# What Will Your Child Look Like? DNA-Net: Age and Gender Aware Kin Face Synthesizer

Pengyu Gao
Southeast University
pi_1412@163.com

Siyu Xia
Southeast University
xia081@gmail.com

Joseph Robinson
Northeastern University
robinson.jo@husky.neu.edu

Junkang Zhang
University of California San Diego
juz007@eng.ucsd.edu

Chao Xia
ShangHai Jiao Tong University
xiabc612@gmail.com

Ming Shao
University of Massachusetts Dartmouth
mshao@umassd.edu

YUN FU
Northeastern University
yunfu@ece.neu.edu

## Abstract

*Visual kinship recognition aims to identify blood relatives from facial images. Its practical application– like in law-enforcement, video surveillance, automatic family album management, and more– has motivated many researchers to put forth effort on the topic as of recent. In this paper, we focus on a new view of visual kinship technology: kin-based face generation. Specifically, we propose a two-stage kin-face generation model to predict the appearance of a child given a pair of parents. The first stage includes a deep generative adversarial autoencoder conditioned on ages and genders to map between facial appearance and high-level features. The second stage is our proposed DNA-Net, which serves as a transformation between the deep and genetic features based on a random selection process to fuse genes of a parent pair to form the genes of a child. We demonstrate the effectiveness of the proposed method quantitatively and qualitatively: quantitatively, pre-trained models and human subjects perform kinship verification on the generated images of children; qualitatively, we show photo-realistic face images of children that closely resemble the given pair of parents. In the end, experiments validate that the proposed model synthesizes convincing kin-faces using both subjective and objective standards.*

## 1. Introduction

The goal of automatic kinship recognition is to determine whether or not people are related, and furthermore if so, the type of relationship shared. In the visual do-
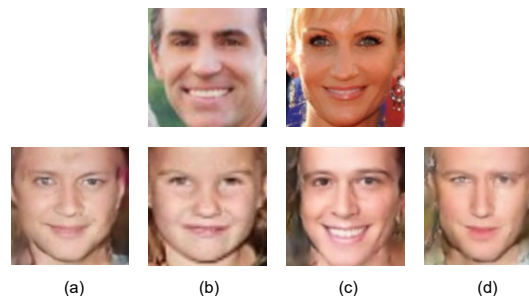


Figure 1: From faces of parents (top row), which face resembles their child the most (bottom row)? Three of the faces are generated, while one is real. Can you guess which one?

main, faces are typically used as the cue to determine kinship. This technology can be applied to mine social relationship [1], build a family tree [2], aid criminal investigations, do nature-based studies [3], and more. From this, kinship recognition has gained the interest of vast researchers nowadays.

In this work, we tackle a different task than is traditionally addressed in of kinship recognition, *i.e.* kin-face generation. Our aim to predict the appearance of a child from a pair of parents conditioned on high-level features (*i.e.* age and gender), which provides control over the desired characteristics.

The biological mechanisms that drive the visual resemblance of parents and their children inspired our efforts, and thus ability, to automatically understand kin-faces [4, 5].

Daly and Wilson [6] hypothesized that face similarity is sufficient evidence for kinship. Naini and Moss [7] claimed to have cracked the code for finding the most critical genetic features, which they quantified as "relatedness". More recently, researchers generated heritability maps that link facial landmarks to specific phenotypes of twins [8]]. The generated maps were from high-resolution faces (*i.e.* 4,096 landmarks) of 954 twins captured by expensive 3D cameras, which the authors identified genetic correspondents in the face variations of twins.

Typically, two directions are followed to recognize kin-faces: hand-crafted features and metric-based learning. Nowadays, deep models, especially Convolutional Neural Network (CNN), have shown promising discriminative power when used to encode faces for kinship recognition, pushing the state-of-the-art in the verification (*i.e.* one-to-one) task [9, 10].

Out of the many recent works in automatic kinship, only a few have attempted the kinship generation problem. Ertugrul *et al*. [11] focused on generating the facial dynamics of a child (*e.g.* smile) from a video of a parent showing different facial expressions. Ozkan *et al*. [12] generated a child's face, given a parent via adversarial training with constraints on the gender class and cycle consistency. Note, existing approaches that generate kin-faces, although unique in their ways, share a common flaw– only a single parent used to predict faces of children. These methods are unable to incorporate information from a pair of parents– the results are ineffective when compared to true child. Furthermore, they do not properly mimic nature (*i.e.* it takes two to reproduce).

In summary, the process of inheritance can be generalized in two main steps: (1) the local traits and global shape of the face are mostly determined by genes controlling the production of proteins at the micro-level and (2) genes of an offspring are inherited from one parent or the other by a random selection and combination process. Thus, children are not identical to a single parent but tend to resemble both parents in various ways. The practical significance of predicting the appearance of a child from a parent pair should be acknowledged, and the existing methods based on single inputs should be christened limited and unrealistic.

