

# Facelet-Bank for Fast Portrait Manipulation

Ying-Cong Chen<sup>1</sup> Huaijia Lin<sup>1</sup> Ruiyu Li<sup>1</sup> Xin Tao<sup>1</sup>  
 Michelle Shu<sup>3</sup> Yangang Ye<sup>2</sup> Xiaoyong Shen<sup>2</sup> Jiaya Jia<sup>1,2</sup>

<sup>1</sup>The Chinese University of Hong Kong <sup>2</sup>Tencent YouTu Lab <sup>3</sup>Johns Hopkins University

{ycchen, linhj, ryli, xtao}@cse.cuhk.edu.hk Mshu1@jhu.edu

yangangye@tencent.com goodshenxy@gmail.com leo{jia9@gmail.com}

## Abstract

Digital face manipulation has become a popular and fascinating way to touch images due to the prevalence of smart phones and social networks. With a wide variety of user preferences, facial expressions and accessories, a general and flexible model is necessary to accommodate different types of facial editing. In this paper, we propose such a model based on an end-to-end convolutional neural network that supports fast inference, edit-effect control, and quick partial-model update. In addition, this model learns from unpaired image sets with different attributes. Experimental results show that our framework can handle a wide range of face expressions, accessories, and makeup effects, producing high-resolution and high-quality results in high speed.

## 1. Introduction

Digital face manipulation aims to change semantically expressive and meaningful attributes, such as smiling and mourning, or add virtual makeup/accessories to human faces, including mustache and eyeglasses. With the increasing popularity of smart phones and digital cameras, the demand of a practical and fast system rises drastically. Face manipulation has attracted great interests in computer vision and graphics [14, 3, 6, 4, 33, 31, 28]. Previous works have been devoted to face beautification [19, 8], de-beautification [10], expression manipulation [28] and age progression [14], to name a few.

Given these many solutions, it is common knowledge that different face makeup or attribute changes require special manipulation operations. For example, face beautification or de-beautification processes skin color and texture while face expression manipulation focuses more on 2D or 3D geometry. With this fact, most existing approaches are developed particularly for individual tasks, and such specialization requires many expert efforts and domain knowledge to establish new methods for effect generation.



Figure 1. Illustration of face manipulation effect with our model.

In the following, we elaborate on the underlying problems of face manipulation when seeking a data-driven framework to unify many face effects. It is followed by introducing the intriguing work to construct our system with this pursuit.

### 1.1. Possible Solutions and Their Problems

**Direct Regression** The most straightforward way to learn face editing operations from external data is to directly regress the input (before edit) and the ground truth image after edit [10, 5]. However, this procedure requires labelled paired data, which in many cases is unavailable or requires intensive human labor to create. For any effects that does not exist in real world, these operations cannot be used in the system.

**Generative Adversarial Network** Recently, Generative Adversarial Network (GAN) has shown its capacity in set-to-set unsupervised learning [36]. It uses a cycle consistent loss to preserve image content, and an adversarial loss to transfer the attribute of one set to another. Although the concept is neat and results are quite phenomenal, it is difficult to train especially for new effects that require part system modification. It has to keep a balance of generator and discriminator. Also we notice non-optimal training

would result in unrealistic manipulation, easily noticeable on visually-sensitive human faces.

**Deep Feature Interpolation** Deep feature interpolation [29] provides another solution to learn image attribute change from two different sets. It is based on the deep feature of two image sets. However, it is not an end-to-end framework and thus cannot be optimized globally. In addition, it is compute-intensive even during testing because of its heavy involvement of hundreds of face warping and convolution operations.

## 1.2. Our Solution

We pursue a general, flexible and high-quality-output network for face manipulation. Fig. 1 shows the result generated by our method. Our work follows the an encoder-decoder architecture rather than the recently popular GAN. Inspired by the Style-Bank [9] that learns replaceable style transfer layers, we propose a Facelet-Bank framework that models face effects with respective middle-level convolutional layers. So interestingly, in order to generate different effects, instead of redesigning the framework completely, only these middle-level convolutional layers need to be updated.

Also, considering the lack of ground truth data for many face manipulation tasks, we instead leverage the result of [29] to produce pseudo targets to learn the facelet-bank layers. The local receptive field of our facelet bank naturally provides regularization, and thus it can capture the correct relation between visual patterns and certain network operations, despite the fact that the pseudo targets are usually noisy.

Finally, we show that these layers can automatically attend to the most important regions so that face manipulation can be performed in an end-to-end fashion. Our method is specially designed to allow users to control the level of effect, and thus it enables interactive face manipulation. Our overall contribution is multifold.

- We propose a set-to-set CNN framework for face manipulation. It does not requires paired data in training;
- The framework is flexible to generate different effects and their levels by simply updating a few convolutional layers, which makes the system developer-friendly.
- Our method naturally benefits from the local prior of convolutional layers, which regularizes noisy labels.
- Experiments show that our approach can handle a wide range of face effects in high speed.

## 2. Related Work

**Face Editing** Our work can be categorized as face editing or manipulation, which has been extensively studied in computer vision and graphics [31, 3, 33, 4, 6, 14]. Traditional face edit includes face relighting [31, 3], expression editing [33], face morphing [4], attribute editing [6], and face aging [14]. However, those models are designed for specific tasks, and rely heavily on domain knowledge. Ours method differs from them as it is a data-driven framework designed to handle general face effects and manipulations.

