# Real-time Memory Efficient Large-pose Face Alignment via Deep Evolutionary Network

Bin Sun, Ming Shao, *Member, IEEE,* Siyu Xia, *Member, IEEE,* and Yun Fu, *Fellow, IEEE*

*Abstract*—There is an urgent need to apply face alignment in a memory-efficient and real-time manner due to the recent explosion of face recognition applications. However, impact factors such as large pose variation and computational inefficiency, still hinder its broad implementation. To this end, we propose a computationally efficient deep evolutionary model integrated with 3D Diffusion Heap Maps (DHM). First, we introduce a sparse 3D DHM to assist the initial modeling process under extreme pose conditions. Afterward, a simple and effective CNN feature is extracted and fed to Recurrent Neural Network (RNN) for evolutionary learning. To accelerate the model, we propose an efficient network structure to accelerate the evolutionary learning process through a factorization strategy. Besides, we also propose a fast recurrent module to replace the traditional RNN for real-time regression. Extensive experiments on three popular alignment databases demonstrate the advantage of the proposed models over the state-of-the-art, especially under large-pose conditions. Notably, the computational speed of our model is 6 times faster than the state-of-the-art on CPU and 14 times on GPU. We also discuss and analyze the limitations of our models and future research work.

## I. INTRODUCTION

Face recognition and related application become increasingly popular, especially with the advances of deep learning. To name a few, face identification/verification [2], gaze detection [3], virtual face make-up [4], age synthesis [5], etc. Almost all of them heavily rely on face alignment that automatically locates predefined key points on a face. It has been treated as one of the fundamental problems in real-world face recognition systems. Nonetheless, there are still many barriers to face alignment techniques, involving unconstrained facial poses, complex expressions, and variable lighting conditions [6]. Among these factors, large facial pose is of great impact [7] and possible reasons are: (1) Face alignment relies on discriminant features, and these features become unreliable when they are from hidden landmarks [1]; (2) Conventional 2D models become vulnerable in modeling invisible landmarks due to its less flexibility; (3) Labeled 2D data is not reliable because invisible landmarks may be roughly estimated only according to the principle of symmetry.

In recent years, deep learning networks have shown great ability to overcome this problem by utilizing 3D information [8], [9], [1]. Our previous work [1] also shows promising performance on faces of large pose variations. However, the

Bin Sun and Yun Fu are with the Department of Electrical and Computer Engineering, Northeastern University, Boston, MA, USA. Ming Shao is with the Institute of CIS, College of Engineering, University of Massachusetts Dartmouth, North Dartmouth, MA, USA. Siyu Xia is with School of Automation, Southeast University, Nanjing, China. (e-mail:sun.bi@husky.neu.edu; mshao@umassd.edu; xia081@gmail.com; yunfu@ece.neu.edu )

time complexity and space complexity of deep convolution methods often beyond the capabilities of many mobile and embedded devices [10], [11], [12]. Therefore, reducing the computational cost and storage size of deep networks becomes an essential and challenging task for further application.

In this work, we mainly focus on two challenges in recent face alignment research: *invisible features* and *high computational cost*. To summarize, the contributions of this paper are:

- We propose a simple yet robust alignment feature learning paradigm using 3D Diffusion Heap Maps (DHM) and CNN to create high-level, reliable features containing both 2D and 3D information. Note that our DHM is calculated from the 3D model and only has three channels while 3DFAN [9] has 68 channels. This reduces the computation cost significantly.
- We cast face alignment to a deep evolutionary model with both 2D texture and 3D structure. Specifically, we use RNN to model the dynamics of the least square system. The system overview can be found in Figure 1.
- To achieve real-time performance, the output of each iteration in RNN module [1] is changed from parameters for 3D faces to normalized coordinates of landmarks. We further investigate factorized convolution structure [10] to accelerate our feature extraction process and keep the performance from a recession.
- We propose a fast recurrent module for our evolutionary learning paradigm as the light-weight replacement of traditional RNN module. Comparing with previous work [1], the fast recurrent module achieves even better result using much less computation cost and parameters.
- We conduct extensive experiments and improve the performance on a few benchmarks. We outperform the state-of-the-art by a large margin and show the robustness on the original AFLW2000-3D dataset [13] and LS3D-W [9]. Besides, a detailed comparison of parameters and speed is demonstrated in this work to show the efficiency of our model.
- Although this paper is an extension of our previous work published in [1], there exist some significant differences between the two works. Compared with the previous one, this framework is more efficient by utilizing the depthwise separable CNN structure and a fast RNN structure which are firstly proposed in this paper. To further reduce the computation cost, the output of this framework is landmarks whereas the output of previous work [1] is the parameters for the BFM model [14].
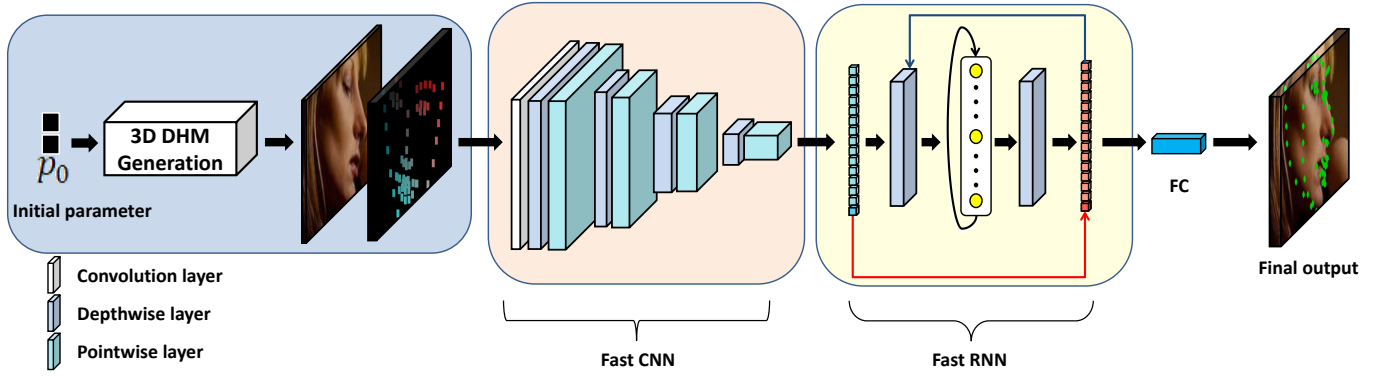
Fig. 1. Framework overview. Different from our previous work [1], this framework only uses RNN as its evolutionary regression, and its output is the normalized coordinates of landmarks. Through this improvement, the new framework in the paper could significantly accelerate our previous work [1] while reducing model size. A **fast CNN** module constructed by pointwise convolution and depthwise convolution is utilized in the framework. We also propose a **fast RNN** structure to further reduce the time complexity and space complexity of the whole framework.

