

# Heredity-aware Child Face Image Generation with Latent Space Disentanglement

Xiao Cui, Wengang Zhou, Yang Hu, Weilun Wang and Houqiang Li, *Fellow, IEEE*

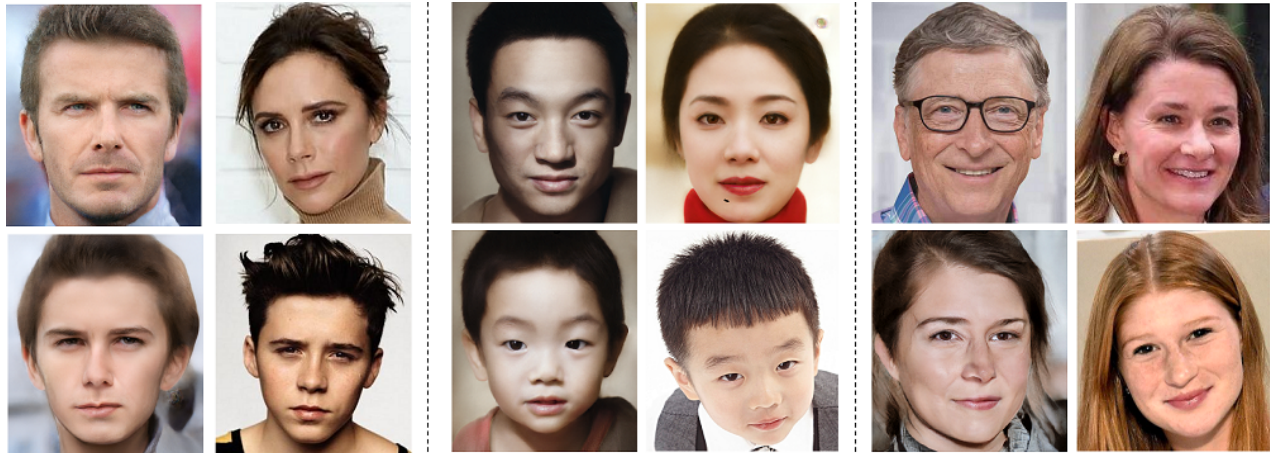


Fig. 1: Examples of children images generated by the proposed method. In each group, the images in the first row are the parents, and the second row shows the child generated (left) and the real image of their child (right).

**Abstract**—Generative adversarial networks have been widely used in image synthesis in recent years and the quality of the generated image has been greatly improved. However, the flexibility to control and decouple facial attributes (e.g., eyes, nose, mouth) is still limited. In this paper, we propose a novel approach, called ChildGAN, to generate a child’s image according to the images of parents with heredity prior. The main idea is to disentangle the latent space of a pre-trained generation model and precisely control the face attributes of child images with clear semantics. We use distances between face landmarks as pseudo labels to figure out the most influential semantic vectors of the corresponding face attributes by calculating the gradient of latent vectors to pseudo labels. Furthermore, we disentangle the semantic vectors by weighting irrelevant features and orthogonalizing them with Schmidt Orthogonalization. Finally, we fuse the latent vector of the parents by leveraging the disentangled semantic vectors under the guidance of biological genetic laws. Extensive experiments demonstrate that our approach outperforms the existing methods with encouraging results.

**Index Terms**—Child Face Image Generation; Generative Adversarial Networks; Semantics Learning; Latent Space Disentanglement

## I. INTRODUCTION

Child image generation aims at synthesizing child face image given the images of parents. This is a very challenging

Xiao Cui, Wengang Zhou, Yang Hu, Weilun Wang and Houqiang Li are with the Department of Electrical Engineering and Information Science, University of Science and Technology of China, Hefei, 230027 China (e-mail: cuixiao2001@mail.ustc.edu.cn, zhgw@ustc.edu.cn, eeyhu@ustc.edu.cn, wwlustc@mail.ustc.edu.cn, lihq@ustc.edu.cn).

Corresponding authors: Wengang Zhou and Houqiang Li.

task since the generated child face should not only resemble the parents, but inherit the attributes following the known genetic laws. Besides, the children born to the same parents may look quite different, which means there is no unique solution to the problem of child image generation. There are many interesting applications on child image generation, such as enabling a couple to preview the appearance of their children, kinship verification, *etc.*

There are only a few works studying this child generation problem. KinshipGAN [24] uses a deep face network to generate a child’s face based on one-to-one relationship. DNA-Net [10] propose to use a deep generative Conditional Adversarial Autoencoder for this task. Although some success has been achieved, those methods suffer three non-trivial issues. First, those methods cannot explicitly control the facial attributes in the generated faces, which significantly limits their application scenarios. Second, the generated images are usually blur and of low quality. Third, they ignore the inheritance law from genetic basis. For instance, thin upper lip is controlled by a dominant gene and is very likely to be inherited. Without considering such prior, the generated child images fail to reflect the inherited attributes.

In the past two years, significant progress have been made on semantic face editing. As an effective tool for editing, latent space is a representation of compressed data where each vector in it corresponds to one orientations. Researchers find that some orientations in the latent space of GANs, such as PGGAN [16] and StyleGAN [17], encode meaningful semantics of human faces. Moving the latent code vectors along these orientation will change the corresponding facial

attributes. These results are helpful to develop a child face generator so as to apply genetics knowledge when blending the parents’ faces and flexibly control over the attributes of the generated faces. Nevertheless, this idea is difficult to implement as we have to manipulate attributes that are more fine-grained than those edited by previous works.

In this paper, we propose a new framework, *i.e.*, ChildGAN, to generate the face image of the child according to the parents’ images under the guidance of genetic laws. The framework consists of two fusion steps: macro fusion and micro fusion. Before these main steps, we first project the parents’ images into the pre-trained latent space of StyleGAN. Then after some preprocessing for gender and age alignment as well as background decoupling, we mix the parents’ latent codes through a macro fusion module. In order to allow the child to inherit attributes of the parents in a micro way under the guidance of genetic laws, we propose a novel method to identify disentangled semantic directions in the latent space. Our semantic learning method is based on gradient estimation from a large number of samples as well as irrelevant factor reweighting. It finds the important and decoupled semantic vectors in the latent space without the need for manual labels. We also prove that it is feasible to orthogonalize the semantic vectors in the  $\mathbb{W}$  space of StyleGAN, which allows conditional manipulation of real images. It is notable that our method does not rely on any external datasets, but only use the images generated by a pre-trained StyleGAN generator and some real images of celebrities to conduct our experiment.

Our main contributions are summarized as follows,

- We propose a framework ChildGAN to generate the child’s face according to the parents’ faces. Our framework consists of two fusion steps, *i.e.*, macro fusion and micro fusion, which not only integrates the parents’ faces from the holistic perspective, but also involves processing at the micro level.
- To leverage genetic laws for the generation, a novel method is proposed to disentangle the latent space, which can extract meaningful and decoupled semantic vectors (*e.g.*, semantics corresponding to the size of eyes, nose, mouth, jaw, eyebrows and the thickness of lips) without external labels.
- We conduct extensive experiments to verify the effectiveness of our semantic learning module and child generation method.

## II. RELATED WORK

In this section Generative Adversarial Networks, semantic face editing and face generation are reviewed in detail, since these are three critical ingredients in our face generation approach.

### A. Generative Adversarial Networks

Generative Adversarial Networks (GANs) use the strategy of adversarial training to learn the probabilistic distribution function of training samples for novel sample generation. To date, GANs have been successfully applied to various computer vision tasks, such as semantic attribute editing

[20], [26], [21], image to image translation [3], [14], [38], [7] and video generation [33], [36]. In recent years, many different network structures for GANs have emerged, which significantly improve the stability of the network and the quality of the generated images. For example, PGGAN [16] increases the scales of the generator and the discriminator step by step, leading to more stable training and higher quality of generated images. StyleGAN [17] injects style vectors at each convolutional layer with adaptive instance normalization (AdaIN) to achieve subtle and precise style control, which improves the interpretability and controllability of the generation process. Our work further explores and decouples the latent space of StyleGAN and applies it to child image generation in combination with biological evidence.

### B. Semantic Face Editing

In semantic face editing, a general idea is to first project the image into some latent space and then edit the latent code. There are two major ways to embed examples from image space into latent space: learning an encoder that maps the given image to the latent space [19], [30], or starting with a random initial latent code and optimizing it by a gradient descent method [1], [2], [9], [15].