To incorporate the concepts of genetics into the kinship generation problem, we utilize an encoder-decoder structure [13, 14] to mimic the process of inheritance in facial appearance by transforming genes from parents-to-child. Previously, the encoder-to-decoder structure has been incorporated into Generative Adversarial Network (GAN) [15] and Variational Autoencoder (VAE) [16] to generate photo-realistic faces [17], where mappings between facial images and high-level personal features were established. Similarly, in our kinship generation task, the facial traits of parents (*i.e.* an image pair) are translated into genes by the encoder. Then, the child's genes can be generated by simu-

lating the random selection and combination process on the gene-encodings of the parents. Finally, the face of the child is generated by decoding the genes.

We propose a kinship generation model with a two-step learning procedure inspired by the genetic process. Step one: a deep generative Conditional Adversarial Autoencoder (CAAE) [18] is trained on a large-scale face dataset to learn to map facial appearance to high-level features with knowledge of age and gender. Step two: a novel DNA-Net, trained on a smaller kinship dataset, transforms high-level features to genes, *i.e.* translates genes of a parent pair to a child. Figure 1 depicts the inputs and outputs of the proposed model. Can you determine which are the real children (bottom row) of the parents (top row)?

There are two main contributions in this paper.

1. We introduce DNA-Net to transfer features from parents to child by simulating the genetic process, while combining it with the CAAE model to realize child facial image generation from the images of parents.

2. We are able to generate multiple siblings by manipulating the gene codes in DNA-Net, which allow for changes to be made to the generated child in both age and gender.

*Beyond these contributions, we plan to promote our methodology with broader impacts through crowd-sourcing: Given enough data, our model will be able to reveal the mechanism and hidden factors of gene combination from parents, which is less random and more governed by natural laws.*

## 2. RELATED WORK

### 2.1. Kinship Verification

The task of kinship verification is to determine whether a face pair is related (*i.e.* KIN or NON-KIN). Evaluations are typically done separately for different relationship, like parent-child, siblings, and sometimes grandparent-grandchild. Research in both psychology and computer vision revealed that different kin relations render different familial features, which motivated researchers to model different relationship types independently. Existing methods for the kinship verification can generally be split into either metric learning based [19, 20] or feature based methods [11]. In metric learning methods, either a distance measure or feature transformation is learned to reduce distances between kin pairs and push away non-kin pairs. Feature based methods use hand-crafted features or learn more discriminative representations.

Recently, deep neural networks have achieved state-of-the-art in kinship verification. [21] proposed a method to discover the optimal features and metrics that relate a parent to offspring via gated autoencoders. [9] utilized CNNs
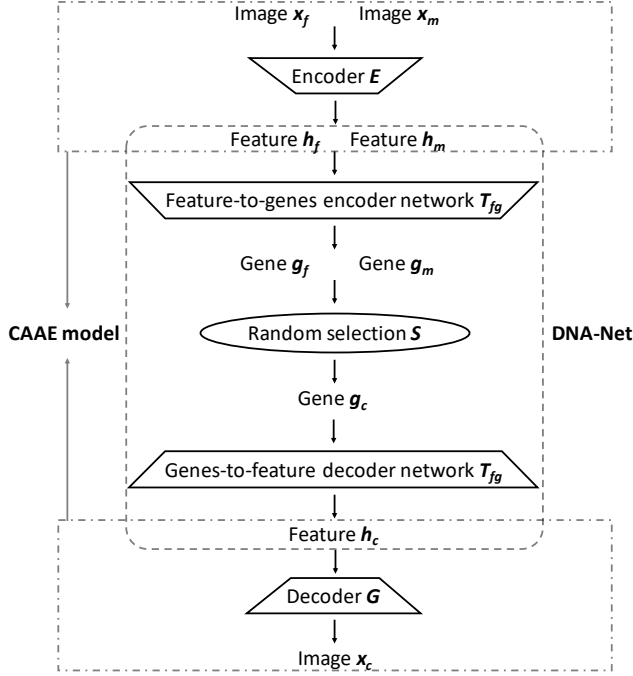
Figure 2: Flowchart of genetic model. Note that the variables used are consistent with that in Eq. 1-14 and Figure 3.

as a feature extractor for kinship verification. [22] integrated the triple ranking loss into CNN model to learn more discriminative representations.

Some methods incorporate deep metric learning for better performance. [23] proposed a denoising auto-encoder based on marginalized metric learning to preserve the structure of data and simultaneously endow the discriminative information into the learned features. [10] developed a discriminative deep multi-metric learning method to jointly learn multiple neural networks to better use the commonality of multiple feature descriptors. See past challenges for various other methods and task specific information [24].

