (GAN). It is well known that the semantic latent space learned by GAN is difficult to be perfectly disentangled [28,18]. We argue that this defect is likely to leave some spatial as well as sequential manipulation traces unveiling sequential facial manipulations. Based on this observation, to detect such two types of manipulations traces, we cast detecting Seq-Deepfake as a specific image-to-sequence (*e.g.* image captioning) task and thus propose a concise yet effective Seq-DeepFake Transformer (**SeqFakeFormer**). Two key parts are devised in SeqFakeFormer: **Spatial Relation Extraction** and **Sequential Relation Modeling with Spatially Enhanced Cross-attention**. Given a manipulated image, to adaptively capture subtle spatial manipulation regions, SeqFakeFormer feeds the image into a deep convolutional neural network (CNN) to learn its feature maps. Then we extract the relation of spatial manipulations captured in feature maps using the self-attention modules of transformer encoder, obtaining features of spatial relation, i.e. spatial manipulation traces. After that, the decoder of SeqFakeFormer models the sequential relation of extracted features of spatial relation via cross-attention modules in an autoregressive mechanism, contributing to the detection of sequential manipulation traces, and thus detecting the facial manipulation sequences. To enable more effective cross-attention given limited annotations of facial manipulation sequences in Seq-DeepFake, SeqFakeFormer further integrates a Spatially Enhanced Cross-Attention (SECA) module in the decoder. This module enriches the spatial information of annotations of manipulation sequences by learning a spatial weight map. After fusing the spatial weight map with the cross-attention map, a spatially enhanced cross-attention can be achieved.

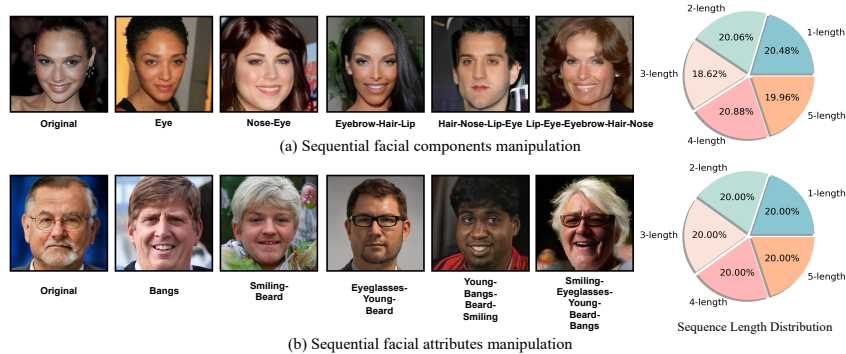Main contributions of our paper can be summarized as follows:

– We introduce a new research problem named Detecting Sequential Deepfake Manipulation (**Seq-DeepFake**), with the objective of detecting sequences of facial manipulations, which expands the scope and poses a new challenge for deepfake detection.
– We contribute the Sequential Deepfake Dataset with sequential manipulated face images using two different facial manipulation techniques. Corresponding annotations of manipulation sequences are provided.
– We propose a powerful Seq-DeepFake Transformer (**SeqFakeFormer**). A comprehensive benchmark is built and rigorous evaluation protocols and metrics are designed for this novel research problem. Extensive quantitative and qualitative experiments demonstrate its superiority.

## 2   Related Work

**Deepfake detection.** Current deepfake detection methods can be roughly categorized into spatial-based and frequency-based deepfake detection. The majority of spatial-based deepfake detection methods focus on capturing visual cues from spatial domain. Face X-ray [20] is proposed to detect the blending boundary left in the face forgery process as visual cues for real/fake detection. A multi-attentional deepfake detection network is proposed in [36] to integrate low-level textural features and high-level semantic features. Zhu *et al.* [39] introduce 3D decomposition into forgery detection and propose a two-stream network to fuse decomposed features for detection. Pair-wise self-consistency learning (PCL) [37] is introduced to detect inconsistency of source features within the manipulated images. Inconsistencies in semantically high-level mouth movements are captured in  [9] by fine-tuning a temporal network pretrained on lipreading. On the other hand, some methods pay attention to the frequency domain for detecting spectrum artifacts. There exist distinct spectrum distributions and characteristics between real and fake images in the high-frequency part of Discrete Fourier Transform (DFT) [5,6]. Qian *et al.* [26] propose a $F^3$-Net to learn local frequency statistics based on Discrete Cosine Transform (DCT) to mine forgery. Liu *et al.* [22] present a Spatial-Phase Shallow Learning method to fuse spatial image and phase spectrum for the up-sampling artifacts detection. A two-stream model is devised in [24] to model the correlation between extracted high-frequency features and regular RGB features to learn generalizable features. A frequency-aware discriminative feature learning framework [19] is introduced to integrate metric learning and adaptive frequency features learning for face forgery detection.

So far, several deepfake datasets have been released to public, such as FaceForensics++ [27], Celeb-DF [21], Deepfake Detection Challenge (DFDC) [4], and DeeperForensics-1.0 (DF1.0) [12]. However, only binary labels are provided in most of existing deepfake datasets, and thus most of the above works are trained to carry out binary classification, which results in performance saturation and poor generalization.

**Facial editing.** Several methods have been proposed for editing facial components (*i.e.* eye, nose, month). Lee *et al.* [17] present a geometry-oriented

Fig. 3: Illustration of Seq-Deepfake dataset. Samples of Seq-Deepfake are provided with annotations of manipulation sequences. We also show sequence length distribution.

face manipulation network MaskGAN for diverse and interactive face manipulation guided by semantic masks annotations. A semantic region-adaptive normalization (SEAN) [38] is proposed to facilitate manipulating face images by encoding images into the per-region style codes conditioned on segmentation masks. StyleMapGAN [16] introduces explicit spatial dimensions to the latent space and manipulates facial components by blending the latent spaces between reference and original faces. Moreover, some works target editing specific facial attributes such as age progression [35], and smile generation [34]. Some recent works discover semantically meaningful directions in the latent space of a pretrained GAN so as to carry out facial attributes editing by moving the latent code along these directions [28,29,40,31,30]. InterFaceGAN [28,29] tries to disentangle attribute representations in the latent space of GANs by searching a hyperplane, of which a normal vector is used as the editing direction. Fine-grained facial attributes editing is achieved by [13] through searching a curving trajectory with respect to attribute landscapes in the latent space of GANs.

## 3   Sequential Deepfake Dataset

To support the novel research problem, we generate a large-scale Sequential Deepfake (Seq-Deepfake) dataset consisting of sequential manipulated face images based on two representative facial manipulation techniques, facial components manipulation [16] and facial attributes manipulation [13]. Unlike most of existing deepfake datasets [27,11] only providing binary labels, the proposed dataset contains annotations of manipulation sequences with diverse sequence lengths. Details of generation pipelines based on the two facial manipulation techniques are as follows.