As for latent code modification, Upchurch Paul [34] proved that complex attribute transformation can be realized by linear interpolation in depth feature space. Due to the nice properties of the latent space of StyleGAN, some recent works [31], [1], [8], [5], [27] perform semantic editing based on the StyleGAN model. As a representative method, InterFaceGAN [31] trains SVMs to find the directions in the latent space corresponding to attributes including age, gender, pose, expression and the presence of eyeglasses. However, for children image synthesis, we need to edit attributes that are more challenging and fine-grained, such as the size of eyes, nose and mouth.

### C. Face Generation

In recent years, great progress has been made with GANs in the field of face generation [32], [37], [13], [6], [17], [16], [4], [18]. Different from these general face image generation work, we focus on a new interesting problem: kin face generation.

Generating face images of children from images of parents is a relatively new research problem. To our best knowledge, there are only two related works on this topic *i.e.*, KinshipGAN [24] and DNA-Net [10]. KinshipGAN [24] generates child’s face image based on one-to-one relationships (father-to-daughter, mother-to-son, *etc*). The generator follows an encoder-decoder structure, where the encoder extracts the main features of the father or mother, where the decoder uses these extracted features to generate the child’s faces. Differently, DNA-Net [10] uses a two-to-one relationship where both images of the father and the mother are used for the generation. However, the model of DNA-Net works mostly in a “black-box” fashion. It is difficult to impose knowledge of genetics to guide the generation and control or adjust the results on demand. Besides, the resolution and quality of the images generated by the above two works are worse than those generated by native GANs. In our work, we blend the characteristics

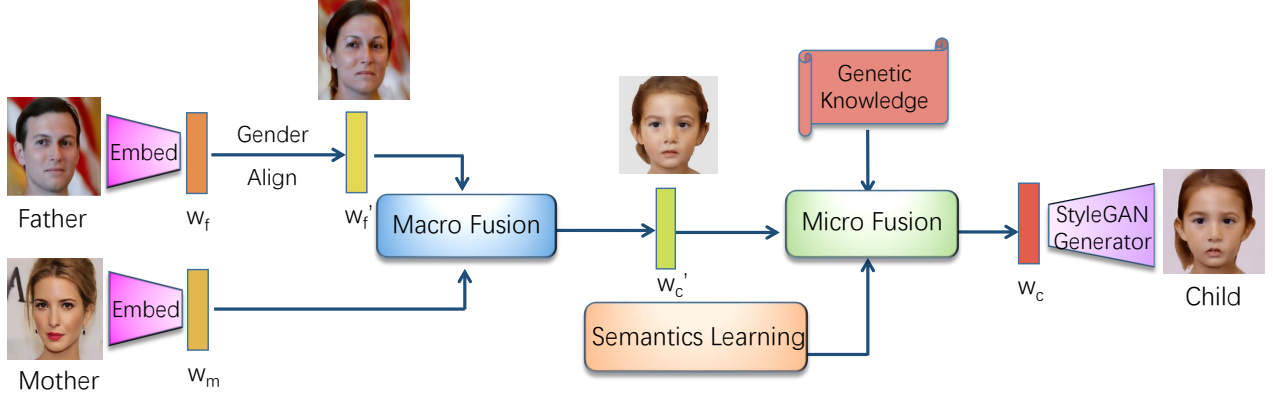


Fig. 2: The flowchart of our generation process. The images of the parents are first embedded to the latent space of StyleGAN. Then after some preprocessing, the latent codes of the parents are mixed through Macro Fusion. We propose an effective method to identify disentangled semantic directions in the latent space, which allow the child to inherit attributes of the parents in a micro way under the guidance of genetic laws. Finally, the child image is generated from the child’s latent code by a pre-trained StyleGAN generator.

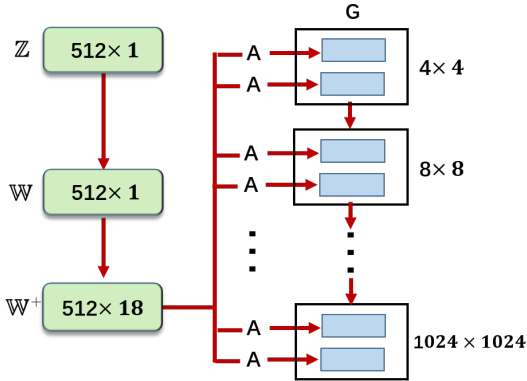


Fig. 3: Framework overview of StyleGAN[17] generator.  $Z$ ,  $W$  and  $W^+$  are three latent spaces,  $G$  is the synthesis network,  $A$  represents the affine transformation.

of the parents in the latent space of StyleGAN. To make the generation more scientifically sound and controllable, we propose a method to find disentangled semantics in this space, and employ rules of genetics to guide the generation.

### III. METHOD

The framework of the proposed ChildGAN is shown in Figure 2. The face images of the parents are first embedded to the latent space of StyleGAN. Then after some preprocessing, the latent codes of the parents are mixed through Macro Fusion. We propose an effective method to identify disentangled semantic directions in the latent space, which allow the child to inherit attributes of the parents in a micro way under the guidance of genetic laws. Finally, the child image is generated from the child’s latent code by a pre-trained StyleGAN generator.

In this section, we first conduct an overview StyleGAN and Image2StyleGAN, then we introduce the method of macro

fusion to achieve the inheritance of macro and relatively rough characteristics. Next we discuss the scientific knowledge for child generation. After that, we introduce the extraction and decoupling of important semantics. Finally, we present how to use the genetic knowledge and orthogonal semantic vectors for micro fusion.

#### A. An overview of StyleGAN and Image2StyleGAN

As a preparation for the following discussion, we make a brief introduction to StyleGAN [17] and Image2StyleGAN [1]. As shown in Figure 3, in StyleGAN, a non-linear mapping network  $f$  first maps a vector  $z \in \mathbb{R}^{512}$  to a vector  $w \in \mathbb{R}^{512}$  in the intermediate space  $W$ . After that, we make 18 copies of  $w$  vector and pass them through an affine transformation  $A$ , whose output is further fed into synthesis network  $G$  to control the style of its each layer. There are 18 convolution layers in  $G$ , two for each resolution, *i.e.*, 9 resolution layers in total. The larger the layer number, the higher the resolution. To use the pre-trained StyleGAN generator, the 18 copies of the  $w$  vector can be extended to 18 different  $w$  vectors whose concatenation constitutes the  $W^+$  space. The  $Z$  Space,  $W$  Space and  $W^+$  Space are called latent space (latent space specifically refers to  $W^+$  space in the following text). We choose to operate on the  $W^+$  space because it is suitable for processing real images and has better ability of disentanglement. In order to enable semantic image editing operations to be applied to existing photos, Image2StyleGAN [1] proposed an efficient algorithm, where a random initial vector is chosen and optimized using gradient descent, to embed a given image into the latent space of StyleGAN [17].

#### B. Macro Fusion

In order to generate a child’s image, we start by making a rough mix of the parents’ faces to get a preliminary image of the child. We call this process *macro fusion*. Before our fusion process, we map the parents’ images to vectors in the

latent space, which called latent codes. Given an image of the father, we first crop and align the face in it. Then similar to Image2StyleGAN [1], we find the optimal latent code  $\mathbf{w}_f$  by minimizing the reconstruction loss between the image generated from  $\mathbf{w}_f$  and the real image. The latent code  $\mathbf{w}_m$  for the mother is obtained in the same way. To produce a child with a specific gender and reduce the background artifacts, we change the gender character of one parent by moving the corresponding latent code along a pre-learned orientation in the latent space which mainly controls the gender attribute. That is to say, if we want the generated child to be a girl, we need to move the father’s latent code vector forward in this orientation (about 2 units in length), and if we want the child to be a boy, we need to move the mother’s latent code vector backward in this orientation.

There are two alternatives in macro fusion. First, as a simple method, we can use linear combination:  $\mathbf{w}'_c = (1 - \lambda) \mathbf{w}_f + \lambda \mathbf{w}_m$ , where  $\lambda$  is a parameter between 0 and 1. In this way, every resolution layer feature of the child will be a mixture of the parents. Alternatively, we can control different resolution layers with different  $\mathbf{w}$ . For example, if we want rough features such as posture, hairstyle, and facial contour to be inherited from the father, while more subtle features such as facial components are inherited from the mother on a macro level, we can set the first two resolution layers be controlled by  $\mathbf{w}_f$ , while the other layers are controlled by  $\mathbf{w}_m$ . At the end of the macro fusion, we just need to adjust  $\mathbf{w}'_c$  along the pre-learned age vector to generate the child of an expected age.