## 2.2. Deep Generative Models

VAE and GAN are two of the most renowned image generation models. Both methods can generate images from latent codes that follow certain prior distributions. In recent years, multiple variants of these two have emerged. Some of them adopt an encoder-decoder structure that can also map images into latent codes which can be considered as features. In [25], Isola proposed pix2pix, an image-to-image translation method based on conditional GAN (cGAN) [26]. Pix2pix can be seen as learning two mappings, image to features and features to image. Then came inverted conditional GAN (IcGAN) [27], a two-step image-to-image translation method which focuses on face attributes editing, like transforming smiling face to non-smiling face. In IcGAN, an

additional encoder is trained to map a image into latent codes/features and conditional representation after a cGan was trained first. After the training of additional encoder, face attributes can be changed by manipulation of latent codes. In [28], a tag mapping net was proposed which maps tags (labels) of image to features which are encoded from image, making it possible to adjust the attributes of generated image by adjusting the tag. [29] proposed an image-to-image translation model which focuses on face attributes editing and can deal with multiple face attributes simultaneously. These works give us inspiration that the mapping between face image and face features can be learned in deep generative models [25, 27], even mapping between features and features can be learned (tags can be seen as kind of features) [28], and image content can be manipulated with latent codes [27].

A special variant of VAE and GAN is the combination of the two, with VAE/GAN [14] and AAE [13] being amongst the most popular. When used together, these models inherit the ability of inference from VAE and the tendancy to generate sharp pictures of GAN. Also, VAE/GAN and AAE have encoder-decoder structures.

## 3. APPROACH

First, we use neural network terminology to model the genetic process. Then, a CAAE model adapts to establish two-way mappings between facial images and face features. Finally, our DNA-Net establishes two-way mappings between face features and genes, *i.e.* analogous to inheritance.

### 3.1. Genetic Model

Research in genetics revealed that multiple genes could contribute to a single facial trait, for instance, 16 genes were found to effect eye color [30]. From this, translating from genes to face appearance is modeled as

$$x_k = F(g_{k1}, g_{k2}, ..., g_{kn}), \quad k \in \{f, m, c\} \qquad (1)$$

where $f$, $m$, $c$ stand for father, mother, and child, respectively, $x_k$ is the appearance of $k$, and $g_{ki}$ denotes the genes responsible for the facial features. $F(\cdot)$ produces a face based on gene(s). As shown in nature, a child genetically inherits genes from both parents via random selection. This random selection can be expressed as follows:

$$
\begin{aligned}
x_c &= F(g_{c1}, g_{c2}, ..., g_{cn}) \\
&= F(S(g_{f1}, g_{m1}), S(g_{f2}, g_{m2}), ..., S(g_{fn}, g_{mn})),
\end{aligned}
$$
(2)

where $S(\cdot)$ simulates the process of obtaining the gene of a child $g_{ci}$ through a random selection over the corresponding genes of the two parents, which is thus defined as

$$
\begin{aligned}
g_{ci} &= S(g_{fi}, g_{mi}) \\
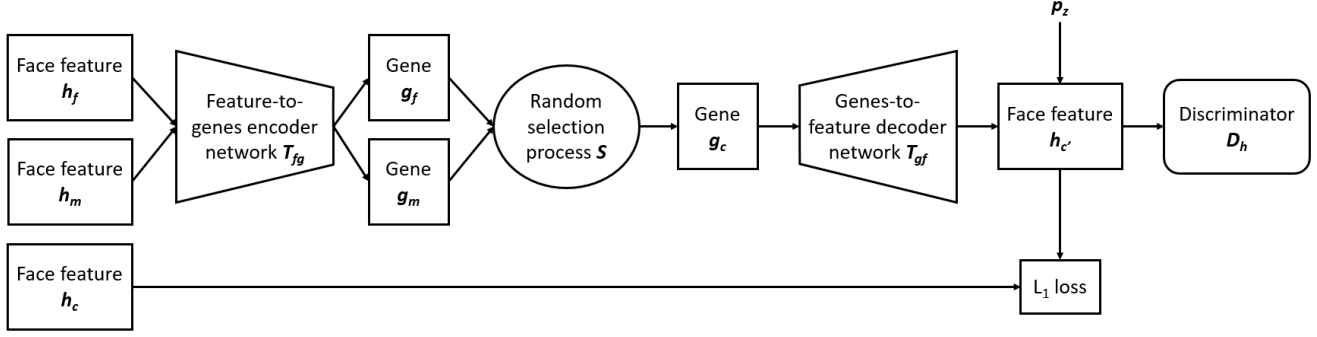&= r_i \cdot g_{fi} + (1 - r_i) \cdot g_{mi}, \quad r_i \in \{0, 1\},
\end{aligned}
$$
(3)

Figure 3: Architecture of DNA-Net. The encoder network $T_{fg}$ maps the extracted face feature $h$ to gene $g$. Random selection process $S(\cdot)$ transfer genes from parents to child. Decoder network $T_{gf}$ maps gene $g$ to feature $h$. The discriminator $D_h$ imposes the uniform distribution on $h$ and $p_z$ is a prior distribution. The network updated based on the $L_1$ loss between the input face feature $h_c$ and generated face feature $h_{c'}$ of child. Note that $f$, $m$, $c$ in the figure means father, mother, child respectively.

where $r_i$ is a value randomly assigned.

To incorporate the process of Eq. (1)-(3) into our generative model, we design a genetic model that generates a face of a child from faces of a pair of parents. This model contains three main stages. Figure 2 depicts this genetic model.

