# DEEP KINSHIP VERIFICATION VIA APPEARANCE-SHAPE JOINT PREDICTION AND ADAPTATION-BASED APPROACH

*Heming Zhang[1], Xiaolong Wang[2], C.-C. Jay Kuo[1]*

[1] University of Southern California
Los Angeles, CA, USA

[2] Samsung Research America
Mountain View, CA, USA

## ABSTRACT

Kinship verification aims to identify the kin relation between two given face images. It is a very challenging problem due to the lack of training data and facial similarity variations between kinship pairs. In this work, we build a novel appearance and shape based deep learning pipeline. First we adopt the knowledge learned from general face recognition network to learn general facial features. Afterwards, we learn kinship oriented appearance and shape features from kinship pairs and combine them for the final prediction. We have evaluated the model performance on a widely used popular benchmark and demonstrated the superiority over the state-of-the-art.

***Index Terms***— Kinship verification, Deep learning

## 1. INTRODUCTION

Human facial image analysis has attracted lots of interests from image processing and computer vision community for a long time, e.g., face recognition [1, 2], face detection [3, 4], facial attributes perception (age[5], gender [6], etc), landmark detection [7]. Compared to general facial image analysis, kinship verification starts to attract our attention recently. There are many potential applications associated with it, such as finding missing children and social data mining [8]. However, identifying the underlying kin relation from facial images is very challenging, even for humans [9].

Recently, with the advancement of deep learning technology and big data, we have observed significant improvements in these general facial image analysis problems. Nevertheless, the performance of kinship verification is still not satisfied. One reason is because of the challenge of collecting large kinship data. Unlike general face recognition database collection, we need to collect pairs of facial images with kin relations. Another reason is due to similarity variations among people with kin relations. Even for the same family, e.g., the similarity between the father-son and father-daughter is usually different.

Most initial work mainly utilizes hand-crafted features, and then learn the prediction with metric learning approach. For example, in [8], a projection based metric with large margin (NRML) is learned and hand-crafted features such as LBP,

HOG, and SIFT are utilized. These features are widely used in many applications, but it is very hard to catch these underlying kin relations.

Since deep learning has gained success in many fields [10, 11, 12], recent work tend to utilize features from deep neural network (DNN). However, extracting DNN features from limited kinship data is not trivial. It results in either shallow network architectures[13], or stacked auto-encoder network with hand-crafted features as inputs [14]. In either cases, the power of DNN is degraded which makes it hard to capture these disriminative kin relation features. In [15], they relabel the kinship data and feed it into a deep neural network to fine-tune the face recognition model, in which large-scale face recognition data is utilized to help the training. Their experiments demonstrated an improvement when combing deep facial features with NRML [8].

As we know, the data used in kinship verification are also face images [8]. There should be common features between general face recognition and kinship verification. Training deep neural face recognition network via enough face data can utilize the advantage of deep networks in extracting comprehensive facial appearance features. These features obtained from face recognition task are aimed for person identification. Therefore, they mainly contain salient appearance information that help differentiate people.

In our case, we not only need the salient appearance information, but also the global appearance and shape information that helps find the underlying kin relation. Therefore, we enforce the feature learning capability via kinship oriented learning framework which has demonstrated the improvement [13, 14]. Moreover, we also combine the shape information with the appearance cues together. We are motivated by the observation that facial geometry information is also shared between most kinship face pairs, such as eye shape, nose structure, etc..

In this work, we utilize both appearance and shape information for kinship verification. Specifically, we design a convolutional neural network for deep appearance and shape feature extraction, comparison and prediction. To overcome the issue with limited kinship data, we propose an adaptation-based two-phase training scheme which utilizes large-scale face recognition data.
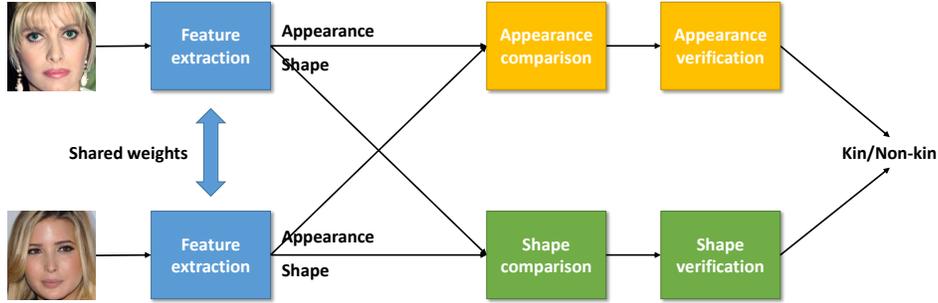
**Fig. 1**: Pipeline of our proposed method

## 2. PROPOSED METHOD

### 2.1. Overview of proposed method

In this section, we will talk about the proposed kinship verification scheme in details. The whole framework is demonstrated in Figure 1. Given two facial images, we first extract appearance and shape features for each image separately. Then we extract the facial appearance and shape features for each pair of images by comparing their individual features. Afterwards, we fuse appearance and shape results together to predict the kinship relation for the input image pairs.