### C. Inheritance Prior for Child Generation

We adopt some genetic evidence in biology [22] to make our results more scientific (here we only consider Mendelian inheritance).

- 1) Skin color [11]: The skin color of a child always follows the natural law of “neutralizing” the skin colors of the parents.
- 2) Eyes: Big eyes are inherited in a dominant way, so as long as one parent has big eyes, the child is more likely to have big eyes.
- 3) Nose: Generally speaking, large, high noses and wide nostrils are in dominant inheritance. If one of the parents has a big nose, it is likely to be inherited by the child.
- 4) Jaw: As dominant inheritance, if one parent has a prominent big chin, the child will be more likely to grow into a similar chin.
- 5) Lip thickness: A thin upper lip is a dominant inheritance, while a thicker lower lip is a dominant inheritance.
- 6) Baldness [23]: Alopecia is caused by an autosomal dominant gene. Bald men may be heterozygous ( $Bb$ ) or homozygous ( $BB$ ), while bald women are homozygous ( $BB$ )

If we use  $A$  to present dominant gene and  $a$  for recessive gene, then the genotype with dominant trait is  $Aa$  or  $AA$  (in our work, we assume that they’re equally likely), and the genotype with recessive trait is  $aa$ . If one parent presents a recessive trait and the other presents a dominant trait, the probability of the child presenting a recessive trait is:

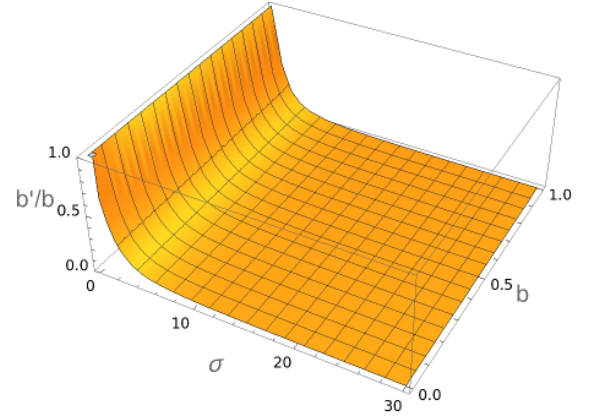


Fig. 4: A 3D plot of  $b'/b$  at different  $b$  and  $\sigma$ , where  $\sigma$  is the variance of  $\frac{u_{i,m}-u_{j,m}}{u_{i,l}-u_{j,l}}$ ,  $b$  is the mean of  $\frac{u_{i,m}-u_{j,m}}{u_{i,l}-u_{j,l}}$ , and  $b'$  is the mean of  $\mathbf{v}_l^e$  in the direction of  $\mathbf{v}_m$  extracted using improved method. In other words,  $b$  and  $b'$  are the mean values of the components of  $\mathbf{v}_l^e$  in the weighed direction when using the basic method and the improved method respectively.

$\frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$ , while the probability of presenting a dominant trait is  $\frac{3}{4}$ . If both parents display dominant traits, the probability of the child presenting a recessive trait is:  $\frac{1}{4} \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{16}$ , while the probability of presenting a dominant trait is  $\frac{15}{16}$ . If both parents present recessive traits, their child will surely present a recessive trait.

In order to classify facial attribute values corresponding to biological traits, we compare the value of the attribute with a threshold value. We get the value for each attribute by first detecting the face landmarks in the image and then computing the distances between corresponding landmarks. Each threshold is set as the average value of the attribute in a large number of samples. For attributes that cannot be classified by size, we use the projection value of the latent vector on the attribute vector’s orientation instead of distance difference as the value of the attribute.

### D. Semantics Learning

As we can see from Section III.B, a genetic rule usually describes how a certain facial attribute is passed down from the parents to the child. To generate the child face according to the genetic laws, we need to identify these attributes in the latent space  $\mathbb{W}+$  of StyleGAN, in which we fuse the faces of the parents. However, this is not readily available, since each  $\mathbf{w}$  vector usually relates to multiple attributes, *i.e.*, we cannot fuse an attribute by simply interpolating a chosen  $\mathbf{w}$ . Some previous works [31], [5], [27] have shown that there are directions in the latent space of StyleGAN that correspond to different attributes of a face. In this section, we propose an effective method to identify semantic directions (or semantics in short) in the  $\mathbb{W}+$  space that separately correspond to the attributes covered by the heredity laws. To ensure that moving a latent vector along one semantic direction affects other attributes as little as possible, we further make these semantic directions orthogonal to each other.

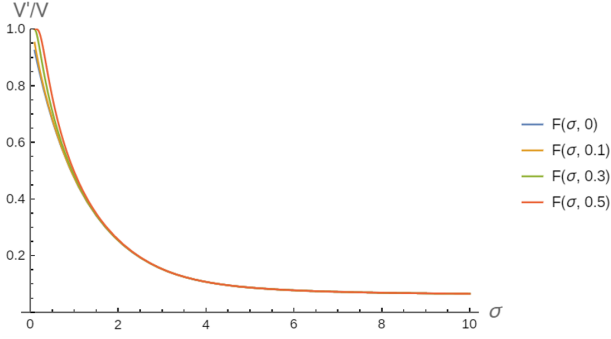


Fig. 5:  $V'/V - \sigma$  diagram under different  $b$ , where  $b$  is the mean of  $\frac{u_{i,m}-u_{j,m}}{u_{i,l}-u_{j,l}}$ ,  $V'$  is the variance of  $\mathbf{v}_l^e$  extracted by our improved method,  $V$  is the variance of  $\mathbf{v}_l^e$  extracted by our basic method, and  $\sigma$  is the variance of  $\frac{u_{i,m}-u_{j,m}}{u_{i,l}-u_{j,l}}$ .

It is widely observed that we can change the semantics contained in a synthesis continuously by linearly interpolating two latent codes. Let  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  be a set of vectors that we ultimately want, each representing a semantic direction in the latent space. The difference between the latent codes of two images can be represented as a linear combination of the semantic vectors  $\{\mathbf{v}_k\}$ . The weight for each  $\mathbf{v}_k$  is proportional to the change of the value for the corresponding attribute. Without loss of generality, we have:

$$\mathbf{w}_i - \mathbf{w}_j = \sum_k (u_{i,k} - u_{j,k}) \times \mathbf{v}_k, \quad (1)$$

where  $\mathbf{w}_i$  and  $\mathbf{w}_j$  are the latent codes of two images.  $u_{i,k}$  and  $u_{j,k}$  denote the values of the face attribute corresponding to  $\mathbf{v}_k$ . If  $u_{i,k} - u_{j,k} \neq 0$ , we have:

$$\frac{\mathbf{w}_i - \mathbf{w}_j}{u_{i,l} - u_{j,l}} = \mathbf{v}_l + \sum_{k \neq l} \frac{u_{i,k} - u_{j,k}}{u_{i,l} - u_{j,l}} \times \mathbf{v}_k. \quad (2)$$

With a large number of image pairs available, we can compute the expectation as:

$$E\left(\frac{\mathbf{w}_i - \mathbf{w}_j}{u_{i,l} - u_{j,l}}\right) = \mathbf{v}_l + \sum_{k \neq l} E\left(\frac{u_{i,k} - u_{j,k}}{u_{i,l} - u_{j,l}}\right) \times \mathbf{v}_k, \quad (3)$$

where  $E(p)$  represents the statistical expectation of  $p$ . As the distribution of  $u_{i,k} - u_{j,k}$  is symmetric with respect to 0,  $E(u_{i,k} - u_{j,k}) = 0$ , since  $u_{i,l} - u_{j,l} \neq 0$ , in the same way we have  $E\left(\frac{1}{u_{i,l} - u_{j,l}}\right) = 0$ . So if different semantics are independent of each other,  $E\left(\frac{u_{i,k} - u_{j,k}}{u_{i,l} - u_{j,l}}\right) = 0$ , then we have:

$$E\left(\frac{\mathbf{w}_i - \mathbf{w}_j}{u_{i,l} - u_{j,l}}\right) = \mathbf{v}_l. \quad (4)$$