**Sequential facial components manipulation.** We adopt the StyleMapGAN proposed in [16] for facial components manipulation. Facial components manipulation is carried out based on original images from CelebA-HQ [23,14] and corresponding facial component masks from CelebAMask-HQ [17] dataset.

Facial components manipulation aims to transplant some facial components of a reference image to an original image with respect to a mask that indicates the components to be manipulated. Specifically, we project the original image and the reference image through the encoder of StyleMapGAN to obtain stylemaps, which are intermediate latent spaces with spatial dimensions. Then, the facial components manipulation is carried out by blending the stylemaps extracted from reference and original faces based on facial component masks. Due to the inevitable appearance of degraded images in the generation process, we adopt the Generated Image Quality Assessment (GIQA) algorithm [8] to quantitatively evaluate the quality of each generated image and then filter out some low-quality ones based on the pre-defined threshold. Fig. 3 (a) shows some samples with corresponding annotations of sequential facial components manipulation. Through this data generation pipeline, we can finally generate 35,166 manipulated face images annotated with 28 types of manipulation sequences in different lengths (including original). As illustrated in Fig. 3 (a), the proportions of 1-5 different lengths of manipulation sequences are: 20.48%, 20.06%, 18.62%, 20.88%, 19.96%.

**Sequential facial attributes manipulation.** Unlike facial components manipulation methods that swap certain local parts from a reference image to an original image, facial attributes manipulation approaches directly change specific attributes on the original face image without any reference images. To take this manipulation type into consideration, we utilize the fine-grained facial editing method proposed by [13]. This method aims to learn a location-specific semantic field for each editing type on the training set, then edit this attribute of interest on the given face image to a user-defined degree by stepping forward or backward on the learned curve in latent space. Based on this idea, we further generate face images with sequential facial attributes manipulation by performing the editing process in a sequential manner. Specifically, we first sample latent codes from the StyleGAN trained on FFHQ dataset [15] to generate original images. Then according to pre-defined attribute sequences, we progressively manipulate each attribute on the original face to another randomly chosen degree using the above method. After generating the final manipulation results, we also perform GIQA algorithm to filter out low-quality samples. Using this pipeline, we generate 49,920 face images with 26 manipulation sequence types, with the length of each sequence ranging from 1 to 5. Since this generation pipeline is more controllable than facial components manipulation, we construct a more balanced dataset, as shown in Fig. 3 (b).

## 4   Our Approach

### 4.1   Motivation

Most current facial manipulation applications are constructed using algorithms of Generative Adversarial Network (GAN). However, it is a well known fact that due to imperfect semantic disentanglement in the latent space of GAN [28,18], manipulating one facial component or attribute is likely to affect

**Eye-Nose**                    **Nose-Eye**        **Bangs-Smiling**              **Smiling-Bangs**
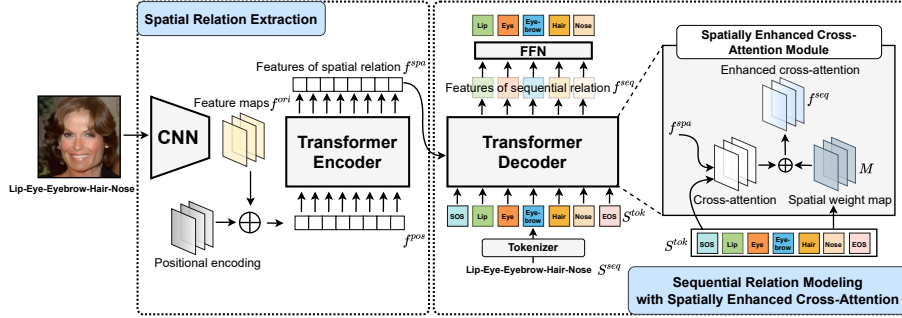(a) Facial components manipulation                  (b) Facial attributes manipulation

**Fig. 4:** Effect of different sequential order for facial manipulation. Switching the sequential order of manipulations between (a) eye and nose and (b) bangs and smiling results in different facial manipulations.

the others. As shown in the first row of Fig. 2, manipulating the nose in the step of 'Eye-Nose' simultaneously results in some little modification on the eye and mouth components compared to the former step 'Eye', which alters the overall **spatial relation** among facial components. We can thus discover some **spatial manipulation traces** from the spatial relation. Furthermore, as illustrated in Fig. 4, switching the sequential order of manipulations (*e.g.,* manipulation order between eye and nose in (a) and bangs and smiling in (b) in Fig. 4) causes different facial manipulation results (*e.g.,* distinct gazes in (a) and distinct amount of bangs in (b) in Fig. 4), which indicates that when changing the sequential order of manipulations, the above overall spatial relation of facial components altered by manipulations will also be changed. This means there exists sequential information from spatial relation that reflects the sequential order of manipulations, which corresponds to the facial manipulation sequence. That is, we can extract the spatial relation among facial components to unveil the **spatial manipulation traces** and model their **sequential relation** to detect the facial manipulation sequence. We thus regard the sequential relation as **sequential manipulation traces**.

## 4.2   Overview

Based on the above observation, we cast detecting Seq-Deepfake manipulation as a specific image-to-sequence task, where inputs are manipulated/original images and outputs are facial manipulation sequences. Three challenges will be encountered when addressing the task. 1) From Fig. 2 and 3, it can be seen that distinguishing original faces and sequential manipulated faces is extremely hard. Besides, with respect to different people, differences in face contour cause diverse manipulation regions for the same type of facial components/attributes manipulation. Thus, given indistinguishable and diverse facial manipulations, how to adaptively capture subtle manipulation regions and model their spatial relation accurately is quite challenging. 2) Based on the spatial relation of manipulated components/attributes, how to precisely model their sequential relation so as to detect the sequential facial manipulation is another challenge. 3) Compared to normal image-to-sequence task (*e.g.* image captions), the annotations of manipulation sequences are much shorter and thus less informative in our task. Therefore, how to effectively learn the sequential information of facial manipulations given limited annotations of manipulation sequences should also be considered.

**Fig. 5:** Overview of proposed Seq-DeepFake Transformer (**SeqFakeFormer**). We first feed the face image into a CNN to learn features of spatial manipulation regions, and extract their spatial relation via self-attention modules in the encoder. Then sequential relation based on features of spatial relation is modeled to detect the sequential facial manipulation. A spatial enhanced cross-attention module is integrated into the decoder, contributing to a more effective cross-attention.

To cope with the above three challenges, as shown in Fig. 5, we propose a Seq-DeepFake Transformer (**SeqFakeFormer**), which is composed of two key parts: **Spatial Relation Extraction**, **Sequential Relation Modeling with Spatially Enhanced Cross-attention**. To capture spatial manipulation traces, features of subtle manipulation regions are first adaptively captured by a CNN and their spatial relation are extracted via self-attention modules in the transformer encoder. After that, we capture sequential manipulation traces by modeling sequential relation based on features of spatial relation through cross-attention modules deployed in the decoder with an auto-regressive mechanism. To achieve more effective cross-attention given limited annotations of manipulation sequences, a spatially enhanced cross-attention module is devised to generate different spatial weight maps for corresponding manipulations to carry out cross-attention. In the following subsections, we describe all components in detail.