The feature extraction module consists of a convolutional neural network as the backbone and the verification modules involve multiple fully connected layers as kinship classifiers. The comparison modules will be explained in details in Section 2.2. To train this modulized network with limited kinship data, we propose an adaptation-based two-phase training scheme in Section 2.3.

### 2.2. Appearance & Affine-invariant Shape Comparison

In this section, we explain the proposed two feature comparison modules for appearance and shape features. In these modules, we construct features for image pairs using comparison between features individually extracted from single images.

**Appearance Comparison (AC)**
The appearance features extracted from the previous module mainly contain essential identification information. To compare appearance features from two different people, we choose to perform element-wise multiplication on two appearance features. It resembles the weighted correlation when combined with fully connected layers in the later stage, and easily fits into current deep learning framework.

**Affine-invariant Shape Comparison (AISC)**
The shape of a face can be represented by a set of facial landmarks. However, the geometry of facial landmarks is very sensitive to pose and view variations. Obviously we do not want these changes to affect the comparison performance between two facial shapes.

Inspired by [16, 17], we use the Grassmann representation as the affine-invariant shape representation. A Grassmann manifold $G_{m,k}$ is the space in which the points are k planes in $\mathbb{R}^m$ [18]. We can associate an $m \times k$ orthogonal matrix $U$ to each $k$-plane $\nu$ in $\mathbb{R}^m$, such that the columns of $U$ form an orthonormal basis for $\nu$.

Given a matrix of coordinates of m facial landmarks $S = [(x_1, y_1); (x2, y2); \cdots; (x_m, y_m)]$, we can obtain an affine-transformed matrix $S'$ by applying a $2 \times 2$ full rank affine transformation matrix $A$ on the right, i.e. $S' = SA$. It is worth noting that column spaces spanned by $S$ and $S'$ are the same. In other words, the 2-D subspace spanned by $S' = SA$ is invariant to $A$, and thus all $S'$ maps to the same point on the Grassmann manifold. Consequently, the Grassmann representation of $S$ is invariant to affine transformations.

Using the Grassmann representation, the comparison between two shape matrices becomes analyzing the geodesic between two points on the Grassmann manifold. The geodesic between two subspaces represented by projectors $P_0$ and $P_1$ on the Grassmann manifold has the form [19]:

$$P_1 = \exp(tX)P_0\exp(-tX). \tag{1}$$

The matrix $X$ can be derived by the eigen-decompsition of $B = P_0 - P_1$. Therefore $B$ contains all the information of the geodesic between the two shapes. We can thus use $B$ as the feature for affine-invariant shape comparison as in [16, 17].

**Forward computation of AISC.** Given two shape matrices $S_0, S_1 \in \mathbb{R}^{m \times 2}$, where each shape matrix contains the $x$ and $y$ coordinates of $m$ landmark points, we compute their SVD as