**Image Attribute Manipulation** Altering semantic attributes of images is an important topic. Image analogy [11, 7, 30] transfers the attribute or appearance changes from a pair of images to a new one. Recently, generative adversarial networks [12, 36, 25, 27, 22] have become popular in image generation and image-to-image transform. Specifically, the method in [12] learns attribute transformation from two sets of paired data. This idea is further generalized to handle unpaired data [36, 27, 25]. These approaches alter the attributes by training a generator to simulate an online updating domain discriminator. However, when designing and training a system, the balance between generator and discriminator is difficult to deal with. Scaling up image sizes would make training even more difficult. Besides, these approaches do not allow control over the degree of editing, e.g., how bright the face is, which is very important to real-world applications.

**Feature Interpolation** Our work is also related to feature interpolation [17, 23, 29, 32, 16, 34]. Variational autoencoder (VAE) [17, 34] learns a latent space such that image manipulation can be done with simple arithmetic computation in the deep space. But it implicitly assumes that the target attribute is well isolated. Given unpaired data for face manipulation, this assumption is usually not satisfied. The method of [29] alleviates this problem by proposing a K-NN based approach to isolate the target attribute, yet only a subset of samples are used for estimation, which leads to undesired artifact in resulting images. In addition, the computational cost is inevitably high, which may be intolerable for fast face manipulations. Our work is different from these approaches. We use convolutional neural networks to learn the shifting direction, so that each transformation is associated with a corresponding visual pattern. We show that the resulting direction is better isolated towards the target attribute. In addition, our method is an end-to-end framework that supports fast inference.

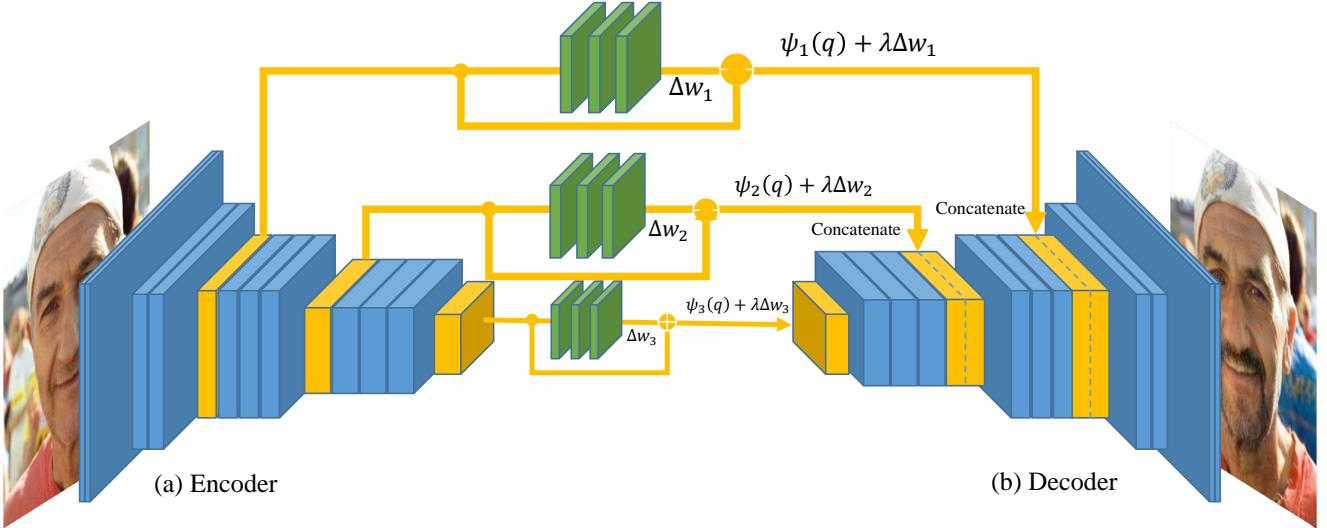


Figure 2. Illustration of our framework. (a) is the encoder  $\mathcal{E}(\cdot)$ ; (b) are convolutional layers of our facelet bank  $\mathcal{V}(\cdot)$ ; (c) is the decoder  $\mathcal{D}(\cdot)$ . The structure of our facelet bank is Conv-ReLU-Conv-ReLU-Conv, where all Convs are with  $3 \times 3$  kernels. Also, all Convs of the facelet bank do not change the height, width and number of channels given the previous input.

### 3. Proposed Method

Suppose there are two face domains  $\mathcal{X}$  and  $\mathcal{Y}$  with different properties. Our goal is to shift images of domain  $\mathcal{X}$  towards domain  $\mathcal{Y}$  without guidance from any paired data. As for face manipulations, we want the algorithm to produce intermediate results so that users can control the strength of operations, such as smoothness of the skin.

This process can be represented as

$$z = \mathcal{O}(x, \lambda), \quad (1)$$

where  $x$  and  $z$  are the input and output images,  $\lambda$  controls the strength, and  $\mathcal{O}$  denotes the operation that changes  $x$  from domain  $\mathcal{X}$  toward domain  $\mathcal{Y}$ .