Also we can compute the variance as:

$$\begin{aligned} V\left(\frac{u_{i,k} - u_{j,k}}{u_{i,l} - u_{j,l}}\right) &= V(u_{i,k} - u_{j,k}) V\left(\frac{1}{u_{i,l} - u_{j,l}}\right) \\ &+ E^2(u_{i,k} - u_{j,k}) V\left(\frac{1}{u_{i,l} - u_{j,l}}\right) \\ &+ V(u_{i,k} - u_{j,k}) E^2\left(\frac{1}{u_{i,l} - u_{j,l}}\right), \end{aligned} \quad (5)$$

where  $V(p)$  represents the statistical variance of  $p$ . According to the above, the last two terms can be dropped because they are 0. So with many sample pairs we have:

$$V\left(\frac{1}{n} \sum_{i,j=i_1,j_1}^{i_n,j_n} \frac{u_{i,k} - u_{j,k}}{u_{i,l} - u_{j,l}}\right) = \frac{V(u_{i,k} - u_{j,k}) V\left(\frac{1}{u_{i,l} - u_{j,l}}\right)}{n}. \quad (6)$$

As we choose to discard the sample pairs with  $u_{i,l} - u_{j,l} = 0$ , the minimum distance difference is 1,  $V(u_{i,k} - u_{j,k})$  and  $V\left(\frac{1}{u_{i,l} - u_{j,l}}\right)$  are finite values, the value of the right side of the above equation goes to zero as  $n$  gets very large, that makes:

$$\lim_{n \rightarrow \infty} V\left(\frac{1}{n} \sum_{i,j=i_1,j_1}^{i_n,j_n} \frac{u_{i,k} - u_{j,k}}{u_{i,l} - u_{j,l}}\right) = 0. \quad (7)$$

So when there are enough samples, we estimate  $\mathbf{v}_l$  through:

$$\mathbf{v}_l^e = \frac{2}{N \times (N-1)} \sum_{i=1}^{N-1} \sum_{j=i+1}^N \frac{\mathbf{w}_i - \mathbf{w}_j}{u_{i,l} - u_{j,l}}, \quad (8)$$

where  $N$  is the total number of images available for learning, and  $\mathbf{v}_l^e$  is an estimation of  $\mathbf{v}_l$ .

In fact, the distribution of values of different attributes in sample pictures is not independent of each other. For example, people with larger eyes have larger mouths on average. That makes  $E\left(\frac{u_{i,k} - u_{j,k}}{u_{i,l} - u_{j,l}}\right) > 0$ . According to Equation 3, the vector  $\mathbf{v}_l^e$  we find based on Equation 8 will have components not only in the  $\mathbf{v}_l$  direction, but also in other directions, which can be written as  $\mathbf{v}_l + \sum_{k \neq l} \frac{u_{i,k} - u_{j,k}}{u_{i,l} - u_{j,l}} \times \mathbf{v}_k$  as the sample size approaches infinity. This means that when we want to change the  $l$ -th attribute, the rest of the attributes will change as well. In order to reduce the proportion of irrelevant components  $\mathbf{v}_k$  ( $k \neq l$ ) in the extracted vector and reduce the influence of other attributes when changing the  $l$ -th feature, we add an additional weight to the terms in Equation 8:

$$\mathbf{v}_l^e = \frac{\sum_{i=1}^{N-1} \sum_{j=i+1}^N \frac{\mathbf{w}_i - \mathbf{w}_j}{u_{i,l} - u_{j,l}} \times e^{-\left|\frac{u_{i,m} - u_{j,m}}{u_{i,l} - u_{j,l}}\right|}}{\sum_{i=1}^{N-1} \sum_{j=i+1}^N e^{-\left|\frac{u_{i,m} - u_{j,m}}{u_{i,l} - u_{j,l}}\right|}}, \quad (9)$$

where  $m$  is the index of the semantic being conditioned due to the large value of  $E\left(\frac{u_{i,m} - u_{j,m}}{u_{i,l} - u_{j,l}}\right)$ . By introducing this weight factor,  $\mathbf{w}_i - \mathbf{w}_j$  would be weighted heavier if the difference of  $m$ -th attribute between the two images is small, and  $\mathbf{w}_i - \mathbf{w}_j$  would be given a lower weight if the  $m$ -th attribute of the two images are quite different. This also makes the variance of the restricted directional component smaller and accelerated the convergence speed. Equation 9 can be further extended if more than one attribute needs to be conditioned. We only need to multiply the weight factors corresponding to these attributes together.

We will prove below that the addition of the weight factor will reduce the expected value of  $\mathbf{v}_l^e$  in the  $\mathbf{v}_m$  direction, so that  $\mathbf{v}_l^e$  is closer to the real  $\mathbf{v}_l$ . Only considering the two-

dimensional space composed of  $\mathbf{v}_l$  and  $\mathbf{v}_m$ , Equation 9 can be rewritten as:

$$\mathbf{v}_l^{e'} = \frac{\sum_{i < j} \frac{u_{i,m} - u_{j,m}}{u_{i,l} - u_{j,l}} e^{-\left| \frac{u_{i,m} - u_{j,m}}{u_{i,l} - u_{j,l}} \right|} \mathbf{v}_m + e^{-\left| \frac{u_{i,m} - u_{j,m}}{u_{i,l} - u_{j,l}} \right|} \mathbf{v}_l}{\frac{N(N-1)}{2}}, \quad (10)$$

$$\mathbf{v}_l^e = \frac{\frac{N(N-1)}{2}}{\sum_{i=1}^{N-1} \sum_{j=i+1}^N e^{-\left| \frac{u_{i,m} - u_{j,m}}{u_{i,l} - u_{j,l}} \right|}} \mathbf{v}_l^{e'}, \quad (11)$$

so we can calculate the expectation as:

$$E(\mathbf{v}_l^{e'}) = E\left(\frac{u_{i,m} - u_{j,m}}{u_{i,l} - u_{j,l}} \times e^{-\left| \frac{u_{i,m} - u_{j,m}}{u_{i,l} - u_{j,l}} \right|}\right) \mathbf{v}_m + E\left(e^{-\left| \frac{u_{i,m} - u_{j,m}}{u_{i,l} - u_{j,l}} \right|}\right) \mathbf{v}_l, \quad (12)$$

$$E(\mathbf{v}_l^e) = \frac{E\left(\frac{u_{i,m} - u_{j,m}}{u_{i,l} - u_{j,l}} \times e^{-\left| \frac{u_{i,m} - u_{j,m}}{u_{i,l} - u_{j,l}} \right|}\right)}{E\left(e^{-\left| \frac{u_{i,m} - u_{j,m}}{u_{i,l} - u_{j,l}} \right|}\right)} \mathbf{v}_m + \mathbf{v}_l. \quad (13)$$

Let's assume that  $\frac{u_{i,m} - u_{j,m}}{u_{i,l} - u_{j,l}}$  satisfies a normal distribution with mean  $0 < b < 1$  and variance  $\sigma^2 > 0$  (it should actually be written as a superposition of a series of different normal distributions, but it makes no difference to the proof). With the weight factor,

$$E(\mathbf{v}_l^e) = \frac{\sqrt{\frac{1}{2\pi}} \frac{1}{\sigma} \int_{-\infty}^{\infty} e^{-\frac{(x-b)^2}{2\sigma^2}} - |x| dx}{\sqrt{\frac{1}{2\pi}} \frac{1}{\sigma} \int_{-\infty}^{\infty} e^{-\frac{(x-b)^2}{2\sigma^2}} - |x| dx} \mathbf{v}_m + \mathbf{v}_l = b' \mathbf{v}_m + \mathbf{v}_l, \quad (14)$$

while with the basic method,

$$E(\mathbf{v}_l^e) = \frac{\sqrt{\frac{1}{2\pi}} \frac{1}{\sigma} \int_{-\infty}^{\infty} e^{-\frac{(x-b)^2}{2\sigma^2}} x dx}{\sqrt{\frac{1}{2\pi}} \frac{1}{\sigma} \int_{-\infty}^{\infty} e^{-\frac{(x-b)^2}{2\sigma^2}} dx} \mathbf{v}_m + \mathbf{v}_l = b \mathbf{v}_m + \mathbf{v}_l \quad (15)$$

Figure 4 shows the relationship between  $b'/b$  and  $b$  and  $\sigma$ . We can see that  $b'/b$  is less than 1 and approaches 0 when  $\sigma$  is large. This indicates that after adding the weight factor, when we try to change the  $l$ -th attribute, the influence on the  $m$ -th attribute will be smaller, these two attribute can be better decoupled.