$$S_0 = U_0 D_0 V_0^T, \quad S_1 = U_1 D_1 V_1^T. \tag{2}$$

The corresponding projection matrices are

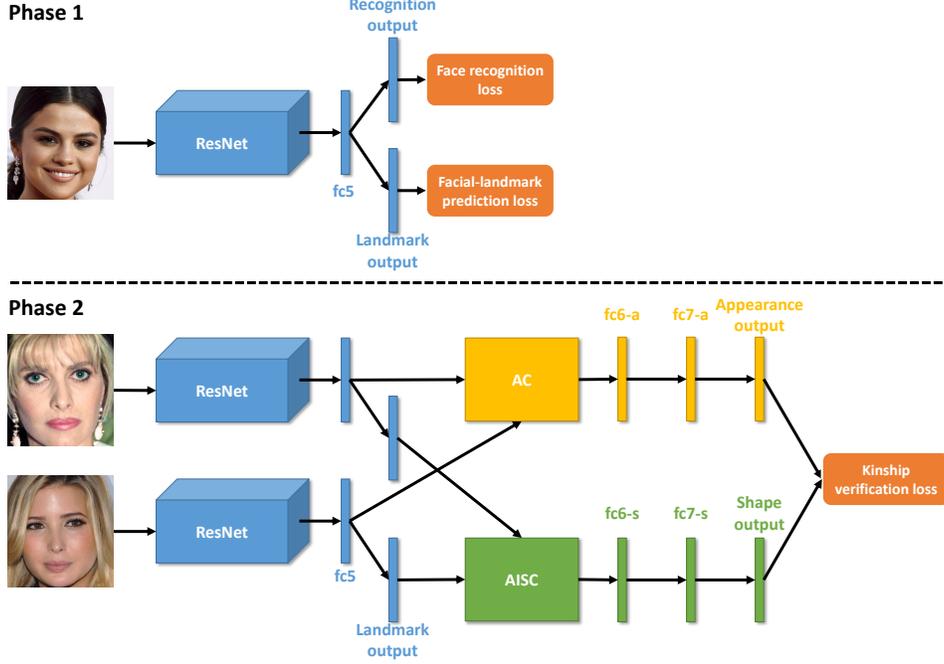$$P_0 = U_0 U_0^T, \quad P_1 = U_1 U_1^T. \tag{3}$$

**Fig. 2**: Our proposed two-phase training scheme. In phase 1, we utilize large-scale face recognition data to train the feature extraction module. In phase 2, the entire network is retrained on small-scale kinship data.

The shape comparison is conducted as $B = P_0 - P_1$.

**Backward computation of AISC.** Given the gradient of the loss function $\mathcal{L}$ respective to $B$, the gradient computations for the Grassmann block are given in Eqs. 4 to 9.

$$\frac{\partial \mathcal{L}}{\partial P_0} = \frac{\partial \mathcal{L}}{\partial B}, \quad \frac{\partial \mathcal{L}}{\partial P_1} = -\frac{\partial \mathcal{L}}{\partial B}. \tag{4}$$

Since the procedures of computing $P_0$ and $P_1$ from $S_0$ and $S_1$ in the forward pass are the same, respectively, we omit the subscript 0 and 1 in the following equations for simplicity.

$$\frac{\partial \mathcal{L}}{\partial U} = 2\frac{\partial \mathcal{L}}{\partial P}U. \tag{5}$$

As proved in [20], the Jacobian of the SVD can be computed as

$$\frac{\partial U}{\partial S_{ij}} = U\Omega_U^{ij}, \quad \frac{\partial V}{\partial S_{ij}} = -V\Omega_V^{ij}, \tag{6}$$

where $S_{ij}$ is the element on the $i$-th row and $j$-th column of $S$, $\Omega_U^{ij}$ and $\Omega_V^{ij}$ are given by

$$\Omega_U^{ij} = U^T\frac{\partial U}{\partial S_{ij}}, \quad \Omega_V^{ij} = \frac{\partial V}{\partial S_{ij}}^T V. \tag{7}$$

The elements of the matrices $\Omega_U^{ij}$ and $\Omega_V^{ij}$ can be computed by solving the linear systems in Eq. 8.

$$\begin{cases} D_l\Omega_{U,kl}^{ij} + D_k\Omega_{V,kl}^{ij} = U_{ik}V_{jl} \\ D_k\Omega_{U,kl}^{ij} + D_l\Omega_{V,kl}^{ij} = -U_{il}V_{jk} \end{cases}, \tag{8}$$

where $D_l$ is the $l$-th diagonal element in $D$, $\Omega_{U,kl}^{ij}$ is the element on the $k$-th row and $l$-th column of $\Omega_U^{ij}$ and $U_{ik}$ is the element on the $i$-th row and $k$-th column of $U$.

From Eqs. 6 through 8, we can derive the gradient for the shape matrix as

$$\frac{\partial \mathcal{L}}{\partial S_{ij}} = \sum_m \sum_n \Gamma_{mn},$$
$$\Gamma = \frac{\partial \mathcal{L}}{\partial U} \circ (U\Omega_U^{ij}), \tag{9}$$
$$\Omega_{U,kl}^{ij} = \frac{D_lU_{ik}V_{jl} + D_kU_{il}V_{jk}}{D_l^2 + D_k^2},$$

where $\circ$ denotes the Hadamard product.