Note that domains  $\mathcal{X}$  and  $\mathcal{Y}$  usually differ *semantically*, so  $\mathcal{O}$  can be highly complex. In order to simplify  $\mathcal{O}$ , we instead define the operation  $\mathcal{O}'$  in a deep space, and change Eq. (1) to

$$\psi(z) = \mathcal{O}'(\psi(x), \lambda), \quad (2)$$

where  $\psi(\cdot)$  denotes the deep space. To one extreme, if  $\psi(\cdot)$  captures enough rich semantic information,  $\mathcal{O}'$  can be further simplified as linear shift in the deep space [20], which leads to

$$\psi(z) = \psi(x) + \lambda\Delta v, \quad (3)$$

where  $\Delta v$  denotes the direction from domain  $\mathcal{X}$  to domain  $\mathcal{Y}$  in the deep space.

As indicated in [35], a network trained on large-scale data captures semantic attributes in the convolutional layers. This idea is validated in the work of [20, 29], which uses a VGG network [26] to encode the semantic information for deep feature interpolation or style transfer. We

follow this idea to use the 5 convolutional layers to construct  $\psi(\cdot)$ . To represent face attributes at different levels, ReLU3\_1, ReLU4\_1 and ReLU5\_1 are jointly used to represent  $\psi(\cdot)$ .

**Reversing  $\psi(\cdot)$**  Following [20, 9], we train a fixed decoder network  $\mathcal{D}$  to decode  $\psi(z)$  to obtain the final output. The decoder  $\mathcal{D}$  has reverse architecture of the encoder  $\mathcal{E}$ , except for concatenating the manipulated features to the corresponding layers. We use the following loss function to train  $\mathcal{D}$ :

$$\mathcal{L} = \sum_{i=1}^n \|z_i - x_i\|_2^2 + \omega \sum_{i=1}^n \|\mathcal{E}(x_i) - \mathcal{E}(\mathcal{D}(\mathcal{E}(x_i)))\|_2^2, \quad (4)$$

where the first and second terms impose consistencies in pixel and feature space respectively, and  $\omega$  is the weight to balance the two losses.<sup>1</sup>

**Overall Structure** The overall network architecture is illustrated in Fig. 2, which is composed of an encoder  $\mathcal{E}$  that transforms images to a deep space, a ConvNet  $\mathcal{V}$  that estimates the domain direction shift  $\Delta v$ , and a decoder that transforms the operated deep feature back to an image. With this effective pipeline, the encoder and decoder are responsible for general operations, while  $\mathcal{V}$  is the key component

<sup>1</sup>We also observe that pretraining the decoder with  $\mathcal{L} = \sum_{i=1}^n \|(\mathcal{E}(x_i) + \lambda\Delta v_i) - \mathcal{E}(\mathcal{D}(\mathcal{E}(x_i) + \lambda\Delta v_i))\|_2^2$  are helpful, where  $\Delta v_i$  are pseudo labels of different attributes computed with Eq (6).

that determines the specific effect required for a face manipulation.

Now, realizing differences across all types of face effects, from a smile on the face to adding mustache, no longer needs frequent redesign of the framework. By changing different  $\mathcal{V}$ s, these apparently dissimilar face operations can be accomplished accordingly. Thus, we name the collection of  $\mathcal{V}$  as **Facelet Bank**, representing the collection of different face effects we can achieve.

### 3.1. Learning Facelet Bank

Estimating  $\Delta v$  requires cancellations of all factors except for the target attributes. Without paired data, a straightforward approach is to compute the average difference between  $\mathcal{X}$  and  $\mathcal{Y}$  [17, 34], namely

$$\Delta v \leftarrow \frac{1}{m} \sum_{i=1}^m \psi(y_i) - \frac{1}{n} \sum_{i=1}^n \psi(x_i), \quad (5)$$

where  $m$  and  $n$  are the training samples of domains  $\mathcal{X}$  and  $\mathcal{Y}$  respectively. However, this implies a strong assumption that all training samples have similar attributes to the queried one. When this assumption is not satisfied, the model does not produce realistic results.

**Generating Pseudo Labels** We relax this assumption by adopting a query-adaptive attribute transform model  $\mathcal{V}(x)$ , where  $x$  is a query sample.

Note that learning  $\mathcal{V}(\cdot)$  is difficult as we do not have paired-data to infer  $\Delta v$  for each training sample. Inspired by the method in [29], we average the neighbors of each training sample of *both* domains to construct the attribute vector. Specifically, for sample  $x$ , the corresponding  $\Delta v$  is computed as

$$\Delta v^*(x) = \frac{1}{K} \sum_{i \in N_y^K(x)} \psi(y_i) - \frac{1}{K} \sum_{i \in N_x^K(x)} \psi(x_i), \quad (6)$$

where  $N_y^K(x)$  and  $N_x^K(x)$  refer to  $K$  nearest neighbors of  $x$  among sets  $\mathcal{Y}$  and  $\mathcal{X}$  respectively. To reduce the influence of pose, viewpoint and rotation, all face images are aligned to a frontal face template in advance<sup>2</sup>. Here the “average” operation suppresses undesired noise and thus  $\Delta v^*$  tends to cancel other factors except those that divide domains  $\mathcal{X}$  and  $\mathcal{Y}$ .