### 4.3   Spatial Relation Extraction

To adaptively capture subtle and various facial manipulation regions, we exploit a CNN to learn feature maps of the input image. Given an input image $x \in \mathrm{R}^{3 \times H' \times W'}$, we first feed it into a CNN [10] to extract its visual feature maps $f^{ori} = \mathrm{CNN}(x)$, $f^{ori} \in \mathrm{R}^{C \times H \times W}$, where $H', W'$, and $H, W$ are the height and width of the input image and its corresponding feature maps, respectively. $C$ is the number of channels of feature maps.

Since the transformer architecture is permutation-invariant, we supplement original visual features maps $f^{ori}$ with fixed positional encodings [25,1], resulting in feature maps denoted as $f^{pos}$. Since transformer encoder accepts a sequence as input, we reshape the spatial dimensions of $f^{pos}$ to one dimension, generating reshaped features $f^{pos} \in \mathrm{R}^{C \times HW}$. After feeding into the transformer encoder, $f^{pos}$ conducts self-attention by generating the key, query,

and value features $K, Q, V$ so as to extract the relations among all spatial positions. Through this self-attention operation on CNN features, spatial relation of manipulation regions are exploited and thus spatial manipulation traces can be extracted. To further facilitate spatial relation extraction, this paper adopts multi-head self-attention which splits features $f^{pos}$ into multiple groups along the channel dimension. The multi-head normalized attention based on dot-product is as follows:

$$f_i^{spa} = \text{Softmax}(K_i^T Q_i / \sqrt{d}) V_i, f^{spa} = \text{Concat}(f_1^{spa}, ..., f_D^{spa}) \tag{1}$$

where $K_i, Q_i, V_i$ denote the $i$-th group of the key, query, and value features, $d$ is dimension of queries and keys, and total $D$ groups are generated. We then concatenate all the groups to form the features of spatial relation $f^{spa}$ as the output of encoder.

## 4.4    Sequential Relation Modeling with Spatially Enhanced Cross-Attention

Given features of spatial relation $f^{spa}$ extracted from the encoder, we propose to model the sequential relation among them to detect the facial manipulation sequences. To this end, we carry out cross-attention between features of spatial relation $f^{spa}$ and corresponding annotations of manipulation sequences in an auto-regressive manner. To achieve this, we send original annotations of manipulation sequences $S^{ori} \in R^{C \times N}$ (*e.g., N*=5 in Fig. 5 before a Tokenizer) into a Tokenizer, where we transform each manipulation in the sequence into one token and insert Start of Sentence (SOS) and End of Sentence (EOS) tokens into the beginning and end of sequence. After that, we obtain tokenized manipulation sequences $S^{tok} \in R^{C \times (N+2)}$ to be cross-attended with features of spatial relation $f^{spa}$. With the auto-regressive mechanism, the decoding process of facial manipulation sequence in the transformer decoder (aided by cross-attention) is triggered by SOS token and will be automatically stopped once the EOS token is predicted. In this way, we can predict facial manipulation sequences with adaptive lengths.

Normally, cross-attention between tokenized sequences $S^{tok}$ and features of spatial relation $f^{spa}$ should be performed directly. However, as mentioned above, compared to the normal image-to-sequence task, annotations of manipulation sequences are much shorter and thus less informative ($S^{tok}$ only has $(N + 2)$-length and maximum of $N$ is 5). To effectively cross-attend features of spatial relation with limited annotations of manipulation sequences, inspired by [7], we propose a sequential relation modeling with spatially enhanced cross-attention. We argue that each manipulation in $S^{tok}$ corresponds to one specific facial component/attribute which has a strong prior of spatial regions, thus we can enrich the information of manipulation sequences guided by this prior. To this end, we generate the spatial weight map for each manipulation by dynamically predicting the spatial center and scale of

each manipulation component/attribute in annotations of manipulation sequences as follows:

$$t_h, t_w = \text{sigmoid}(\text{MLP}(S^{tok})), r_h, r_w = \text{FC}(S^{tok}) \qquad (2)$$

where $t_h, t_w$ and $r_h, r_w$ are estimated 2-dimensional coordinates corresponding to spatial centers and scales of specific manipulations in the sequences, respectively. Then the Gaussian-shape spatial weight map can be generated as:

$$M(h, w) = \exp\left(-\frac{(h - t_h)^2}{\lambda r_h^2} - \frac{(w - t_w)^2}{\lambda r_w^2}\right) \qquad (3)$$

where $(h, w) \in [0, H] \times [0, W]$ are 2-dimensional coordinates of the spatial weight map $M$, and $\lambda$ is a hyper-parameter to modulate the bandwidth of the Gaussian-shape distribution. From Eq. 3, it can be seen that spatial weight map $M$ can assign higher importance to spatial regions closer to the centers and lower weights to locations farther from the centers. Moreover, as analyzed before, since diverse manipulation regions are presented for different people, the above dynamically learned scales can further tune the height/width ratios of spatial weight map based on each manipulation, contributing to a more adaptive spatial weight map. Based on this idea, we can enhance the cross-attention between features of spatial relation and annotations of manipulation sequences with generated spatial weight map $M$ as follows:

$$\begin{aligned} S &= \text{FC}(S^{tok}), K, V = \text{FC}(f^{spa}), \\ f_i^{seq} &= \text{Softmax}(K_i^T Q_i \sqrt{d} + logM)V_i, \\ f^{seq} &= \text{Concat}(f_1^{seq}, ..., f_D^{seq}) \end{aligned} \qquad (4)$$

where FC denotes a single fully-connected layer, and $f_i^{seq}$ denotes features of sequential relation. The cross-attention of the $i$-th head is further element-wise added with logarithm of spatial weight map $M$, which contributes to spatially enhanced cross-attention. Furthermore, to model the sequential relation of facial manipulation, the auto-regressive mechanism is integrated into the above cross-attention process. This is implemented by masking out (setting to $-\infty$) all values in the input of the Softmax function in Eq. 4 which correspond to illegal connections. Through concatenation of features of sequential relation from all cross-attention heads, we can obtain the final features of sequential relation $f^{seq}$ as the output of decoder.

The features of sequential relation are then fed into a Fast Forward Network (FFN) and transformed to a class score for each manipulation. Finally, we jointly train the CNN, transformer encoder and decoder by minimizing the cross-entropy loss between each class score and corresponding annotation of manipulation in the sequence.

## 5    Experiments

### 5.1    Experimental Setup

**Implementation Details.** We choose two different CNNs, ResNet-34 [10] and ResNet-50 [10] pre-trained on ImageNet [3] dataset in our paper. To be comparable in the number of parameters, we adopt a transformer model with 2 encoder and 2 decoder layers with 4 attention heads. For the training schedule, we employ 20 epochs warm-up strategy and train for 170 epochs with a learning rate drop to 10% in every 50 epochs. The initial learning rates are set as $1e-3$ for transformer part and $1e-4$ for CNN part. We set $\lambda = 4$.