In Figure 5, we calculate the variance ratio of the  $\mathbf{v}_l^e$  extracted by the improved method and the basic method in the  $m$ -th direction. It can be seen that the addition of the weight factor can significantly reduce the variance, so we only need a relatively small amount of data to make the results converge.

In addition to focusing on the overall characteristics of latent space, differences in different resolution layers in latent space facilitate further decoupling. We found that the first resolution layer of StyleGAN mainly controls camera elevation and horizontal angles, while the last four resolution layers mainly control color and background. In order to decouple the extracted facial semantics from attributes that we are not interested in, we can choose to only work on the middle three resolution layers.

The selection of the resolution layers and the estimation of the semantic vectors  $\{\mathbf{v}_l\}$  have to a large extent disentangled



Fig. 6: The results of changing the latent codes that control different resolution layers of StyleGAN [17]. For each person, the first picture is the original image, and remaining ones are the results of replacing the latent codes which control the  $4 \times 4, 8 \times 8, \dots, 1024 \times 1024$  resolution layers with a zero vector, respectively.

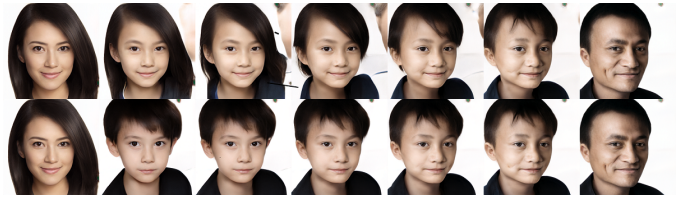


Fig. 7: Examples of gender align. The first and second rows are the results with and without gender align, respectively.

the semantics. However, there are still some coupling of attributes in  $\{\mathbf{v}_l\}$ . To achieve more precise control, we further orthogonalize these vectors. In our work, we use the Gram-Schmidt process to make the semantic vectors orthogonal to each other. Start with  $\mathbf{n}_1 = \mathbf{v}_1$ , we have:

$$\mathbf{n}_l = \mathbf{v}_l - \sum_{i=1}^{l-1} \frac{\langle \mathbf{v}_l, \mathbf{n}_i \rangle}{\langle \mathbf{n}_i, \mathbf{n}_i \rangle} \mathbf{n}_i, \quad (16)$$

where  $\mathbf{v}_l$  is the semantic vector found by Equation 9 and  $\mathbf{n}_i$  represents the orthogonal vector.

In InterFaceGAN [31], the authors failed to decouple some semantics through orthogonalization in the  $\mathbb{W}$  space of StyleGAN. For example, for the ‘‘age’’ and ‘‘eyeglasses’’ attributes, they found an eyeglasses-included age direction that is some-

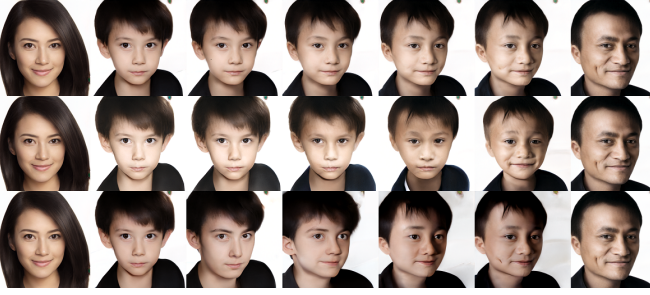


Fig. 8: Examples of macro fusion. The first row is the result of weighted mixing. The second and the third rows are the results of feeding the parents’ latent codes to different resolution layers, where in the second row the mother’s latent code controls layers with low resolutions and the father’s code controls layers with high resolutions and the opposite in the third row.

how orthogonal to the eyeglasses direction itself. So it can not remove eyeglasses from the age direction by orthogonalization. In our experiment, this situation does not happen. In other words, our method can identify semantic directions that are more disentangled than those found by InterFaceGAN. We give an example in Figure 12.

With the disentangled semantics identified by the method shown in this section, we can decompose the latent vectors of the parents by projecting them onto these semantic directions. Then an attribute of the child will be determined by picking a point in each semantic direction based on the genetic laws. We call this process the micro fusion of the parents.

#### E. Micro Fusion

Now that we have obtained the decoupled semantic vectors that correspond to key attributes of the face, we can inherit face components of the parents according to the genetic laws. Based on the preliminary child latent code obtained after macro fusion, we further adjust it in the semantic directions. For each semantic vector, we first project the parents’ and preliminary child’s latent codes onto it. For the case that one parent presents a dominant trait while the other parent presents a recessive trait, we get the child’s phenotype according to probability, and move the child’s latent code to the father or mother’s projection. If both parents show dominant traits but the child should show the recessive character according to probability, we move the child’s latent code across the less obvious dominant side and move on until it becomes recessive. In the case of parents and child all showing dominant traits, we make the child’s latent code move randomly under the restriction of parents’ projection in this direction (the same for both parents with recessive traits or when there is no clear genetic rule to guide the semantic).

After dealing with each semantic direction in accordance with the above methods, we resynthesize  $w_c$  by:

$$w_c = \hat{w}'_c + \sum_l p_l v_l, \quad (17)$$

where  $\{v_l\}$  are the semantic vectors,  $p_l$  is the projection component on each semantic direction and  $\hat{w}'_c$  is what’s left of  $w'_c$  after been decomposed. After that, we send  $w_c$  into the StyleGAN Generator to obtain the final child image.

Attribute	Index of landmarks
length of eyebrow	(18, 22), (23, 27)
width of eyes	(38, 42), (45, 47)
length of eyes	(37, 40), (43, 46)
width of nose	(34, 36), (32, 34)
thickness of upper lip	(52, 63), (51, 62), (53, 64)
thickness of lower lip	(58, 67), (59, 68), (57, 66)
width of mouth	(52, 58), (51, 59), (53, 57)
length of mouth	(49, 55)
chin shapeness	(8, 10), (7, 11)

TABLE I: The face attributes we consider and the corresponding landmark pairs we use to get the labels for the attributes.

## IV. EXPERIMENTS

We use the extended  $\mathbb{W}$  space in StyleGAN, namely  $\mathbb{W}^+$  space, which is a concatenation of 18 different 512-dimensional  $w$  vectors, to carry out our experiment. Previous works [1] have shown that the  $\mathbb{W}^+$  space is more effective than the original  $\mathbb{W}$  space. In this section, we will first reveal the influence of gender alignment and two kinds of macro fusion, and then show the effects of semantics learning for micro fusion, including the capability of different extraction methods for semantic vectors and the effectiveness of orthogonalization in  $\mathbb{W}^+$  space. After that, we will present the results of child generation and show how the genetic laws work. Finally we will present some quantitative and visual results.

#### A. Evaluation of Macro Fusion

We explore the effects of the latent codes for different resolution layers by replacing each latent code with a zero vector, respectively. As shown in Figure 6, the code for first resolution layer mainly affects the face orientation and elevation angle, the latent codes for the last four resolution layers mainly affect the color and illumination. That enables us to decouple these irrelevant attributes by keeping the resolution layers mentioned above unchanged.

Figure 7 shows that there is coupling between gender and background in StyleGAN’s  $\mathbb{W}$  space. Under different gender attributes, the same background corresponds to different latent codes, resulting in background artifacts after macro fusion. In order to remove the background confusion, we adjust the gender of one of the parents before the macro fusion.