**Network Design** Since  $\Delta v^*$  relates strongly to the deep feature, we reuse the extracted deep feature by setting the input of  $\mathcal{V}(\cdot)$  to  $\mathcal{V}(\mathcal{E}(q))$ . We stack 3 fully convolutional

<sup>2</sup>A face alignment tool [13] in DLIB [2] is used in our experiment.

layers with ReLU activation to capture the non-linearity relation. In this way,  $\mathcal{V}$  can be learned as

$$\mathcal{L}_{\mathcal{V}} = \sum_{i=1}^n \|\mathcal{V}(\mathcal{E}(x_i)) - \Delta v^*(x_i)\|_2^2, \quad (7)$$

where  $x_i$  is training samples of set  $\mathcal{X}$ .  $\Delta v^*(x_i)$  is defined in Eq. (6).

### 3.2. More Analysis of the Facelet Bank

Although our facelet bank learns with the help of pseudo paired data  $\Delta v^*$  defined in Eq. (7), it has additional advantages compared with directly using  $\Delta v^*$ .

**Noise Resistance** In practice, our facelet bank consistently outperforms the pseudo label  $\Delta v^*$ . To demonstrate it, we compute the heat map of  $\Delta v$  as

$$H_{i,j} = \sum_{k=1}^c v_{i,j,k}^2, \quad (8)$$

where  $i = \{1, 2, \dots, h\}$ ,  $j = \{1, 2, \dots, w\}$ , and  $k = \{1, 2, \dots, c\}$  denotes the row, column and channel indexes respectively. Fig. 3 shows the heat map of “add beard” operation. Intuitively, adding beard is only relevant to the mouth region, while other place should not be updated. However, the pseudo target  $\Delta v^*$  activates in many regions as shown in the heat map. This may cause undesirable changes in wrong places as shown in Fig 3(c). Note that the pseudo label is computed by only a subset of samples<sup>3</sup>, not sufficient to cancel all noise. On the contrary, our facelet bank has a much cleaner heat map, which only activates in the correct region.

Although this may look counterintuitive at first, it will make sense if we consider the built-in regularization of the convolutional neural network. Note that CNN has local receptive fields, which force the system to capture the relation between certain visual pattern (the mouth in this case) and its corresponding manipulation operation (“add beard” in this case). By training on a large set of samples, unrelated activations can be well suppressed, and only the relevant region (mouth in this case) is activated.

**Implicit Attention Mechanism** Another advantage of our Facelet Bank is the implicit attention mechanism. Note that computing the pseudo target requires faces to be aligned in advance. This is not needed for our facelet bank during testing. As shown in Fig. 4, the trained operations

<sup>3</sup>2K = 200 samples are used in our method, which is less than 1% compared with the training set.

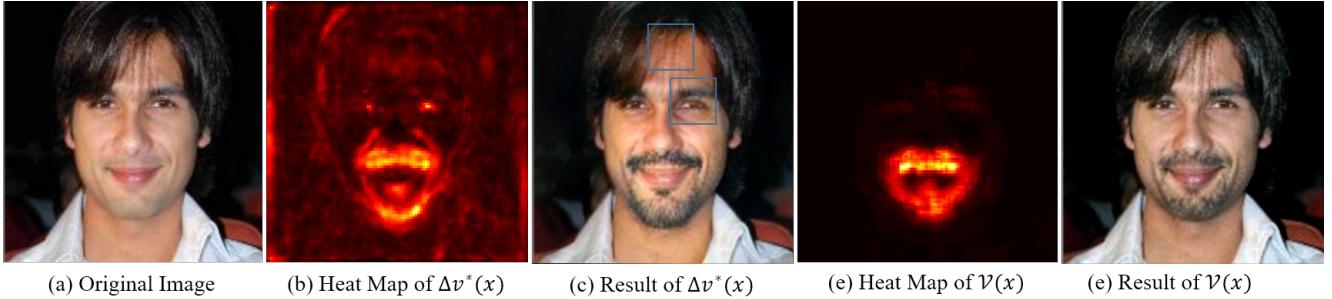


Figure 3. Illustration of the noise resistance effect. (a) The original image. (b) The heat map of pseudo shifting direction computed by Eq. (3). The blue rectangles mark the undesired change regions. (c) Corresponding result of Eq. (3). (d) Heat map of our estimated direction shift. (e) Our result.



Figure 4. Visualization of the attention region. The attention mask is computed by Eq. (8). (a), (b) and (c) correspond to the operations of adding beards, making smiles and changing ages. Note that for the beard effect, our facelet bank focuses on the mouth area; for the smile effect, it attends to smile-related facial muscles; as for the age changing effect, the attention region covers the whole face. These results match our intuition.

are adaptive to the position and pose of the face, and keep staying in the correct region. This is another major advantage stemming from the fully convolution architecture. Essentially, convolution layers of our facelet bank capture the relation of certain visual patterns and corresponding  $\Delta v$ , which are not dependent of the position. As a result, the facelet bank is capable of detecting certain patterns (e.g., mouth for beard), and add the effects accordingly.

**Speed** Calculating  $\Delta v^*$  with Eq. (3) is computationally expensive. The cost is mostly on face warping, deep feature extraction and nearest-neighboring search. The total testing time for an  $448 \times 448$  image is 1.09 minute<sup>4</sup>. But for our Facelet Bank, computing  $V(x)$  only requires network forwarding, which achieves only 0.0194 seconds.

<sup>4</sup>The system is implemented with Pytorch, and run on server with a Titan XP GPU and a E5-2623 v4 CPU.