**First stage**. Genes of parents are predicted from their appearances. Specifically, we encode faces to represent $x_k$, $k \in \{f, m\}$ and generate personal facial features $h_k$ with encoder $E$ through

$$h_k = E(x_k), \quad k \in \{f, m\}. \quad (4)$$

Feature vectors $h_k$ will then be translated to gene vectors $g_k$ by another gene encoder $T_{fg}$ as

$$g_k = [g_{k1}, g_{k2}, ...] = T_{fg}(h_k), \quad k \in \{f, m\}. \quad (5)$$

**Second stage**. We derive the gene vector of the child $g_c$ from the genes of the parents via a random selection process over corresponding gene elements. This can be expressed as

$$g_c = [g_{c1}, g_{c2}, ...] = [S(g_{f1}, g_{m1}), S(g_{f2}, g_{m2}), ...]. \quad (6)$$

**Third stage**. We predict the facial appearance of the child $x_c$ from genes $g_c$ output from two decoders. Specifically, the personal facial feature $h_c$ is decoded from gene $g_c$ by a gene decoder $T_{gf}$ as

$$h_c = T_{gf}(g_c). \quad (7)$$

Then, the facial image is generated by another decoder $G$:

$$x_c = G(h_c). \quad (8)$$

Eq. (1) can be represented as $x_k = F(g_k) = G(T_{gf}(g_k))$.

We use CAAE [18] to train the image-feature encoder $E$ and decoder $G$. Then, a novel neural network dubbed DNA-Net was designed to model the mappings between extracted features and genes via $T_{fg}$ and $T_{gf}$ as well as the random selection process $S(\cdot, \cdot)$. One reason to use separate networks $T_{gf}(\cdot)$ and $G(\cdot)$, instead of a single network for $F(\cdot)$ (and vice versa) is that, the limited amount of data labeled for kinship recognition is less suited to support training of a single larger network that directly maps between images and genes (i.e. prone to overfitting). Instead, we choose to train $E(\cdot)$ and $G(\cdot)$ in the CAAE on a large-scale face dataset, and then train the smaller DNA-Net on the smaller kinship dataset. Besides, we want encoder $E$ and decoder $G$ to capture age and gender information, opposed to DNA-Net, as most genes are age-invariant.

### 3.2. Image-Feature Mapping via CAAE

Next, we discuss the details of CAAE [18]. The input and output of CAAE net are $128 \times 128$ RGB facial images $x \in R^{128 \times 128 \times 3}$. On the one hand, the encoder $E(\cdot)$ preserves the high-level personal features of the input face $x$ in a feature vector $h = E(x) \in R^n$. On the other hand, the decoder $G$ generates a face image $\hat{x} = G(h, l)$ that is conditioned on a certain age and gender. Note that $l$ is a one-hot vector encoding age and gender labels. In the end, the input and output faces aim to be as similar as possible:

$$\min_{E,G} L(x, G(E(x), l)), \quad (9)$$

where $L(\cdot, \cdot)$ denotes euclidean distance.

Additionally, two discriminator networks, $D_z$ and $D_{img}$, are placed after $E$ and $G$, respectively, for the purpose of adversarial training. $D_z$ regularizes the feature vector $h$ to be uniformly distributed to smooth the age transformation. We denote the distribution of the training data as $p_{data}(x)$,
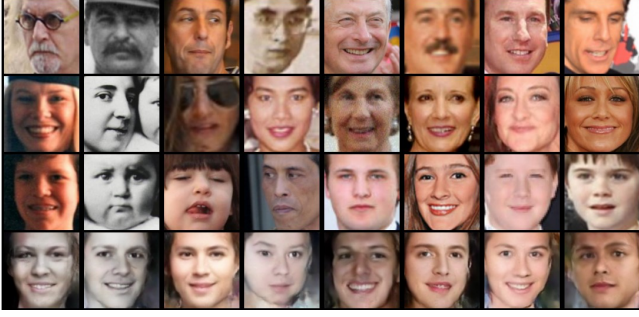
Figure 4: Samples results. Each column corresponds to a family, with faces of fathers on first row, mothers on second, real children on third, and generated children on bottom.

while the distribution of feature $h$ is $q(h|x)$. Also, $p(z)$ is assumed to be a prior distribution, and $z^* \sim p(z)$ denotes the random sampling process from $p(z)$. A min-max objective function can be used to train $E$ and $D_z$ as

$$\min_E \max_{D_z} \mathbb{E}_{z^* \sim p(z)}[logD_z(z^*)] + \\ \mathbb{E}_{x \sim p_{data}(x)}[log(1 - D_z(E(x)))]. \quad (10)$$

Besides, $D_{img}$ forces $G$ to generate photo-realistic and plausible faces for an arbitrary $h$ and $l$, which can be trained along with $G$ by a similar token with Eq. (10). Specifically,

$$\min_G \max_{D_{img}} \mathbb{E}_{x,l \sim p_{data}(x,l)}[logD_{img}(x,l)] + \\ \mathbb{E}_{x,l \sim p_{data}(x,l)}[log(1 - D_{img}(G(E(x),l)))]. \quad (11)$$