**Baseline Methods.** The most straightforward solution for detecting Seq-Deepfake manipulation is to regard it as a multi-label classification problem [32]. It treats all manipulations in the sequences as independent classes and classifies the manipulated images into multiple manipulation classes. Specifically, we design a simple multi-label classification network (denoted as **Multi-Cls**) as one of the baselines. We use ResNet-34 [10] and ResNet-50 [10] pre-trained on ImageNet [3] dataset as backbones for the multi-label classification network, which is concatenated with $N$ single linear-layer branches as $N$ classification heads ($N = 5$ as maximum manipulation steps are 5 in Seq-Deepfake dataset). Moreover, we study a more complex transformer structure modified for our problem. **DETR** [2] is a popular transformer devised for end-to-end object detection. This model detects input images' bounding boxes and corresponding object classes conditioned on object queries. We revise this model by replacing the object queries with annotations of manipulation sequences and only preserve the output of object classes. Furthermore, to examine the performance of existing deepfake detection methods for our research problem, we adapt three state-of-the-art deepfake detection methods, a Dilated Residual Network variant (**DRN**) [33], a two-stream network modeling the correlation between high-frequency features and regular RGB features (**Two-Stream**) [24], and a multi-attentional deepfake detection (**MA**) [36], into multi-label classification setting. To be specific, we replace their binary label classifier with multiple classification heads to classify sequential manipulations. Please note since all of the above baselines are only able to predict the facial manipulation with fixed length ($N = 5$), 'no manipulation' class will be padded into the annotation sequence if its length is shorter than $N$ so that we can keep the same length between predictions and annotation sequences for training.

**Evaluation Metrics.** We propose two evaluation metrics for this new task.

- **Fixed Accuracy (Fixed-Acc):** Given prediction with fixed $N$-length ($N = 5$) by above baselines, as in the training process, the first type of evaluation pads 'no manipulation' class into the annotated manipulation sequences and compares each manipulation class in the predicted sequences with its corresponding annotation to calculate the evaluation accuracy.
- **Adaptive Accuracy (Adaptive-Acc):** Moreover, since the proposed method exploits sequential information to detect facial manipulation

**Table 1:** Accuracy of detecting Seq-Deepfake based on sequential facial components manipulation

| Methods | ResNet-34 | | ResNet-50 | |
|---|---|---|---|---|
| | Fixed-Acc | Adaptive-Acc | Fixed-Acc | Adaptive-Acc |
| Multi-Cls | 69.66 | 50.54 | 69.65 | 50.57 |
| DETR [2] | 69.87 | 50.63 | 69.75 | 49.84 |
| Ours | **72.13** | **54.80** | **72.65** | **55.30** |

**Table 2:** Accuracy of detecting Seq-Deepfake based on sequential facial attributes manipulation

| Methods | ResNet-34 | | ResNet-50 | |
|---|---|---|---|---|
| | Fixed-Acc | Adaptive-Acc | Fixed-Acc | Adaptive-Acc |
| Multi-Cls | 66.99 | 46.68 | 66.66 | 46.00 |
| DETR [2] | 67.93 | 48.15 | 67.62 | 47.99 |
| Ours | **67.99** | **48.32** | **68.86** | **49.63** |

sequences based on the auto-regressive mechanism, predictions will be automatically stopped once predicting the EOS token. Thus, the proposed method can detect facial manipulation sequences with adaptive lengths. To conduct the evaluation in this scenario, the second type of evaluation is devised, which compares predicted manipulations and corresponding annotations within the maximum steps of manipulations ($N \leq 5$) between them. This makes the evaluation focus more on accuracy of manipulations.

More details of two evaluation metrics can be found in **Supplementary Material**.

### 5.2   Benchmark for Seq-Deepfake

We tabulate the first benchmark for detecting sequential facial manipulation based on facial components manipulation and facial attributes manipulation in Tables 1 to 3. We note that, both baselines and the proposed method obtains much higher performance under evaluation metric Fixed-Acc than Adaptive-Acc. This validates that detecting sequential facial manipulation with adaptive lengths is much harder than its simplified version with fixed length. It can be observed from Tables 1 and 2, that the proposed SeqFakeFormer obtains the best performance of detecting facial manipulation sequences compared to all considered baselines in both facial components manipulation and facial attributes manipulation. In addition, SeqFakeFormer also performs better than other baselines with both CNNs (ResNet-34 and ResNet-50), indicating the compatibility of the proposed method with different feature extractors. Specifically, the proposed method has achieved about 3-4% improvement in facial components sequential manipulation and 1-2% improvement in facial attributes sequential manipulation under two evaluation metrics. In particular, there exists a larger performance gap between SeqFakeFormer and other baselines under evaluation metric Adaptive-Acc than Fixed-Acc, which demonstrates that the effectiveness of the proposed method is more significant in the harder case. Moreover, we tabulate the comparison between three SOTA deepfake detection methods and our method in Table 3.

**Table 3:** Accuracy of detecting Seq-Deepfake compared to deepfake detection methods

| Methods | Face Components Manipulation | | Face Attributes Manipulation | |
|---|---|---|---|---|
| | Fixed-Acc | Adaptive-Acc | Fixed-Acc | Adaptive-Acc |
| DRN [33] | 66.06 | 45.79 | 64.42 | 43.20 |
| MA [36] | 71.31 | 52.94 | 67.58 | 47.48 |
| Two-Stream [24] | 71.92 | 53.89 | 66.77 | 46.38 |
| Ours | **72.65** | **55.30** | **68.86** | **49.63** |

**Table 4:** Ablation study of detecting Seq-Deepfake based on sequential facial components manipulation

| Components | | ResNet-34 | | ResNet-50 | |
|---|---|---|---|---|---|
| Auto-regressive | SECA | Fixed-Acc | Adaptive-Acc | Fixed-Acc | Adaptive-Acc |
| ✗ | ✗ | 70.64 | 52.19 | 71.22 | 53.43 |
| ✗ | ✔ | 70.77 | 51.71 | 70.99 | 52.66 |
| ✔ | ✗ | 71.88 | 53.84 | 72.18 | 54.64 |
| ✔ | ✔ | **72.13** | **54.80** | **72.65** | **55.30** |