**Advantage over CycleGAN** Comparing with cycleGAN [36], our framework has the following advantages. First, it is very easy to train. Compared with cycleGAN that jointly optimizes two generators and two discriminators, our system only needs to learn the convolutional layers of facelet bank, which is both easier and faster.

Second, our method can handle higher-resolution images. Note that the original cycleGAN has demonstrated its effectiveness on  $256 \times 256$  images, whereas our method handles much higher resolutions under the same time budget. Third, our work allows changes the level (or strength) of effect conveniently by modifying the value of  $\lambda$  in Eq. (3). In contrast, changing the effect strength of cycleGAN requires the use of different training data and re-training of the model.

## 4. Experiments

**Implementation Details** We implemented our model using PyTorch [1]. Before training the facelet bank, we train our decoder with Adam optimizer [15] with an initial learning rate of 0.0001 and step weight decay. After training, the decoder is fixed. The face effects are trained separately. Each corresponds to a set of convolution layers of the facelet bank. Adam optimizer with default hyper-parameters are used to train the facelet-bank layers.

**Dataset** Celeba [21] is a large face dataset that contains 202,599 images belonging to 10,177 identities. We randomly sample 90% for training and the rest are used for testing. The landmark and attribute information is used for computing the pseudo label  $\Delta v^*$  computed by Eq. (3). but it is not used during testing. In addition to the Celeba data, we also use the Portrait [24] and Helen [18] datasets for testing. The images of [24] and [18] are collected on Flickr with varying ages, colors, clothing, hair styles, etc. We use them to validate the cross-dataset generalization ability.

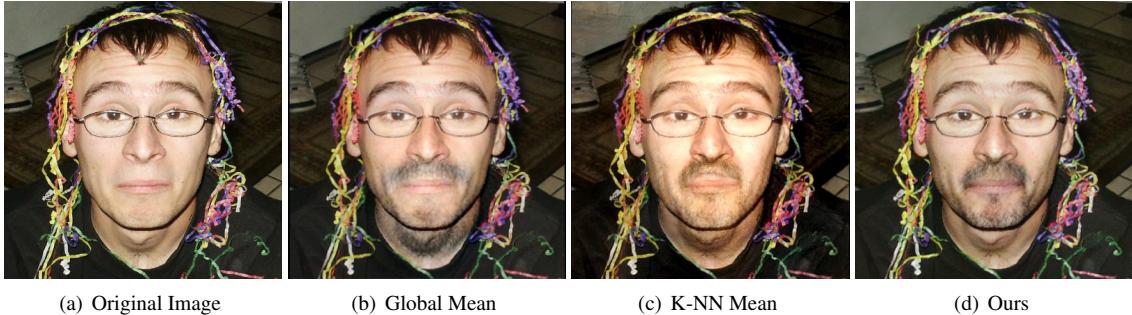


Figure 5. Comparing our facelet-bank method with baseline approaches. The rectangle highlights the region with improper changes. **Please zoom in to see the details.**

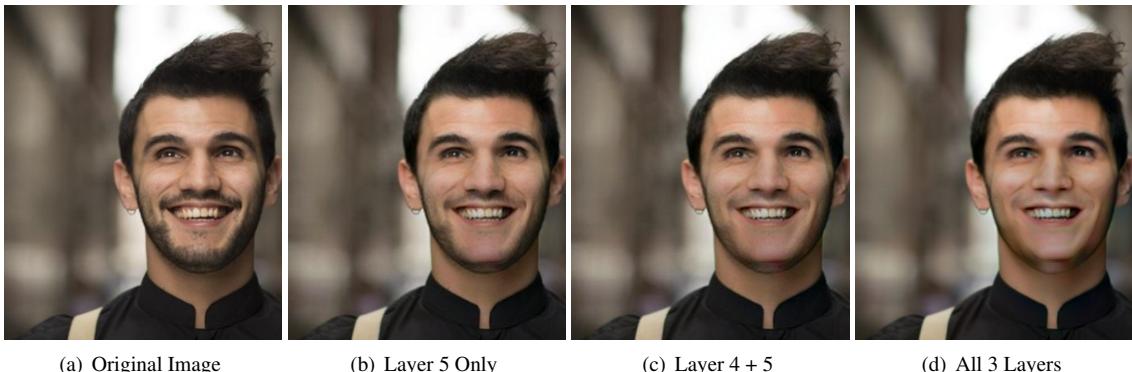


Figure 6. Results of removing facial hair. (a) the original image; (b), (c) and (d) are the result of using layer 5, layer 5+layer 4 and all three layers respectively. **Please zoom up to see details.**

#### 4.1. Evaluation of Our Approach

The critical part of our approach is the facelet bank. In the following, we will split the evaluation of this design into three parts: the effectiveness of the facelet bank, the reason to use multi-layer aggregation, and the flexibility in edit strength control.

##### 4.1.1 Effectiveness of Facelet Bank

Estimation of attribute direction shift  $\Delta v$  is critical to the result of a manipulation. We compare our facelet-bank solution with several baseline approaches. The simplest way to estimate  $\Delta v$  is to average both positive and negative samples, which is used in [17, 23, 32, 34]. As shown in Fig. 5, the beard position could be wrong since the globally computed  $\Delta v$  is not adaptive to the query face. If the face pose is far from the average of the dataset, it fails inevitably. Directly applying  $\Delta v^*$  computed by Eq. (6) alleviates this problem. However, in this case, the skin color is changed since color factor is not fully canceled by the query’s nearest neighbors. Our method learns the relation between semantic visual pattern (mouth) and the corresponding beard add-up operation, which has smaller influence on the other regions.