Figure 8 shows the results of macro fusion using the two different ways respectively. In the case of a weighed combination, with the increase in the proportion of father’s latent code, different attributes of the child are equally shifted from mother-like to father-like. As for feeding the parents’ latent codes to different resolution layers, it can be seen that in the second row, from left to right, the child’s subtle features first change from maternal to paternal, and it is only after the subtle changes are complete that the rough features begin to shift to paternal. As opposed to this, in the third row the

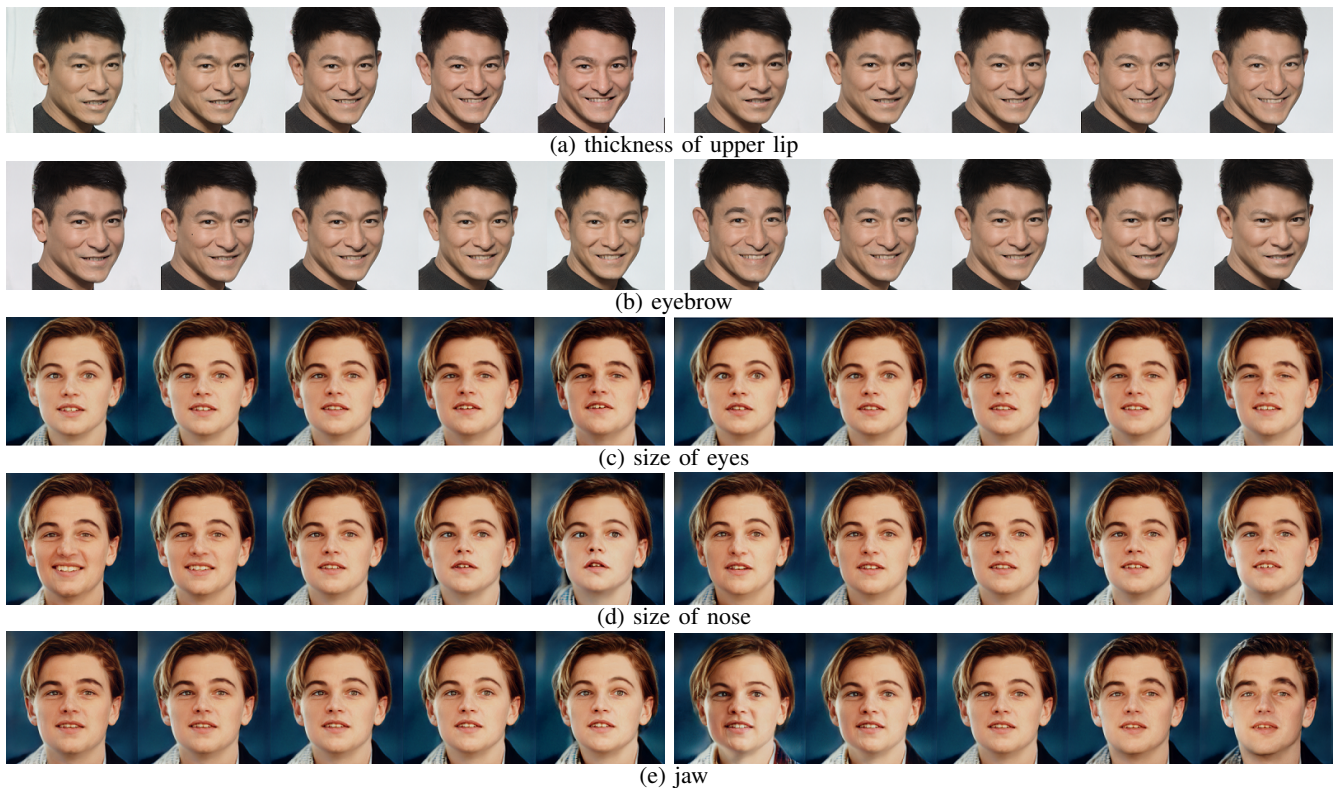


Fig. 9: Single attribute manipulation results. In each row, the five images on the left are the result of using semantic vectors extracted in the basic way, and the five images on the right are the results using semantic vectors extracted in the improved way. For each group of five images, the one in the middle is the original image, and the images on the left and right sides of it are obtained by moving its latent code along the positive and negative directions of the corresponding semantic vector respectively.

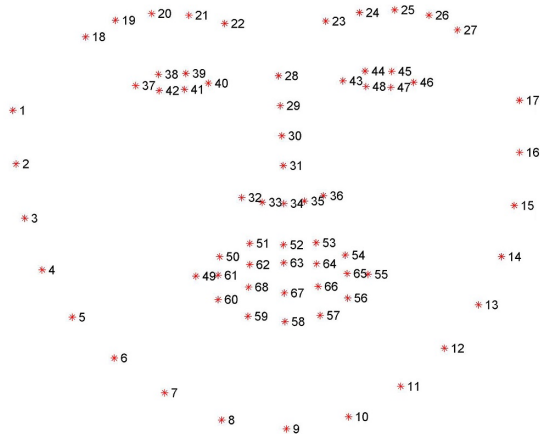


Fig. 10: Positions of the 68 landmarks on a human face. <sup>1</sup>

child’s rough features first change from mother-like to father-like while the subtle features change later.

### B. Evaluation of Micro Fusion

We use 10,000 images generated randomly by the pre-trained StyleGAN to learn the semantic vectors. We compare

<sup>1</sup> From [http://dlib.net/files/shape\\_predictor\\_68\\_face\\_landmarks.dat.bz2](http://dlib.net/files/shape_predictor_68_face_landmarks.dat.bz2)

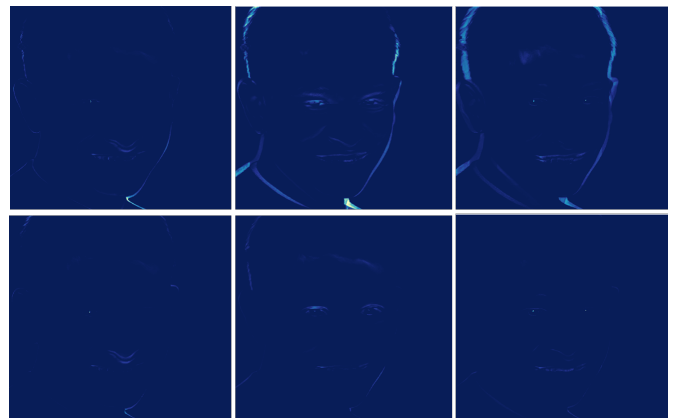


Fig. 11: Heat maps of the mean squared error between the edited outputs and original image. The first and second rows are the results of using the basic method and the improved method, respectively. The edited attributes are the nose size (left), eye size (middle) and upper lip thickness (right).

two versions of our learning method. In the basic method, we use Equation 8 for semantics extraction. In the improved version of our method, we use Equation 9 and keep the first two and the last six  $w$  vectors in the  $\mathbb{W}^+$  space unchanged. As we use the coordinate difference between different landmarks





Fig. 12: Analysis on whether it is effective to do orthogonalization in  $\mathbb{W}$  Space of StyleGAN [17]. The first row and the second row are the results of aging without and with orthogonalization, respectively.

to represent the size or thickness of different facial semantics, we do not need to use external labels or manual markings.

For a given face image, we get the labels for the attributes, such as the size of the eyes, the size of the nose, and the thickness of the lip, by first detecting the face landmarks in the image and then computing the distances between corresponding landmarks pairs. We use a detector <sup>1</sup> which can locate 68 keypoints of face components. The distances between pairs of landmarks are used as quantitative labels for the attributes.

In Figure 10 we show the positions of the 68 landmarks on human face. The landmark pairs selected to measure the sizes of the facial components are displayed in Table I. By calculating the distance between each pair of landmarks, we describe the size or thickness of the main components of the face without using external labels or manual marking. For some attributes, more than one set of landmark points can be used to describe them. We take the average of distances to get the corresponding labels. For example, we use the average of the distances between (38, 42) and between (45, 47) as a label for the width of eyes.

Figure 9 shows that with the semantic vectors extracted by the basic method, we are able to control the important facial attributes, but there are still some coupling problems with other irrelevant attributes such as the change of the pose. With the improved method, these problems have been greatly alleviated. As shown from the heatmaps in Figure 11, our improved method focuses on the component of interest better than our basic method, and it no longer modifies the face contour. These all demonstrate that we can manipulate the facial attributes better with less attribute coupling.

To demonstrate the effect of the orthogonalization operation, like InterFaceGAN [31], we first tested the coupling of age and eyeglass attributes in Figure 12. We observe that the appearance of glasses will be accompanied by the increase of age. However, after the orthogonalization process, the age attribute will no longer be coupled with the eyeglasses semantic. Also, the cosine similarity between the original age vector and the eyeglass vector is 0.13, which is not a very close number to 0. All these show that the orthogonalization of age vector is effective. Other attributes also show this effectiveness. In addition, we find that the order of the semantic vectors in the

Item	Accuracy(Epoch=10)	Accuracy(Epoch=20)
FIW test data	0.696	0.729
Gene(ours)	0.739	0.870
Real	0.714	0.722
Gene(DNA-Net)	0.563	0.563

TABLE II: Kinship verification scores. The first row is the accuracy of the network’s identification of kinship in the FIW [28] DataSet, and the following each row is the proportion of the corresponding child-parent pairs judged by the network to have kinship among all the pairs

Type	Average distance
Real - real pairs	1.194
Gener - real(randomly selected) pairs	1.189
Gener - real(corresponding) pairs	0.943

TABLE III: Similarity scores computed between every two features extracted by Facenet using cosine distance

process of orthogonalization has little effect on the decoupling results, which demonstrates the robustness of the operation.