SeqFakeFormer also outperforms all SOTA deepfake detection methods in both manipulation types. Since all the baselines treat detecting Seq-Deepfake as a multi-label classification problem, only spatial information of manipulated images are extracted. In contrast, SeqFakeFormer is capable of exploiting both spatial and sequential manipulation traces and thus more useful sequential information can be modeled, which is the key to enhance the performance of Seq-Deepfake Detection.
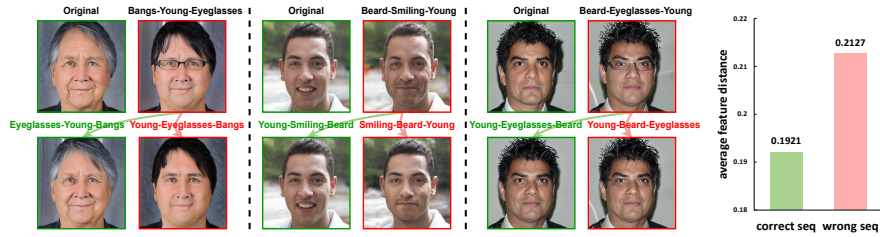
## 5.3  Ablation study

In this sub-section we investigate the impact of two key components in SeqFakeFormer, auto-regressive mechanism and Spatially Enhanced Cross-Attention module (SECA), to the overall performance. The considered components and the corresponding results obtained for each case are tabulated Tables 4 and 5. As evident from Tables 4 and 5, removing either auto-regressive mechanism or SECA will degrade the overall performance. This validates that auto-regressive mechanism facilitates the sequential relation modeling and SECA benefits the cross-attention. These components complement each other to produce better performance for detecting Seq-Deepfake.

## 5.4  Face Recovery

After detecting facial manipulation sequences, we are able to perform more challenging tasks, like recovering the original face from the manipulated face image. Specifically, we formulate the Face Recovery task as: given a sequentially manipulated face image, reverse the manipulation process to get an image as close as possible to the original image. For example, in the facial attributes manipulation case, given an image generated by sequential manipulations on different attributes on the original face, we want to recover the original image. In fact, this task can be seen as an inverse sequential facial attribute manipulation problem, which can be effectively solved by the data generation pipeline described in Section 3 in an inverse manner. Specifically, as

**Table 5:** Ablation study of detecting Seq-Deepfake based on sequential facial attributes manipulation

| Components | | ResNet-34 | | ResNet-50 | |
|---|---|---|---|---|---|
| Auto-regressive | SECA | Fixed-Acc | Adaptive-Acc | Fixed-Acc | Adaptive-Acc |
| ✗ | ✗ | 66.98 | 45.87 | 68.14 | 48.49 |
| ✗ | ✔ | 67.36 | 47.22 | 68.77 | 49.54 |
| ✔ | ✗ | 66.70 | 46.56 | 68.17 | 48.81 |
| ✔ | ✔ | **67.99** | **48.32** | **68.86** | **49.63** |



**Fig. 6:** Face recovery based on correct and wrong facial manipulation sequences.

**Fig. 7:** Identity preservation.

can be observed in Fig. 6, once we detect the correct facial manipulation sequence, *i.e.* correct manipulations ordered with correct manipulation steps, we can recover original face by performing face attribute manipulation based on the inverse order of detected facial manipulation sequence (process with green arrow). Comparatively, recovering the face image with wrongly ordered manipulation sequences may encounter different problems, such as incomplete recovery of age, smile, glasses, etc. (process with red arrow). Fig. 7 evaluates the results using identity preservation metrics as in [13], where smaller feature distance means identity is better preserved. The average feature distance between randomly selected 100 original faces and recovered faces using correct manipulation sequences is clearly smaller than that of the wrongly ordered sequence, indicating that the identity can be better recovered with correct manipulation sequence. Based on the above analysis and experiments, we prove that the detection of facial manipulation sequences is highly useful for face recovery, and we hope it can be applied to more meaningful tasks in the future.

## 6   Conclusion

This paper studies a novel research problem – Detecting Sequential DeepFake Manipulation, aiming to detect a sequential vector of multi-step facial manipulation operations. We also introduce the first Seq-DeepFake dataset to provide sequentially manipulated face images. Supported by this new dataset, we cast detecting Seq-DeepFake manipulation as a specific image-to-sequence task and propose a Seq-DeepFake Transformer (SeqFakeFormer). Two modules, Spatial Relation Extraction and Sequential Relation Modeling with Spatially Enhanced Cross-Attention, are integrated into SeqFakeFormer, complementing

each other. Extensive experimental results demonstrate the superiority of SeqFakeFormer and valuable observations pave the way for future research in broader deepfake detection.

# References

1. Bello, I., Zoph, B., Vaswani, A., Shlens, J., Le, Q.V.: Attention augmented convolutional networks. In: CVPR (2019) 8
2. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: ECCV. pp. 213–229 (2020) 11, 12
3. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: CVPR (2009) 11
4. Dolhansky, B., Howes, R., Pflaum, B., Baram, N., Ferrer, C.C.: The deepfake detection challenge (dfdc) preview dataset. arXiv preprint arXiv:1910.08854 (2019) 4
5. Durall, R., Keuper, M., Pfreundt, F.J., Keuper, J.: Unmasking deepfakes with simple features. arXiv preprint arXiv:1911.00686 (2019) 4
6. Dzanic, T., Shah, K., Witherden, F.: Fourier spectrum discrepancies in deep network generated images. NeurIPS (2020) 4
7. Gao, P., Zheng, M., Wang, X., Dai, J., Li, H.: Fast convergence of detr with spatially modulated co-attention. In: CVPR (2021) 9, 18
8. Gu, S., Bao, J., Chen, D., Wen, F.: Giqa: Generated image quality assessment. In: ECCV (2020) 6
9. Haliassos, A., Vougioukas, K., Petridis, S., Pantic, M.: Lips don't lie: A generalisable and robust approach to face forgery detection. In: CVPR (2021) 4
10. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016) 8, 11
11. He, Y., Gan, B., Chen, S., Zhou, Y., Yin, G., Song, L., Sheng, L., Shao, J., Liu, Z.: Forgerynet: A versatile benchmark for comprehensive forgery analysis. In: CVPR (2021) 2, 5
12. Jiang, L., Li, R., Wu, W., Qian, C., Loy, C.C.: Deeperforensics-1.0: A large-scale dataset for real-world face forgery detection. In: CVPR (2020) 4
13. Jiang, Y., Huang, Z., Pan, X., Loy, C.C., Liu, Z.: Talk-to-edit: Fine-grained facial editing via dialog. In: ICCV (2021) 3, 5, 6, 14, 21
14. Karras, T., Aila, T., Laine, S., Lehtinen, J.: Progressive growing of gans for improved quality, stability, and variation. In: ICLR (2018) 5
15. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: CVPR (2019) 6
16. Kim, H., Choi, Y., Kim, J., Yoo, S., Uh, Y.: Exploiting spatial dimensions of latent in gan for real-time image editing. In: CVPR (2021) 3, 5, 21
17. Lee, C.H., Liu, Z., Wu, L., Luo, P.: Maskgan: Towards diverse and interactive facial image manipulation. In: CVPR (2020) 4, 5
18. Lee, W., Kim, D., Hong, S., Lee, H.: High-fidelity synthesis with disentangled representation. In: ECCV (2020) 3, 6
19. Li, J., Xie, H., Li, J., Wang, Z., Zhang, Y.: Frequency-aware discriminative feature learning supervised by single-center loss for face forgery detection. In: CVPR (2021) 4
20. Li, L., Bao, J., Zhang, T., Yang, H., Chen, D., Wen, F., Guo, B.: Face x-ray for more general face forgery detection. In: CVPR (2020) 4