Figure 7. Adding facial hair for woman; (a) is the original image; (b) is the result of DFI [29]; (c) is our result.

##### 4.1.2 Effectiveness of Multi-layer Aggregation

We found that different levels of layers are complementary to one another, which is consistent with [29]. Therefore, using all of them jointly could yield better results. Fig. 6 shows the result using different layers for the “remove facial hair” task. Using only layer 5 removes a large proportion of facial hair, but a few smaller parts still remain. By incorporating layer 4 and layer 3, the amount of facial hair gradually decreases and is eventually removed completely.

##### 4.1.3 Different Operation Strength

Our model provides a fast and convenient way to control the strength of operations. Both deep feature  $\mathcal{E}(x)$  and the

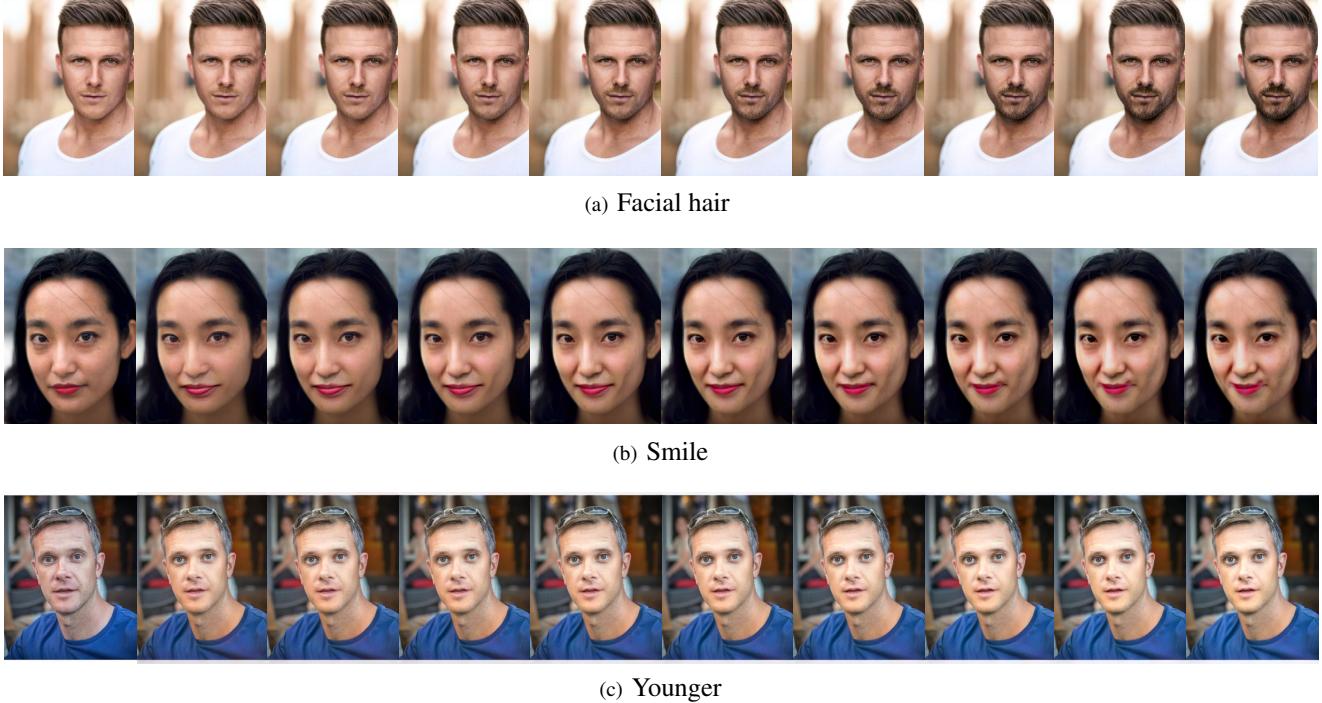


Figure 8. Illustration of different edit strength. (a), (b) and (c) show the results of different edit strength respectively.

estimated attribute direction shift,  $\Delta v$ , are computed only once. After that, changing the strength only requires to forward  $\mathcal{E} + \lambda \Delta v$  to the decoder network. This takes only 9 ms on our server. We believe that it will continue to maintain a high speed on mobile device. We have also shown results of controlling different edit strength in Fig. 8. For each case, we use strength of 0, 0.2, 0.4, 0.6, 0.8, 1, 1.2, 1.4, 1.6, 1.8 to edit the original image. Fig. 8 shows the resulting images accordingly. When  $\lambda > 1$ , the edit can be seen as extrapolation rather than interpolation. This adds stronger edits to the original image, but sometimes yields unnatural result.