### 2.3. Adaptation-based Two-phase Training Scheme

Most previous work tends to train a shallow network with small-scale kinship data [13, 14]. To overcome this issue, inspired by [15], we utilize large-scale face recognition data for training.

In phase 1, we perform a joint training on face recognition and facial landmark prediction on large-scale face recognition data. The aim of this phase is to train a good feature extraction module for both appearance and shape features.

In phase 2, we omit the recognition output and add the feature comparison and verification modules to the pre-trained feature extraction module. The network is then trained on

small-scale kinship data. The proposed adaptation-based two-phase training scheme for our deep network is illustrated in Figure 2.

Comparing to the adaptation approach in [15], our training scheme has two major differences. Firstly, our pre-training in phase 1 involves a joint training on face recognition and facial landmark prediction. This joint training not only helps extract more discriminative appearance features, but also provides extra shape features to further improve the accuracy and robustness of our model.

Secondly, after the pre-training on large-scale face recognition data in [15], the same network is adopted except for the output layer. To further re-train the face recognition network, each positive kin pair of images are manually assigned with a different identity label. In this way, the network is potentially forced to identify different people as the same person, as well as identify the same person as different people. On the contrary, our modulized network can be adapted to different training data and directly trained on kinship data.

## 3. EXPERIMENTS

### 3.1. Experimental Setup

**Datasets.** To have a fair comparison, we follow the same experimental setting as used in [15] where large-scale face recognition dataset CASIA WebFace is used for adaptation. WebFace contains 500,000 images of 10,000 subjects. For the landmark label, we used dlib [21] to predict the facial landmarks as pseudo labels. We conducted our experiments on the small-scale kinship dataset KinFaceW-I [8], which contains 1k images. It provides more than 500 pairs of positive and negative samples of four types of kin relationships: father-son (F-S), father-daughter (F-D), mother-son (M-S) and mother-daughter (M-D).

**Training details.** The feature extraction module of our network is first trained on WebFace jointly for face recognition and landmark prediction with cross-entropy and mean-squared error as loss functions, respectively. We adopted the ResNet [12] architecture and trained it from scratch. Then we add the comparison and verification modules and train the network on KinFaceW-I via 5-fold cross validation.

### 3.2. Comparison with Previous Work

Our experimental results on KinFaceW-I are listed in Table 1 together with results from previous work.

One can observe that our proposed method outperforms humans' results [9] as well as handcrafted features [8, 9]. Comparing with shallow network [13] or inputs with reduced dimension [14], we also demonstrate large improvement.

We achieve the best performance on father-son, father-daughter, mother-son tasks and mean accuracy over four kin relations when comparing with deep network [15].

| Method | F-S | F-D | M-S | M-D | Mean |
|---|---|---|---|---|---|
| Human A [9] | 62.0 | 60.0 | 68.0 | 72.0 | 65.6 |
| Human B [9] | 68.0 | 66.5 | 74.0 | 75.0 | 70.9 |
| MNRML [8] | 72.5 | 66.5 | 66.2 | 72.0 | 69.6 |
| MPDFL [9] | 73.5 | 67.5 | 66.1 | 73.1 | 70.1 |
| DKV [14] | 71.8 | 62.7 | 66.4 | 66.6 | 66.9 |
| SMCNN [13] | 75.0 | 75.0 | 68.7 | 72.2 | 72.7 |
| CFT* [15] | 78.8 | 71.7 | 77.2 | **81.9** | 77.4 |
| Ours | **81.8** | **76.6** | **77.5** | 77.2 | **78.3** |

**Table 1**: Mean verification accuracy (%) on KinFaceW-I dataset.

| Setting | Joint pre-training | Appearance | Shape | Mean |
|---|---|---|---|---|
| A | | √ | | 65.9 |
| B | √ | √ | | 69.6 |
| C | √ | | √ | 68.0 |
| D | √ | √ | √ | **78.3** |

**Table 2**: Ablation study

### 3.3. Ablation Study

To investigate the contribution to the performance of each component, we conducted an ablation study. We provide the results in Table 2, in which the setting D is our proposed setting.

The comparison between setting A and B demonstrates the benefit of joint pre-training with landmark objective, which also enables the utilization of shape information. By comparing settings B, C, D, one can find out that the appearance and shape information are complementary to each other and thus the fusion achieves a significant improvement in performance. Without the utilization of both appearance and shape information, our adaptation-based approach performs worse than the previous adaptation-based approach [15]. The reason is that we used a much deeper backbone network structure compared to that in [15]. Consequently, without training both the appearance and shape branches, the network is easily overfitting to the limited kinship data.