21. Li, Y., Yang, X., Sun, P., Qi, H., Lyu, S.: Celeb-df: A large-scale challenging dataset for deepfake forensics. In: CVPR (2020) 4
22. Liu, H., Li, X., Zhou, W., Chen, Y., He, Y., Xue, H., Zhang, W., Yu, N.: Spatial-phase shallow learning: rethinking face forgery detection in frequency domain. In: CVPR (2021) 4
23. Liu, Z., Luo, P., Wang, X., Tang, X.: Deep learning face attributes in the wild. In: CVPR (2015) 5
24. Luo, Y., Zhang, Y., Yan, J., Liu, W.: Generalizing face forgery detection with high-frequency features. In: CVPR (2021) 4, 11, 13
25. Parmar, N., Vaswani, A., Uszkoreit, J., Kaiser, L., Shazeer, N., Ku, A., Tran, D.: Image transformer. In: ICML (2018) 8
26. Qian, Y., Yin, G., Sheng, L., Chen, Z., Shao, J.: Thinking in frequency: Face forgery detection by mining frequency-aware clues. In: ECCV (2020) 4
27. Rossler, A., Cozzolino, D., Verdoliva, L., Riess, C., Thies, J., Nießner, M.: Faceforensics++: Learning to detect manipulated facial images. In: CVPR (2019) 2, 4, 5
28. Shen, Y., Gu, J., Tang, X., Zhou, B.: Interpreting the latent space of gans for semantic face editing. In: CVPR (2020) 3, 5, 6
29. Shen, Y., Yang, C., Tang, X., Zhou, B.: Interfacegan: Interpreting the disentangled face representation learned by gans. TMPAMI (2020) 5
30. Shen, Y., Zhou, B.: Closed-form factorization of latent semantics in gans. In: CVPR (2021) 5
31. Voynov, A., Babenko, A.: Unsupervised discovery of interpretable directions in the gan latent space. In: ICML (2020) 5
32. Wang, H., Liu, W., Bocchieri, A., Li, Y.: Can multi-label classification networks know what they don't know? NeurIPS (2021) 11
33. Wang, S.Y., Wang, O., Owens, A., Zhang, R., Efros, A.A.: Detecting photoshopped faces by scripting photoshop. In: CVPR (2019) 11, 13
34. Wang, W., Alameda-Pineda, X., Xu, D., Fua, P., Ricci, E., Sebe, N.: Every smile is unique: Landmark-guided diverse smile generation. In: CVPR (2018) 5
35. Yang, H., Huang, D., Wang, Y., Jain, A.K.: Learning face age progression: A pyramid architecture of gans. In: CVPR (2018) 5
36. Zhao, H., Zhou, W., Chen, D., Wei, T., Zhang, W., Yu, N.: Multi-attentional deepfake detection. In: CVPR (2021) 4, 11, 13
37. Zhao, T., Xu, X., Xu, M., Ding, H., Xiong, Y., Xia, W.: Learning self-consistency for deepfake detection. In: ICCV (2021) 4
38. Zhu, P., Abdal, R., Qin, Y., Wonka, P.: Sean: Image synthesis with semantic region-adaptive normalization. In: CVPR (2020) 5
39. Zhu, X., Wang, H., Fei, H., Lei, Z., Li, S.Z.: Face forgery detection by 3d decomposition. In: CVPR (2021) 4
40. Zhuang, P., Koyejo, O., Schwing, A.G.: Enjoy your editing: Controllable gans for image editing via latent space navigation. In: ICLR (2021) 5
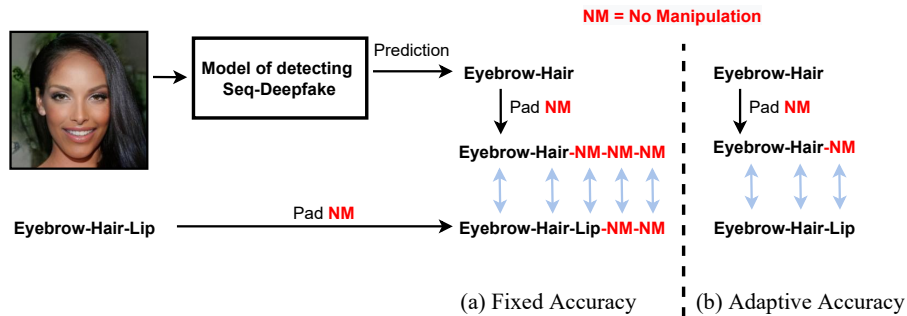
## Supplementary Material

## A    Training Details and Hyper-parameter Setting

Implementation is in PyTorch. For the training schedule, we employ a 20-epochs warm-up strategy. The initial learning rate is set as $1e-3$ for transformer part and $1e-4$ for CNN part, with a decay factor of 10 at 70 and 120 epochs, totally 170 epochs. We use the SGD momentum optimizer with weight decay set as $1e-4$. We use a mini-batch size of 32 per GPU and 4 GPUs in total. Model selection for evaluation is conducted by considering the trained model that has produced the best accuracy on the validation set.

## B    Evaluation Metrics



**Fig. 8:** Comparison between two evaluation metrics (a) Fixed Accuracy and (b) Adaptive Accuracy.

As illustrated in Fig. 8, we elaborate on two evaluation metrics proposed in the experiment for our new task.

– **Fixed Accuracy (Fixed-Acc):** As mentioned in the main paper, in the training process, since all of the baselines are only able to predict the facial manipulation with fixed length ($N = 5$), 'no manipulation' class will be padded into the annotation sequence if its length is shorter than $N$ so that we can keep the same length between predictions and annotation sequences for training. Following this strategy, as shown in Fig. 8, under the evaluation metric of Fixed Accuracy, given the prediction, such as 'Eyebrow-Hair', from the model of detecting Seq-Deepfake, we first pad 'no manipulation' class into it to form the padded prediction sequence as 'Eyebrow-Hair-NM-NM-NM' (NM means 'no manipulation' class) so that we can obtain the prediction with fixed $N$-length ($N = 5$). To keep the

same length between predictions and annotation sequences for evaluation, we pad 'no manipulation' class into the annotation of manipulation sequences as well, denoted as 'Eyebrow-Hair-Lip-NM-NM', and compare each manipulation class in the predicted sequences with its corresponding annotation to calculate the evaluation accuracy.