#### 4.2. Comparison with State-of-the-Arts

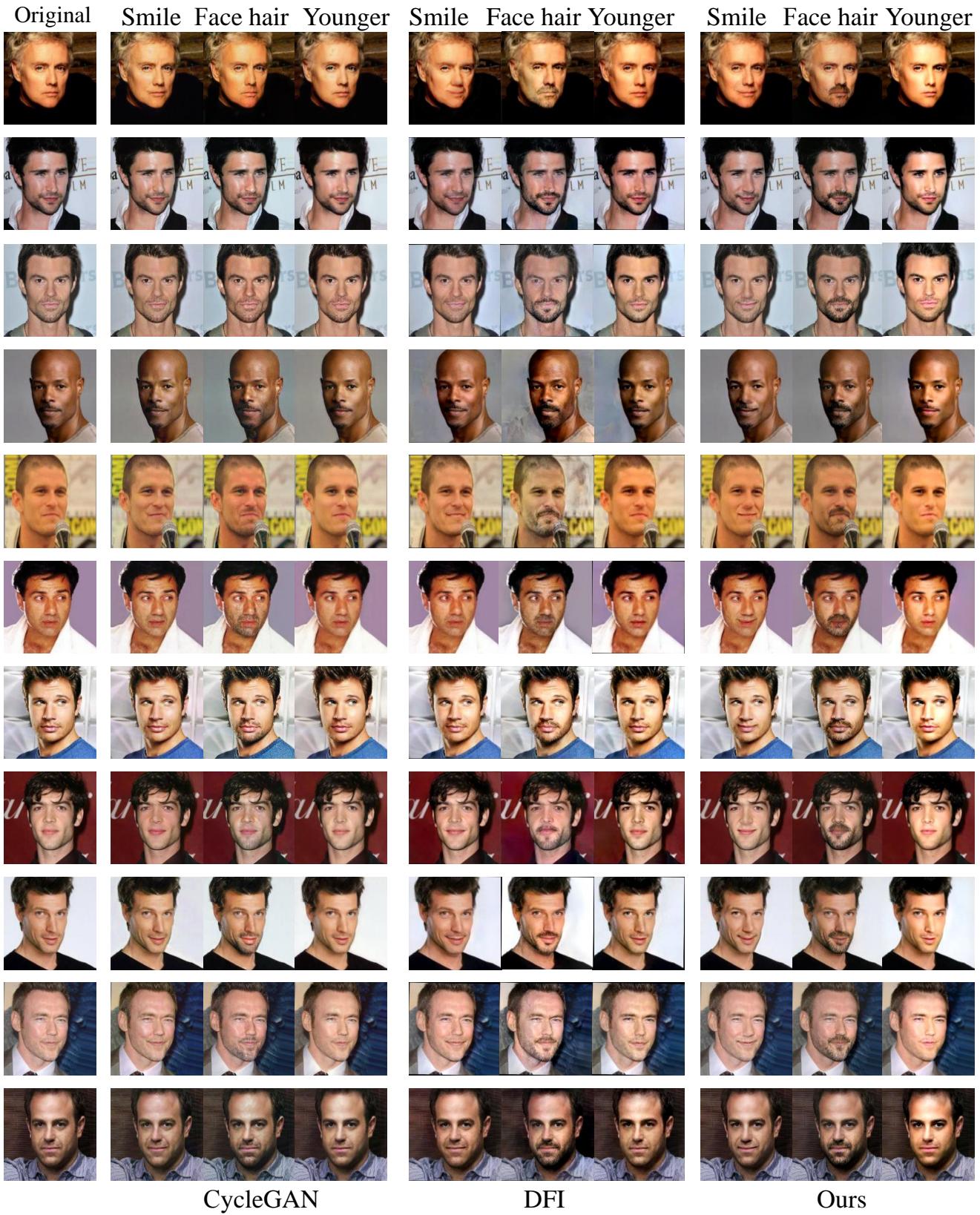
Fig. 9 shows an overall comparison with deep feature interpolation (DFI) [29] and cycleGAN [36], which can also perform set-to-set image attribute transform. Three effects are tested, which includes smiling, adding facial hair, and getting younger. Generally, our method can achieve better results than those of CycleGAN, i.e. achieving stronger effect without introducing too much artifact.

Compared with the DFI, our approach achieves better performance on local effects (e.g., facial hair), since our approach introduces less irrelevant changes, and performs comparably on other effects. However, our method is much lighter than DFI. It does not require face warping, nearest-neighboring computation and backward optimization. As a result, it has achieves 3371 times ( 65.4 second / 0.0194 second) faster than DFI.

**Disentangling Correlated Attributes** We compare the results of “adding beard to woman” in Fig. 7. Intriguingly, our model has done a much better job than DFI did. DFI unavoidably transforms the subject to a man, while our approach maintains most of the subject’s feminine features. Note that DFI intrinsically is an instance-based method, which relies heavily on training samples. Since “beard” always comes with man, it cannot disentangle these two attributes. In contrast, our method captures the relation between “beard” and “mouth”. Since the former is a masculine trait and the latter can be found in all genders, our model can naturally disentangle other similar relations, including the comparison of “beard” with “gender”. Thus, in the above case, our model has better performance.

#### 5. Concluding Remarks

In this paper, we have proposed a general framework for face manipulation. Our framework can learn face attribute shift from two image sets without any paired-example information; thus, it does not need intensive human efforts for labeling. In addition, our framework is highly flexible – each manipulating operation is related to only a few computed convolutional layers. We have proved in our experiments that it yields superior results when compared with other set-to-set image translation models.