Finally the objective function becomes

$$\min_{E,G} \max_{D_z,D_{img}} L(x, G(E(x),l)) \\ + \mathbb{E}_{z^* \sim p(z)}[logD_z(z^*)] \\ + \mathbb{E}_{x \sim p_{data}(x)}[log(1 - D_z(E(x)))] \\ + \mathbb{E}_{x,l \sim p_{data}(x,l)}[logD_{img}(x,l)] \\ + \mathbb{E}_{x,l \sim p_{data}(x,l)}[log(1 - D_{img}(G(E(x),l)))]. \quad (12)$$

### 3.3. Genetic Mappings via DNA-Net

We propose DNA-Net to map face features of a pair of parents to a child (see Figure 3). As mentioned, DNA-Net is made-up of two networks, *i.e.* a feature-to-genes encoder network $T_{fg}$ and a genes-to-feature decoder network $T_{gf}$. During the encoding process, given an input face feature vector $h \in R^n$, $T_{fg}$ produces a gene vector $g \in R^m$, where $n$ and $m$ are dimensions of the feature vector and gene vector respectively. During the decoding process, given a gene vector $g \in R^m$, the decoder $T_{gf}$ will output a feature vector $h \in R^n$. For the complete generation process, $T_{fg}$ predicts the gene vectors for both parents, while $T_{gf}$ maps the genes-to-features for the child.

When the gene vectors of the parents are obtained from $T_{fg}$, there are two ways to implement random selection process in Eq. (3). Since the convergence of a neural network requires a certain structure, the randomness in $S(\cdot)$ should eliminate. This can be done in two ways: (1) use a determined random seed when training; (2) use a determined rule to select which parent will pass down which gene elements to child (*i.e.* the parent for which particular genes of the child are inherited). We follow (2) in this work. Specifically, our selection rule keeps the genes with maximum values of the two parents. During testing, along with the selection rule, the DNA-Net can also use a random 0-1 sequence for genes selection from parents to children to generate additional children (*i.e.* siblings).

The training process of DNA-Net is as follows. Given a triplet set of family images $(x_f, x_m, x_c)$, we first extract facial features $(h_f, h_m, h_c)$ from the trained encoder $E$ in Eq. (4). They are then used as the inputs and ground truth of DNA-Net. The objective of DNA-Net is to generate similar features as $h_c$. Therefore, the loss over the triplet set is defined as

$$\min_{T_{fg},T_{gf}} ||T_{gf}(S((T_{fg}(h_f), T_{fg}(h_m))) - h_c||_2 \quad (13)$$

Due to the uniform distribution constraint on $h$ in CAAE, the output of DNA-Net $h_c$ should also follow the same distribution. So, a discriminator $D_h$ is trained along with $T_{fg}$ and $T_{gf}$. The loss that regularizes DNA-Net's output is defined as

$$\min_{T_{fg},T_{gf}} \max_{D_h} \mathbb{E}_{z^* \sim p(z)}[logD_h(z^*)] + \\ \mathbb{E}_{h_c \sim T(h_f,h_m)}[log(1 - D_h(T(h_f,h_m))], \quad (14)$$

$T(h_f, h_m) = T_{gf}(S((T_{fg}(h_f), T_{fg}(h_m)))$ is the output.

## 4. EXPERIMENTS

This section first introduces the data, and then details the implementation. Also, our model is evaluated qualitatively and quantitatively in several experiments, specifically, conditional face generation, kinship verification, human evaluation, and heritable mappings.

### 4.1. Datasets

**UTKFace** [18] is used to train CAAE model, which divides the images into 10 age groups (*i.e.* 0-5, 6-10, 11-15, 16-20, 21-30, 31-40, 41-50, 51-60, 61-70, and 71-80 years old). For this, a 10-dim one-hot vector is used to represent the age. For the gender, another 10-dim one-hot vector is formulated. UTKFace datasets is a large-scale face dataset with wide age span (ranging from 0 to 116 years old), containing over 20,000 aligned and cropped face images with labels for age and gender.

(a) Across ages (*i.e.* 10, 20, 30 years old from row 4-6, respectfully).



(b) Across gender (*i.e.* male-to-female from row 4-5, respectfully).

Figure 5: First three rows are real families face images which are similar to Figure 4. The last three and two rows are generated face images with different ages (a) and gender (b).

**FIW** [31, 32] contains 1,000 families, over 11,000 persons, and is the largest kinship recognition dataset up to date. This gave us 1,997 father-mother-child face sets selected at random, with 1,600 used for training and the remaining 397 for testing.