– **Adaptive Accuracy (Adaptive-Acc):** Moreover, since the proposed method exploits sequential information to detect facial manipulation sequences based on the auto-regressive mechanism, predictions will be automatically stopped once predicting the EOS token. Thus, the proposed method can detect facial manipulation sequences with adaptive lengths. To conduct the evaluation in this scenario, as illustrated in Fig. 8, the second type of evaluation is devised, which compares predicted manipulations and corresponding annotations within the maximum steps of manipulations ($N = 3$ in Fig. 8 and we just pad one 'no manipulation' class into prediction sequence) between them. This makes the evaluation focus more on accuracy of manipulations.
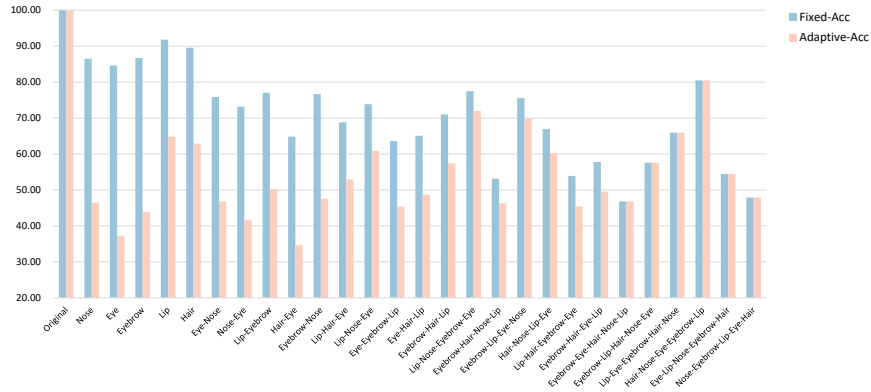
## C    Multi-head version of SECA

Similar to [7], we extend the basic version of Spatially Enhanced Cross-Attention (SECA) introduced in the main paper into multi-heads version, which enhances cross-attention features differently for different cross-attention heads. As mentioned in Eq. 2 in the main paper, the basic version of SECA estimates the 2-dimensional coordinates corresponding to spatial centers $[t_h, t_w]$. Similarly, the multi-head version of SECA estimates a head-shard spatial center $[t_h, t_w]$ and then predicts a head-specific center offset $[\triangle t_{h,i}, \triangle t_{w,i}]$ and corresponding head-specific scales $[r_{h,i}, r_{w,i}]$ for $i$-th cross-attention head. In this way, we generate $i$-th head-specific Gaussian-shape spatial weight map $M_i$ based on the $i$-th head-specific center $[t_h + \triangle t_{h,i}, t_w + \triangle t_{w,i}]$ and scales $[r_{h,i}, r_{w,i}]$ as:

$$M_i(h, w) = \exp\left(-\frac{(h - (t_h + \triangle t_{h,i}))^2}{\lambda r_{h,i}^2} - \frac{(w - (t_w + \triangle t_{w,i}))^2}{\lambda r_{w,i}^2}\right) \quad (5)$$
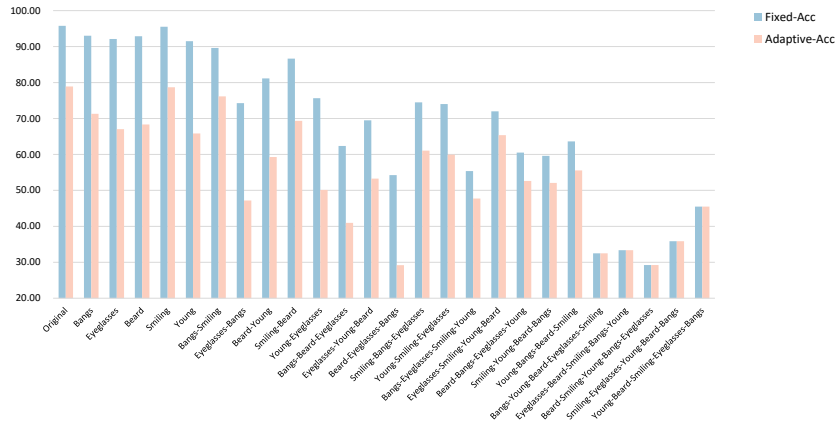
Based on this, we can calculate the features of sequential relation $f_i^{seq}$ from $i$-th cross-attention head enhanced by the $i$-th SECA as follows:

$$f_i^{seq} = \text{Softmax}(K_i^T Q_i \sqrt{d} + log M_i) V_i, \quad (6)$$

Different from basic version of SECA, above Eq. 6 shows that in the multi-head version of SECA, the cross-attention of the $i$-th head is element-wise added with logarithm of $i$-th head-specific spatial weight map $M_i$, which contributes to a more adaptive and specific enhanced cross-attention. Experiments regarding the proposed method in the main paper are all carried out based on the multi-head version of SECA.