## References

- [1] <http://pytorch.org>. 5
- [2] <https://pypi.python.org/pypi/dlib>. 4
- [3] V. Blanz, C. Basso, T. Poggio, and T. Vetter. Reanimating faces in images and video. In *Comput. Graph. Forum*, 2003. 1, 2
- [4] V. Blanz and T. Vetter. A morphable model for the synthesis of 3d faces. In *Ann. Conf. Comput. Graph. Interactive Techn.*, 1999. 1, 2
- [5] A. Brock, T. Lim, J. M. Ritchie, and N. Weston. Neural photo editing with introspective adversarial networks. *arXiv preprint arXiv:1609.07093*, 2016. 1
- [6] C. Cao, Y. Weng, S. Zhou, Y. Tong, and K. Zhou. Facewarehouse: A 3d facial expression database for visual computing. *Trans. Vis. Comput. Graph.*, 2014. 1, 2
- [7] T. Cao, C. Zach, S. Modla, D. Powell, K. Czumek, and M. Niethammer. Registration for correlative microscopy using image analogies. In *Int. Workshop Biomed. Image Registration*, 2012. 2
- [8] C.-W. Chen, D.-Y. Huang, and C.-S. Fuh. Automatic skin color beautification. In *Int. Conf. Arts Technol.*, 2009. 1
- [9] D. Chen, L. Yuan, J. Liao, N. Yu, and G. Hua. Stylebank: An explicit representation for neural image style transfer. *arXiv preprint arXiv:1703.09210*, 2017. 2, 3
- [10] Y.-C. Chen, X. Shen, and J. Jia. Makeup-go: Blind reversion of portrait edit. In *ICCV*, 2017. 1
- [11] A. Hertzmann, C. E. Jacobs, N. Oliver, B. Curless, and D. H. Salesin. Image analogies. In *Ann. Conf. Comput. Graph. Interactive Techn.*, 2001. 2
- [12] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. *arXiv preprint arXiv:1611.07004*, 2016. 2
- [13] V. Kazemi and J. Sullivan. One millisecond face alignment with an ensemble of regression trees. In *CVPR*, 2014. 4
- [14] I. Kemelmacher-Shlizerman, S. Suwajanakorn, and S. M. Seitz. Illumination-aware age progression. In *CVPR*, 2014. 1, 2
- [15] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [16] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2
- [17] A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther. Autoencoding beyond pixels using a learned similarity metric. *arXiv preprint arXiv:1512.09300*, 2015. 2, 4, 6
- [18] V. Le, J. Brandt, Z. Lin, L. Bourdev, and T. S. Huang. Interactive facial feature localization. In *ECCV*, 2012. 5
- [19] T. Leyvand, D. Cohen-Or, G. Dror, and D. Lischinski. Digital face beautification. In *Siggraph*, 2006. 1
- [20] Y. Li, C. Fang, J. Yang, Z. Wang, X. Lu, and M.-H. Yang. Universal style transfer via feature transforms. *arXiv preprint arXiv:1705.08086*, 2017. 3
- [21] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *ICCV*, 2015. 5
- [22] G. Perarnau, J. van de Weijer, B. Raducanu, and J. M. Álvarez. Invertible conditional gans for image editing. *arXiv preprint arXiv:1611.06355*, 2016. 2
- [23] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. 2, 6
- [24] X. Shen, X. Tao, H. Gao, C. Zhou, and J. Jia. Deep automatic portrait matting. In *ECCV*, 2016. 5
- [25] Z. Shu, E. Yumer, S. Hadap, K. Sunkavalli, E. Shechtman, and D. Samaras. Neural face editing with intrinsic image disentangling. In *CVPR*, 2017. 2
- [26] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 3
- [27] Y. Taigman, A. Polyak, and L. Wolf. Unsupervised cross-domain image generation. *arXiv preprint arXiv:1611.02200*, 2016. 2
- [28] J. Thies, M. Zollhöfer, M. Nießner, L. Valgaerts, M. Stamminger, and C. Theobalt. Real-time expression transfer for facial reenactment. *TOG*, 2015. 1
- [29] P. Upchurch, J. Gardner, K. Bala, R. Pless, N. Snavely, and K. Weinberger. Deep feature interpolation for image content changes. *arXiv preprint arXiv:1611.05507*, 2016. 2, 3, 4, 6, 7, 8
- [30] G. Wang, T.-T. Wong, and P.-A. Heng. Deringing cartoons by image analogies. *TOG*, 2006. 2
- [31] Y. Wang, L. Zhang, Z. Liu, G. Hua, Z. Wen, Z. Zhang, and D. Samaras. Face relighting from a single image under arbitrary unknown lighting conditions. *TPAMI*, 2009. 1, 2
- [32] X. Yan, J. Yang, K. Sohn, and H. Lee. Attribute2image: Conditional image generation from visual attributes. In *ECCV*, 2016. 2, 6
- [33] F. Yang, J. Wang, E. Shechtman, L. Bourdev, and D. Metaxas. Expression flow for 3d-aware face component transfer. 2011. 1, 2
- [34] R. Yeh, Z. Liu, D. B. Goldman, and A. Agarwala. Semantic facial expression editing using autoencoded flow. *arXiv preprint arXiv:1611.09961*, 2016. 2, 4, 6
- [35] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *ECCV*, 2014. 3
- [36] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *arXiv preprint arXiv:1703.10593*, 2017. 1, 2, 5, 7, 8