## 4.2. Implementation

The implementation of CAAE is the same as [18]. With the CAAE model trained, the feature vectors $h_x$ of all faces in the father-mother-child sets could be generated and used to train the DNA-Net. In our experiment, dimensions of the feature vectors $h_x \in R^n$ and genes vectors $g_x \in R^m$ are both set to $n = m = 100$. In DNA-Net, $T_{fg}$ and $T_{gf}$ are both 3-layer fully connected networks. CAAE and DNA-Net were optimized using Adam optimizer [33] with a learning rate of 0.0001.
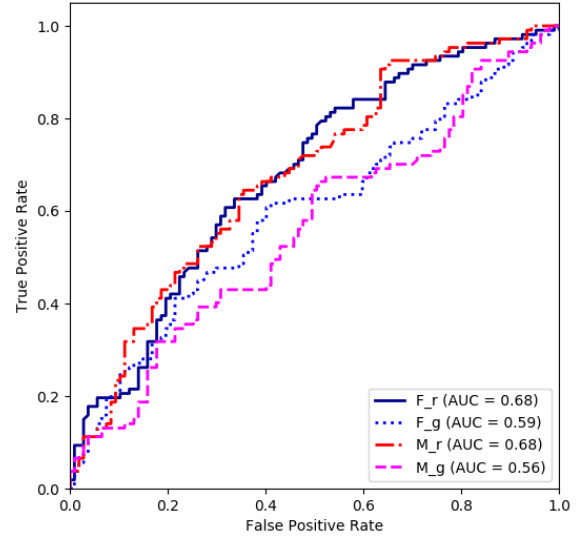


Figure 6: ROC curve for verification evaluation. Legend items translate to father (F) or mother (M) and real (r) or generated (g) children.

## 4.3. Face Generation in Multiple Conditions

Figure 4 shows examples of generated children's face images. As can be seen, the generated images have a high visual quality and clearly resemble one of the parents. For example, the mouth or eyes of generated children's face look like either their father (*e.g.* fourth column) or mother (*e.g.* second column). All results are with high quality, indicating that DNA-Net learns a mapping from feature space.

Benefiting from the novel two-stage generation process, our model can generate children faces at different ages and genders by changing input of age and gender labels. Samples of children with different ages are shown in Figure 5a, and those in different genders are shown in Figure 5b. Clearly, we can observe the aging progress from juvenile to young people to middle-age in row 4-6 (*e.g.* second column, Figure 5a).

We also generate sibling faces of the child by using sequences from the random selection process instead of the determined rule for training. Some generated sibling faces are shown in Figure 7.

## 4.4. Quantitative evaluation

To quantify the performance of the proposed, we evaluated via kinship verification. Both learning model (CNN) and human subject performance are evaluated. Experimental settings and results are described in the following subsections and shown in Fig 9.

Figure 7: Samples of sibling generation. First three rows are real face images of families like in Figure 4. The last four rows are child faces generated with different random seeds.
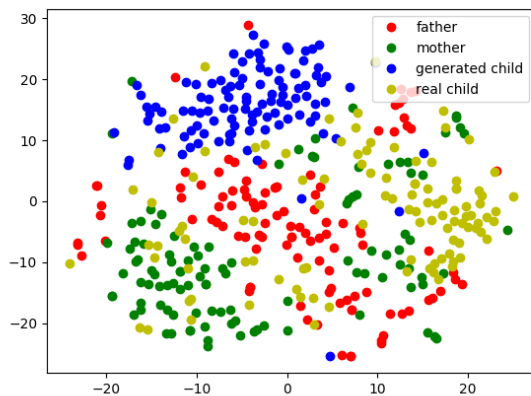


Figure 8: Visualization of facial feature distribution of fathers, mothers, children, and generated ones. Red points represent the feature of fathers, green for mothers, yellow for real children, blue for generated children, respectively. Best viewed in color.

### 4.4.1 Kinship Verification

To evaluate the quality of generated faces, we used a pretrained kinship verification classifier to identify whether the generated child's image can be classified as the child of a given parent. The more generated images that can fool the classifier, the better the performance of the generation method. In this paper, the pre-trained kinship verification classifier is a FaceNet network fine-tuned on FIW [31].

We randomly sampled 100 families, with each consisting of a mother, father, and child. For each set of parents, a child's face was generated using our model. We then evaluated kinship verification accuracy on both the real and generated face images, with another 100 negative samples added to the test set. Thus, the same number of negatives was used for both the real and the generated cases. The generated children faces scored a verification accuracy of 58.89% (with father) and 57.01% (with mother), while the real children achieved 67.29% and 73.83%, respectively. Figure 6 shows the ROC curves for each cases.

To measure the identity similarity between the real child images and the generated ones, we use pre-trained FaceNet model to extract identity features from both faces, where the training data are totally independent from FIW. Then, a similarity score is computed between every two extracted features using cosine distance. The average distance of 100 real-to-generated pairs is 0.90, compared to 0.94 for generated faces and random real faces. This means the generated faces are a little closer to the real ones. In addition, we visualize the low-dimensional distribution of facial features from generated faces, real ones, and parents respectively by t-SNE [34]. Figure 8 shows the distribution of the face features– those of generated children are more clustered, and with small overlap with real ones. This may be due to lack of large training images and complex genetic mechanism. However, we can see that the feature distributions of the generated child face is as close to the faces of the parents as it is to the faces of the real child. This is consistent with the verification results.
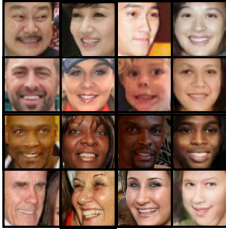
### 4.4.2 Human Evaluation