## 4. CONCLUSION

Kinship verification is a challenging task on which even human cannot perform well. In this work we proposed to use deep learning techniques for its promising performance. We utilized two types of complementary information, namely appearance and shape. Both of them are essential for identifying the underlying kin-relation from facial images. To enable the training of such deep network with limited kinship data, we proposed an adaptation-based two-phase training scheme and utilized large-scale face recognition data. We demonstrated the superiority of our proposed method over the state-of-the-arts.

# 5. REFERENCES

[1] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, et al., "Deep face recognition.," in *BMVC*, 2015, vol. 1, p. 6.

[2] Yandong Wen, Kaipeng Zhang, Zhifeng Li, and Yu Qiao, "A discriminative feature learning approach for deep face recognition," in *European Conference on Computer Vision*. Springer, 2016, pp. 499–515.

[3] Peiyun Hu and Deva Ramanan, "Finding tiny faces," in *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*. IEEE, 2017, pp. 1522–1530.

[4] Zekun Hao, Yu Liu, Hongwei Qin, Junjie Yan, Xiu Li, and Xiaolin Hu, "Scale-aware face detection," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, vol. 3.

[5] Xiaolong Wang, Rui Guo, and Chandra Kambhamettu, "Deeply-learned feature for age estimation," in *2015 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2015, pp. 534–541.

[6] Xiaolong Wang and Chandra Kambhamettu, "Gender classification of depth images based on shape and texture analysis," in *Global Conference on Signal and Information Processing (GlobalSIP), 2013 IEEE*. IEEE, 2013, pp. 1077–1080.

[7] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang, "Facial landmark detection by deep multi-task learning," in *European Conference on Computer Vision*. Springer, 2014, pp. 94–108.

[8] Jiwen Lu, Xiuzhuang Zhou, Yap-Pen Tan, Yuanyuan Shang, and Jie Zhou, "Neighborhood repulsed metric learning for kinship verification," *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 2, pp. 331–345, 2014.

[9] Haibin Yan, Jiwen Lu, and Xiuzhuang Zhou, "Prototype-based discriminative feature learning for kinship verification," *IEEE transactions on cybernetics*, vol. 45, no. 11, pp. 2535–2545, 2015.

[10] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[11] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," in *International Conference on Learning Representations*, 2015.

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[13] Lei Li, Xiaoyi Feng, Xiaoting Wu, Zhaoqiang Xia, and Abdenour Hadid, "Kinship verification from faces via similarity metric based convolutional neural network," in *International Conference Image Analysis and Recognition*. Springer, 2016, pp. 539–548.

[14] Mengyin Wang, Zechao Li, Xiangbo Shu, Jinhui Tang, et al., "Deep kinship verification," in *Multimedia Signal Processing (MMSP), 2015 IEEE 17th International Workshop on*. IEEE, 2015, pp. 1–6.

[15] Qingyan Duan, Lei Zhang, and Wangmeng Zuo, "From face recognition to kinship verification: An adaptation approach," in *Computer Vision Workshop (ICCVW), 2017 IEEE International Conference on*. IEEE, 2017, pp. 1590–1598.

[16] Tao Wu, Pavan Turaga, and Rama Chellappa, "Age estimation and face verification across aging using landmarks," *IEEE Transactions on Information Forensics and Security*, vol. 7, no. 6, pp. 1780–1788, 2012.

[17] Xiaolong Wang and Chandra Kambhamettu, "Leveraging appearance and geometry for kinship verification," in *Image Processing (ICIP), 2014 IEEE International Conference on*. IEEE, 2014, pp. 5017–5021.

[18] Alan Edelman, Tomás A Arias, and Steven T Smith, "The geometry of algorithms with orthogonality constraints," *SIAM journal on Matrix Analysis and Applications*, vol. 20, no. 2, pp. 303–353, 1998.

[19] Anuj Srivastava and Eric Klassen, "Bayesian and geometric subspace tracking," *Advances in Applied Probability*, vol. 36, no. 1, pp. 43–56, 2004.

[20] Théodore Papadopoulo and Manolis IA Lourakis, "Estimating the jacobian of the singular value decomposition: Theory and applications," in *European Conference on Computer Vision*. Springer, 2000, pp. 554–570.

[21] Davis E. King, "Dlib-ml: A machine learning toolkit," *Journal of Machine Learning Research*, vol. 10, pp. 1755–1758, 2009.