### C. Results of Child Generation

Figure 1 and Figure 13 show some of our generated results. As can be seen, all facial attributes of children are combinations of the corresponding parents. The children’s attributes may be biased to one of the parents, which is mainly regulated by the genetic rules. Figure 14 compares our method with DNA-Net [10]. Although the input images given by their paper have low resolution, our results still show great advantages. In Figure 15, we show how the known genetic laws act on specific characteristics. The results demonstrate that the dominant traits of parents are more likely to be passed on to their children, which makes our generation process more scientifically sound and reliable.

In Figure 16 we show the diversity of our generated results. We can generate children with different ages and genders. Also, since the heredity of various characteristics follows the Mendelian inheritance law, it is not deterministic but with a certain probability that different children of the same parents do not look exactly the same in appearance. Our results show this diversity.

### D. Quantitative Evaluation

We quantitatively evaluate the proposed method through kinship verification and similarity calculation, and visualize the feature distribution.

For kinship verification, we used two Resnet50 [12] pre-trained on the VggFace dataset [25] to extract image features respectively and fine-tune the entire network on the FIW [28] dataset. The structure of network is shown in figure 17. For the network, given two pictures of faces, it can output whether the two people in the pictures are related. After 10 and 20 rounds of training, 69.6% and 72.9% accuracy were achieved on the FIW [28] data set respectively, and similar performance was also achieved in other real image samples we selected, indicating that no over-fitting occurred.

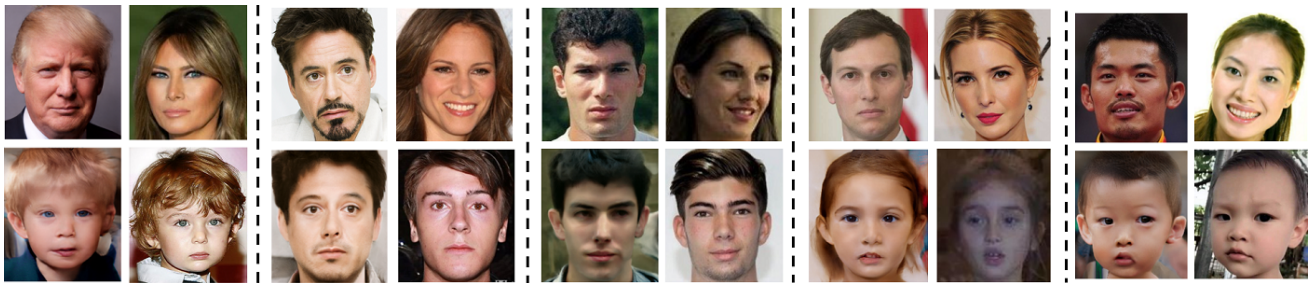
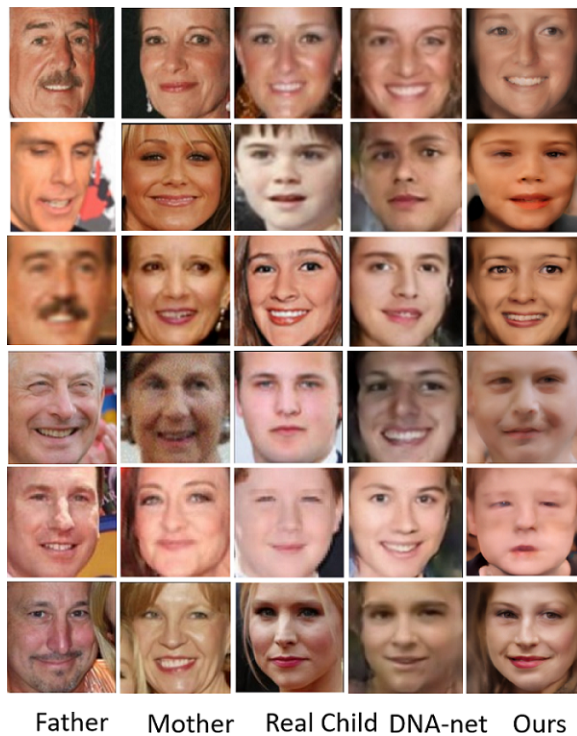


Fig. 13: More results of child generation. In each group, the images in the first row are the parents, while the second row shows the child generated by the proposed method (left) and the real image of the child (right)



Father Mother Real Child DNA-net Ours

Fig. 14: Comparison between our results and the results obtained by DNA-Net [10].

We submitted the generated child images and the images of their corresponding parents to the network for judgment. The higher the proportion of the groups considered by the classifier to have kinship, the better the performance of the generation method. As shown in Table II, in our generation method, the kinship verification network trained for 10 and 20 epochs considers 73.9% and 87.0% of image pairs to have kinship relations respectively, which is close to or even better than the results of real image pairs. And images generated by DNA-Net [10] achieved 56.3% and 56.3%, respectively.

We use a pre-trained Facenet [29] to perform the similarity test. Facenet [29] extracts the image identity features and directly learns the mapping between the image and the points on the Euclidean space of the feature. The distance of the points corresponding to the two images directly corresponds to the similarity of the two images. The training data are totally



Fig. 15: Examples of the role of genetic laws. Each column is the result of considering one genetic factor, and all children inherit the corresponding dominant traits.

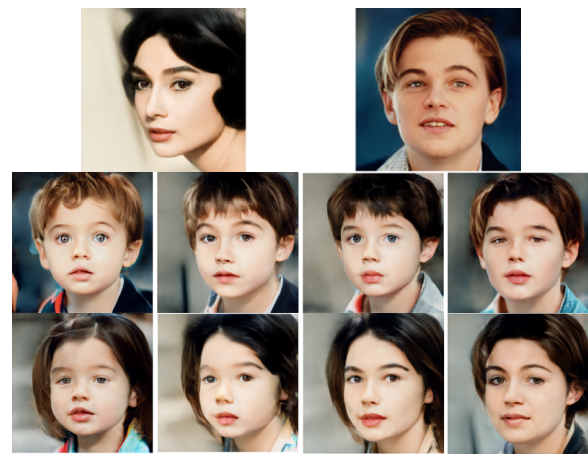


Fig. 16: An example about the diversity of generated results. The first row shows the images of the parents, the second and the third rows are the results of children with different genders, ages, and genetic patterns.

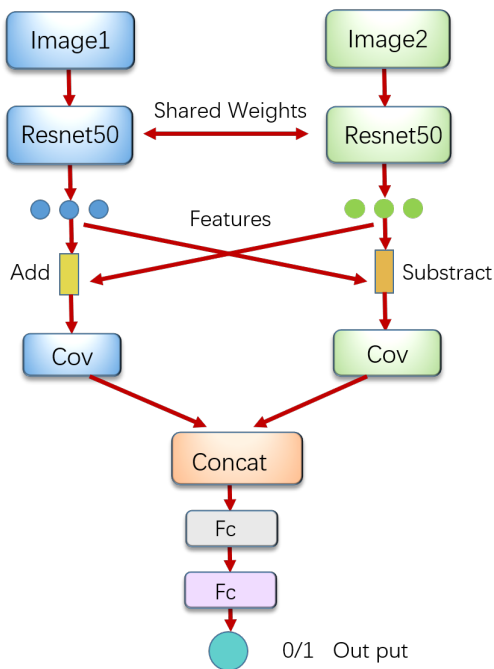


Fig. 17: Kinship verification network architecture.