We asked human participants to vote on child images (real or generated) that were thought to be the true child of a pair of parents. In other words, we randomly selected 30 parent pairs from the verification set. Thus, facial images of each parent pair were shown next to their actual alongside the generated (order of the actual and generated faces were set at random, while the father was proceeded by the mother on the left side). The task was to determine the child that was the descendent of the parent pair. In other words, the volunteers picked the face of the child that resembled the parents more. Hence, each pair included a generated face of the child. We created a Google Form to distribute, and used university email lists and social media for recruiting volunteers. In total, 35 volunteers partook. Note that no volunteers had prior knowledge that some of the faces were generated (*i.e.* we just asked which child is the true descendant).

The generated children obtained more votes than the actual. Specifically, about 60.29% of the generated stumped the user into believing it was the true child, which was mea-

| Pair-Type | CNN | Human |
|---|---|---|
| Father-Real | 67.29 | 38.88 |
| Father-Gene. | 58.89 | 61.12 |
| Mother-Real | 73.83 | 40.55 |
| Mother-Gene. | 57.01 | 59.45 |
| Avg.-Real | 70.56 | 39.71 |
| Avg.-Gene. | 57.95 | 60.29 |

(a)

(b)

(a)

Figure 9: Kinship verification scores (%) for real and generated children (a). Face samples shown are those that CNN and most humans agree (b): parents (columns 1-2 are father and mother, respectfully) and children (columns 3-4 are actual and generated child, respectfully). Top 3 rows are samples of generated children scored highest and accumulated most votes. To the contrary, the bottom row are real children that scored highest and received most votes.

sured by the number of votes. Thus, the faces generated by the proposed appeared more genuine than that of the actual child to humans (see Figure 9b).

### 4.4.3 Heritability Maps

It is evident that the human face consists of complex traits under strong genetic control. To further explore heritability of facial traits, we study the geometric similarity of face image pairs. Here, we compare the shape features of four parts of face, *i.e.* eyes, nose, mouth, and chin, between parents and child. In detail, we select 20 pairs of front faces of real child and parents, generated child and parents, respectively, from above testing images. We detect the landmarks of faces and connect them into lines. After that, Hu invariant moment ( [35]) is computed to represent the shapes of the four facial parts. Accumulative cosine distances are then utilized to represent heritablility. Figure 10a shows the heritability map of generated child face. It can be seen that mouth region has high similarity with parents. For the real child face (see Figure 10b), like the mouth, the nose region is highly similar. Besides, the chin regions are potential evidence for genetics. These results are consistent with findings in genetics [36].

## 5. Conclusion

In this paper, we investigate a multidisciplinary problem of children face generation from their parents which resides in the intersection of computer vision, biology and genetics. We hope to open a gate for visual face modeling for genetic combination and expression. To this end, we propose a novel DNA-Net to construct the transformation and random selection process from parents' genes to child's

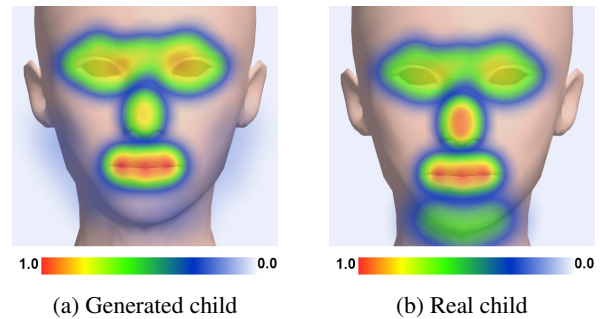(a) Generated child          (b) Real child

Figure 10: Heritability map represents the estimated salience about facial landmarks. Best viewed in color.

ones. Furthermore, our model could generate face images of children of different ages and genders by the leverage of CAAE model. Quantitative and qualitative experimental results show the generated children faces have high similarity with parents as well as similar heritability with real children. Our study could be useful in a varity of applications, ranging from population genetics and gene-mapping studies, to face modeling and reconstruction applications.

## References

[1] Junkang Zhang, Siyu Xia, Ming Shao, and Yun Fu. Family photo recognition via multiple instance learning. In *ACM on ICMR*, 2017.

[2] Chao Xia, Siyu Xia, Yuan Zhou, Le Zhang, and Ming Shao. Graph based family relationship recognition from a single image. In *PRCAI*, 2018.

[3] Daniel JM Crouch, Bruce Winney, Willem P Koppen, William J Christmas, Katarzyna Hutnik, Tammy Day, Devendra Meena, Abdelhamid Boumertit, Pirro Hysi, Ayrun Nessa, et al. Genetics of the human face: Identification of large-effect single gene variants. *Proceedings of the National Academy of Sciences*, 115(4):E676–E685, 2018.

[4] Alexandra Alvergne, Charlotte Faurie, and Michel Raymond. Differential facial resemblance of young children to their parents: who do children look like more? *Evolution and Human behavior*, 28(2):135–144, 2007.