(a) Accuracy on sequential facial components manipulation dataset



(b) Accuracy on sequential facial attributes manipulation dataset

**Fig. 9:** Accuracy for each manipulation sequence.

## D    Accuracy For Each Manipulation Sequence

As mentioned in the main paper, we generate 28 types of manipulation sequences based on facial components manipulation while 26 types of manipulation sequences based on facial attributes manipulation. To provide a more detailed analysis, in this section, we plot accuracy for each manipulation sequence in both facial manipulation methods as shown in Fig. 9. It can be observed that diverse accuracy performance are achieved for different manipulation sequences, ranging from 46.81% to 100% under Fixed-Acc and 34.69% to 100% under Adaptive-Acc in sequential facial components manipulation, while ranging from 29.25% to 95.75% under Fixed-Acc and 29.21% to 78.88% under Adaptive-Acc in sequential facial attributes
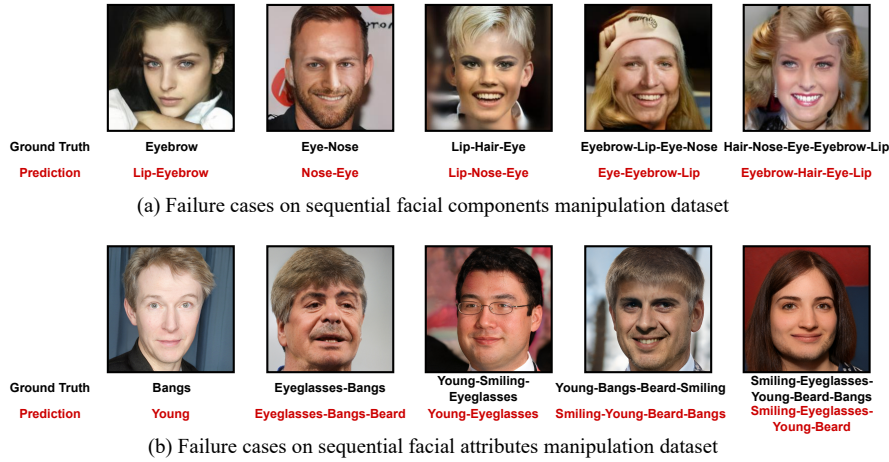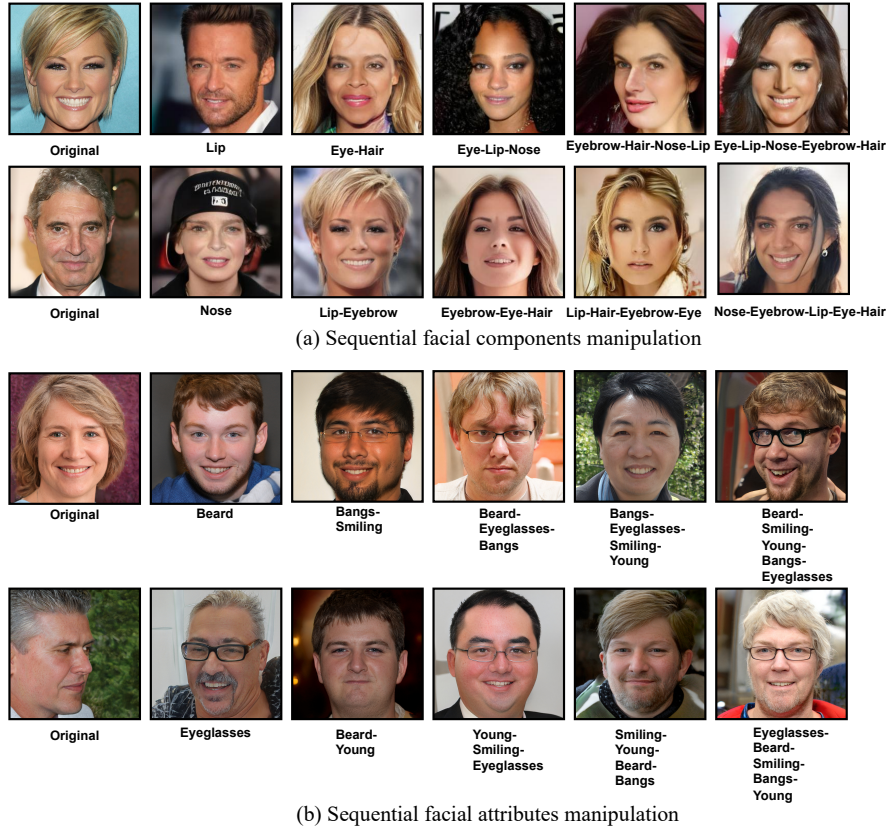
(a) Failure cases on sequential facial components manipulation dataset



(b) Failure cases on sequential facial attributes manipulation dataset

**Fig. 10:** Examples of failure cases.

manipulation. This demonstrates various manipulation sequences are challenging for detection and there exist some extremely hard cases. Therefore, we should further improve our method to cope with all types of manipulation sequences in the future. Furthermore, it can be seen from Fig. 9 that the accuracy gap between two evaluation metrics, Fixed-Acc and Adaptive-Acc, decreases along with the length of sequence increases. This is because the padded 'no manipulation' class is fewer in the longer manipulation sequence when evaluating under Adaptive-Acc, which is closer to the evaluation under Fixed-Acc.

# E    Failure Cases

To provide a deeper understanding for our novel task and method, we display some failure cases produced by the proposed method as illustrated in Fig. 10. From Fig. 10, it can be seen that there exist diverse failure cases, including wrong predictions with respect to manipulation type, sequence order, sequence length, etc. This validates that it is quite difficult for our novel research problem since we need to detect facial manipulation sequences in terms of correct manipulation types, orders and lengths simultaneously from hyper-realistic face images with subtle sequential manipulations. This motivates us to continually improve the performance of the proposed SeqFakeFormer to tackle such a challenging task in the future.

(a) Sequential facial components manipulation



(b) Sequential facial attributes manipulation

**Fig. 11:** Illustration of Seq-Deepfake dataset. Samples of Seq-Deepfake are provided with annotations of manipulation sequences.

## F   Sequential Deepfake Dataset

We display more samples from the generated large-scale Sequential Deepfake (Seq-Deepfake) dataset in Fig. 11. As shown in Fig. 11, based on two different facial manipulation methods, facial components manipulation [16] and facial attributes manipulation [13], various sequential facial manipulations are produced with diverse manipulation steps, expressions, ages, and genders.

## G   Face Recovery

Fig. 12 shows more samples regarding the face recovery based on correct and wrong facial manipulation sequences. As can be observed in Fig. 12, once we detect the correct facial manipulation sequence, *i.e.* correct manipulations ordered with correct manipulation steps, we can recover original face by
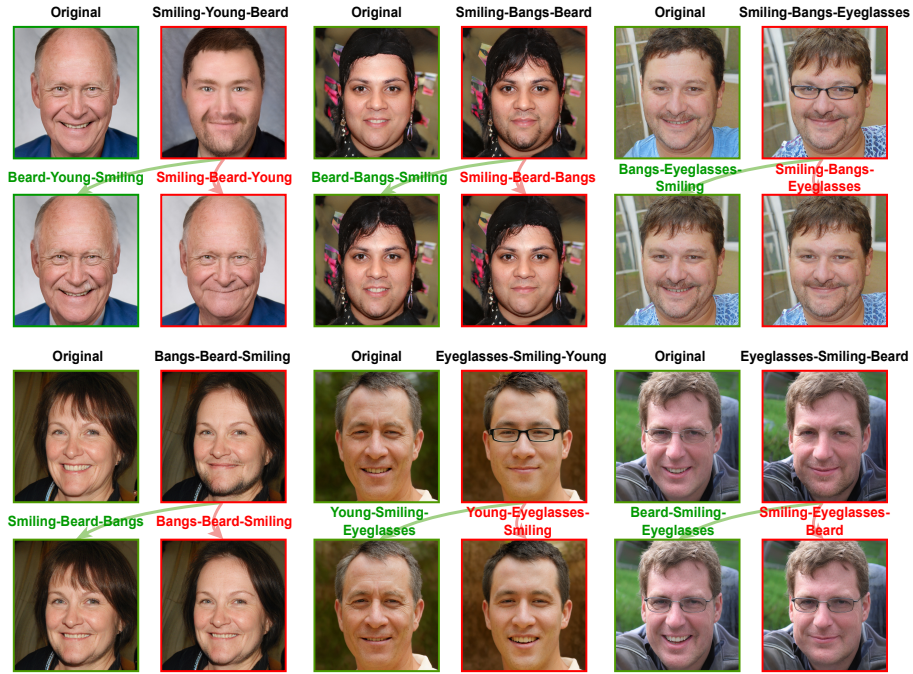
**Fig. 12:** Face recovery based on correct and wrong facial manipulation sequences.

performing face attribute manipulation based on the inverse order of detected facial manipulation sequence (process with green arrow). In contrast, recovering the face image with wrongly ordered manipulation sequences may encounter different problems, such as incomplete recovery of age, smile, glasses, bangs, etc. (process with red arrow).