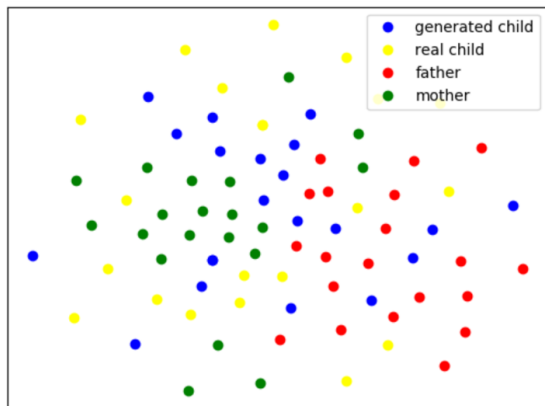


Fig. 18: Visualization of facial feature distribution of fathers, mothers, children, and generated ones. Red points represent the feature of fathers, green for mothers, yellow for real children, blue for generated children, respectively. Best viewed in color.

independent from FIW. As shown in Table III, the average distance of real-to-generated pairs is 0.943, compared to 1.189 for generated faces and random real faces and 1.194 for real-real pairs. This indicates that the child images generated by us have good feature similarity with the corresponding images of parents.

T-Distributed Stochastic Neighbor Embedding [35] is a machine learning algorithm for dimensional reduction. We use it to reduce the dimensionality of high-dimensional facial features to 2-dimensional features for visualization. As shown in Figure 18, unlike in DNA-Net [10] where the features the generated children are concentrated on one side and far away from real ones, the features of children generated by

us are evenly distributed. The distribution of image features generated by us is similar to that of real children’s images, and the generated image features are also very close to parents’ image features (some are closer to father’s features and some are closer to mother’s features), which is consistent with our verification results.

## V. CONCLUSION

We proposed ChildGAN to generate child images from the images of a couple under the guidance of genetic laws. Based on two fusion steps, our approach not only integrates the parents’ faces from the macro perspective, but also processes at the micro level. With a new method for semantic learning, we can precisely and independently control key facial attributes including eyes, nose, jaw, mouth and eyebrows. Extensive experiments suggest that our semantic learning module and child generation method are effective. For further work, more available genetic basis can make our work closer to real genetics. Also, exploring the latent code relationships of the parent-child datasets might lead to some new genetic evidence.

## REFERENCES

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2StyleGAN: How to embed images into the StyleGAN latent space? In *ICCV*, pages 4432–4441, 2019. 2, 3, 4, 7
- [2] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2StyleGAN++: How to edit the embedded images? In *CVPR*, pages 8296–8305, 2020. 2
- [3] Yazeed Alharbi, Neil Smith, and Peter Wonka. Latent filter scaling for multimodal unsupervised image-to-image translation. In *CVPR*, pages 1458–1466, 2019. 2
- [4] Jeff Donahue, Andrew Brock, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *ICLR*, pages 1–35, 2019. 2
- [5] Mohamed Elgharib Ayush Tewari, Florian Bernard Gaurav Bharaj, Hans-Peter Seidel, Patrick Pérez, Michael Zollhöfer, and Christian Theobalt. StyleRig: Rigging StyleGAN for 3D control over portrait images. In *CVPR*, pages 6141–6150, 2020. 2, 4
- [6] WenTing Chen, Xinpeng Xie, Xi Jia, and Linlin Shen. Texture deformation based generative adversarial networks for face editing. *arXiv preprint arXiv:1812.09832*, 2018. 2
- [7] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. StarGAN: Unified generative adversarial networks for multi-domain image-to-image translation. In *CVPR*, pages 8789–8797, 2018. 2
- [8] Edo Collins, Raja Bala, Bob Price, and Sabine Susstrunk. Editing in style: Uncovering the local semantics of GANs. In *CVPR*, pages 5771–5780, 2020. 2
- [9] Antonia Creswell and Anil Anthony Bharath. Inverting the generator of a generative adversarial network. *IEEE Transactions on Neural Networks and Learning Systems*, 30(7):1967–1974, 2018. 2
- [10] Pengyu Gao, Siyu Xia, Joseph Robinson, Junkang Zhang, Chao Xia, Ming Shao, and Yun Fu. What will your child look like? DNA-net: Age and gender aware kin face synthesizer. *arXiv preprint arXiv:1911.07014*, 2019. 1, 2, 9, 10, 11
- [11] G Ainsworth Harrison. The measurement and inheritance of skin colour in man. *The Eugenics review*, 49(2):73, 1957. 4
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 9
- [13] Fisher Yu Huiwen Chang, Jingwan Lu and Adam Finkelstein. PairedCycleGAN: Asymmetric style transfer for applying and removing makeup. In *CVPR*, pages 40–48, 2018. 2
- [14] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, pages 1125–1134, 2017. 2
- [15] Eli Shechtman Jun-Yan Zhu, Philipp Krhenbhl and Alexei A. Efros. Generative visual manipulation on the natural image manifold. In *Lecture Notes in Computer Science*, pages 597–613, 2016. 2

- [16] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. In *ICLR*, pages 1–26, 2018. 1, 2
- [17] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, pages 4401–4410, 2019. 1, 2, 3, 6, 9
- [18] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *CVPR*, pages 8110–8119, 2020. 2
- [19] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, pages 1–14, 2013. 2
- [20] Tingting Li, Ruihe Qian, Chao Dong, Si Liu, Qiong Yan, Wenwu Zhu, and Liang Lin. BeautyGAN: Instance-level facial makeup transfer with deep generative adversarial network. In *ACM MM*, pages 645–653, 2018. 2
- [21] Ming Liu, Yukang Ding, Min Xia, Xiao Liu, Errui Ding, Wangmeng Zuo, and Shilei Wen. STGAN: A unified selective transfer network for arbitrary image attribute editing. In *CVPR*, pages 3673–3682, 2019. 2
- [22] Victor A McKusick. *Mendelian inheritance in man: catalogs of autosomal dominant, autosomal recessive, and X-linked phenotypes*. Elsevier, 2014. 4
- [23] Dorothy Osborn. Inheritance of baldness: Various patterns due to heredity and sometimes present at birth—a sex-limited character—dominant in man—women not bald unless they inherit tendency from both parents. *Journal of Heredity*, 7(8):347–355, 1916. 4
- [24] Savas Ozkan and Akin Ozkan. KinshipGAN: Synthesizing of kinship faces from family photos by regularizing a deep face network. In *ICIP*, pages 2142–2146, 2018. 1, 2
- [25] Omkar M Parkhi, Andrea Vedaldi, and Andrew Zisserman. Deep face recognition. pages 1–12, 2015. 9
- [26] Shengju Qian, Kwan-Yee Lin, Wayne Wu, Yangxiaokang Liu, Quan Wang, Fumin Shen, Chen Qian, and Ran He. Make a face: Towards arbitrary high fidelity face manipulation. In *ICCV*, pages 10033–10042, 2019. 2
- [27] Niloy Mitra Rameen Abdal, Peihao Zhu and Peter Wonka. StyleFlow: Attribute-conditioned exploration of StyleGAN-generated images using conditional continuous normalizing flows. *arXiv preprint arXiv:2008.02401*, 2020. 2, 4
- [28] Joseph P Robinson, Ming Shao, Yue Wu, Hongfu Liu, Timothy Gillis, and Yun Fu. Visual kinship recognition of families in the wild. *IEEE TPAMI*, 40(11):2624–2637, 2018. 9
- [29] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, pages 815–823, 2015. 10
- [30] Bingbing Ni Shanyan Guan, Ying Tai, Feiyue Huang Feida Zhu, and Xiaokang Yang. Collaborative learning for faster StyleGAN embedding. *arXiv preprint arXiv:2007.01758*, 2020. 2
- [31] Yujun Shen, Jinjin Gu, Xiaou Tang, and Bolei Zhou. Interpreting the latent space of GANs for semantic face editing. In *CVPR*, pages 9243–9252, 2020. 2, 4, 6, 9
- [32] Hao Yang Shuyang Gu, Jianmin Bao, Dong Chen, Fang Wen, and Lu Yuan. Mask-guided portrait editing with conditional GANs. In *CVPR*, pages 3436–3445, 2019. 2
- [33] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. MocoGAN: Decomposing motion and content for video generation. In *CVPR*, pages 1526–1535, 2018. 2
- [34] Paul Upchurch, Jacob Gardner, Geoff Pleiss, Robert Pless, Noah Snaveley, Kavita Bala, and Kilian Weinberger. Deep feature interpolation for image content changes. In *CVPR*, pages 7064–7073, 2017. 2
- [35] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *JMACH*, 9(11), 2008. 11
- [36] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *NeurIPS*, pages 613–621, 2016. 2
- [37] Qi Li Yunfan Liu and Zhenan Sun. Attribute-aware face aging with wavelet-based generative adversarial networks. In *CVPR*, pages 11877–11886, 2019. 2
- [38] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. In *NeurIPS*, pages 465–476, 2017. 2