[5] Lisa M DeBruine, Benedict C Jones, Anthony C Little, and David I Perrett. Social perception of facial resemblance in humans. *Archives of sexual behavior*, 37(1):64–77, 2008.

[6] Martin Daly and Margo I Wilson. Whom are newborn babies said to resemble? *Ethology and Sociobiology*, 3(2):69–78, 1982.

[7] Farhad B Naini and James P Moss. Three-dimensional assessment of the relative contribution of genetics and environment to various facial parameters with the twin method. *American Journal of Orthodontics and Dentofacial Orthopedics*, 126(6):655–665, 2004.

[8] Dimosthenis Tsagkrasoulis, Pirro Hysi, Tim Spector, and Giovanni Montana. Heritability maps of human face mor-

phology through large-scale automated three-dimensional phenotyping. *Scientific reports*, 7:45885, 2017.

[9] Kaihao Zhang, Yongzhen Huang, Chunfeng Song, Wu Hong, and Wang Liang. Kinship verification with deep convolutional neural networks. In *BMVC*, 2015.

[10] Jiwen Lu, Junlin Hu, and Yap-Peng Tan. Discriminative deep metric learning for face and kinship verification. *IEEE TIP*, 26(9):4269–4282, 2017.

[11] I. . Ertugrul and H. Dibeklioglu. What will your future child look like? modeling and synthesis of hereditary patterns of facial dynamics. In *FG*, pages 33–40, 2017.

[12] Savas Ozkan and Akin Orkan. Kinshipgan: Synthesizing of kinship faces from family photos by regularizing a deep face network. In *ICIP*, pages 2142–2146. IEEE, 2018.

[13] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015.

[14] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. *arXiv preprint arXiv:1512.09300*, 2015.

[15] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.

[16] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[17] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.

[18] Zhifei Zhang, Yang Song, and Hairong Qi. Age progression/regression by conditional adversarial autoencoder. In *CVPR*, pages 5810–5818, 2017.

[19] Jiwen Lu, Xiuzhuang Zhou, Yap-Pen Tan, Yuanyuan Shang, and Jie Zhou. Neighborhood repulsed metric learning for kinship verification. *IEEE TPAMI*, 36(2):331–345, 2014.

[20] Haibin Yan, Jiwen Lu, Weihong Deng, and Xiuzhuang Zhou. Discriminative multimetric learning for kinship verification. *IEEE TIFS*, 9(7):1169–1178, 2014.

[21] Afshin Dehghan, Enrique G Ortiz, Ruben Villegas, and Mubarak Shah. Who do i look like? determining parent-offspring resemblance via gated autoencoders. In *CVPR*, pages 1757–1764, 2014.

[22] Mengyin Wang, Jiashi Feng, Xiangbo Shu, Zequn Jie, and Jinhui Tang. Photo to family tree: Deep kinship understanding for nuclear family photos. In *Proceedings of the Joint Workshop of the 4th Workshop on Affective Social Multimedia Computing*, pages 41–46. ACM, 2018.

[23] Shuyang Wang, Joseph P Robinson, and Yun Fu. Kinship verification on families in the wild with marginalized denoising metric learning. In *FG*, pages 216–221. IEEE, 2017.

[24] Jiwen Lu, Junlin Hu, Venice Erin Liong, Xiuzhuang Zhou, Andrea Bottino, Ihtesham Ul Islam, Tiago Figueiredo Vieira, Xiaoqian Qin, Xiaoyang Tan, Songcan Chen, et al. The kinship verification in the wild evaluation. In *FG*, volume 1, pages 1–7. IEEE, 2015.

[25] 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.

[26] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.

[27] Guim Perarnau, Joost Van De Weijer, Bogdan Raducanu, and Jose M Álvarez. Invertible conditional gans for image editing. *arXiv preprint arXiv:1611.06355*, 2016.

[28] Chaoyue Wang, Chaohui Wang, Chang Xu, and Dacheng Tao. Tag disentangled generative adversarial networks for object image re-rendering. In *IJCAI*, 2017.

[29] Taihong Xiao, Jiapeng Hong, and Jinwen Ma. Elegant: Exchanging latent encodings with gan for transferring multiple face attributes. In *ECCV*, pages 168–184, 2018.

[30] Désirée White and Montserrat Rabago-Smith. Genotype–phenotype associations and human eye color. *Journal of human genetics*, 56(1):5, 2011.

[31] 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.

[32] Joseph P Robinson, Ming Shao, Yue Wu, and Yun Fu. Families in the wild (fiw): Large-scale kinship image database and benchmarks. In *ACM MM*, pages 242–246. ACM, 2016.

[33] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2015.

[34] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11):2579–2605, 2008.

[35] Ming-Kuei Hu. Visual pattern recognition by moment invariants. *IRE transactions on information theory*, 8(2):179–187, 1962.

[36] Lisa M DeBruine, Finlay G Smith, Benedict C Jones, S Craig Roberts, Marion Petrie, and Tim D Spector. Kin recognition signals in adult faces. *Vision research*, 49(1):38–43, 2009.