## II. RELATED WORK

**2D Face Alignment:** The first milestone work of 2D face alignment is ASM [15], followed by many successful non-deep algorithms [16], [17] that considered the local patches around the facial landmarks as the features. Recently, critical works include tree-based models [18], [19] which improved the speed of face alignment to more than 1000 frames per second. Xiong et al. demonstrated the Supervised Descent Method (SDM) [20] with the cascade of weak regressors for face alignment, and achieved the state-of-the-art performance. Zhu et al. extended the work [21] and presented a new strategy [22] for large poses alignment by searching the best initial shape. Sun et al. [23] firstly employed CNN model for face alignment tasks with a raw face as the input and conduct regression with high-level features. Another extension of SDM called Global Supervised Descent Methods (GSDM) [24] tried to solve the large poses problem by dividing the training space into different descent spaces. All these face alignment methods only use 2D information and most of them use cascade method [18], [20], [21] and local patch features [20], [17], [16], [25]. SAN [26] implemented GAN model to generate training sample for further improvement of the accuracy. DSRN [27] utilized doubly CNN and Fourier embedding to learn the low-rank features of the facial key points. Differently, we suggest an integration of global 2D and 3D deep evolutionary network to overcome the information loss caused by 2D patch features.

**3D Face Alignment:** As 3D face model can maintain the depth information well against pose issues, a bunch of 3D face alignment methods and 3D face datasets become increasingly popular. Dollár et al. [28] estimated the landmarks on large poses through a 3D Morphable model (3DMM) with cascade regression in 2010. Zhu and Jourabloo et al. [13], [29] presented CNN based model fitting a 3DMM to a 2D face through a cascade method, along with the facial key points. A very dense 3D alignment model has been demonstrated by Liu et al. [30], [31] and achieved good performance. A 3D face training dataset called 300W-LP, and a testing dataset called AFLW-20003D were offered in [13] recently. Another 3D

alignment dataset called LS3D-W was published by Bulat et al. [9] with about 230,000 images, and the deep learning models based on this dataset, e.g., HourGlass(HG) [9], [32], [33], have achieved impressive performance very recently. In this paper, we use sparse 3D heat maps together with the original image as input, whereas most of the previous works use dense 3D models. Jointly working with a plain RNN and a simple CNN, our model achieves the state-of-the-art performance.

**Efficient Network:** In recent years, more efforts have been taken to speed up the deep learning models. Faster activation function named rectified-linear activation function (ReLU) was first proposed to accelerate the model [34]. In [35] depthwise separable convolution was initially introduced and was used in Inception models [36], Xception network [37], MobileNet [10], [11], and Shufflenet [12]. Jin et.al. [38] showed the flattened CNN structure was able to accelerate the feedforward procedure. A Factorized Network [39] was designed with a similar philosophy as well as the topological connection. A compression method of a deep neural network was introduced in [40], indicating that in certain cases complicate deep models are equal in performance by small models. Then Hinton et al. extended the work in [41] with the weight transfer strategy. Squeezenet [42] combined such work with a fire module containing many $1 \times 1$ convolutional layers. Another strategy [43], [44] of converting the parameters from float type to binary type can compress the model significantly and achieve an impressive speed. However, the performance will be compromised. In this paper, we use the factorized convolution structure in RNN modeling to reduce the number of parameters and accelerate the model.

## III. ALGORITHMS

In this section, we will detail our new framework (Figure 1) including three exclusive components: (1) 3D DHM generation; (2) deep evolutionary 3D heat maps; (3) fast deep evolutionary framework.
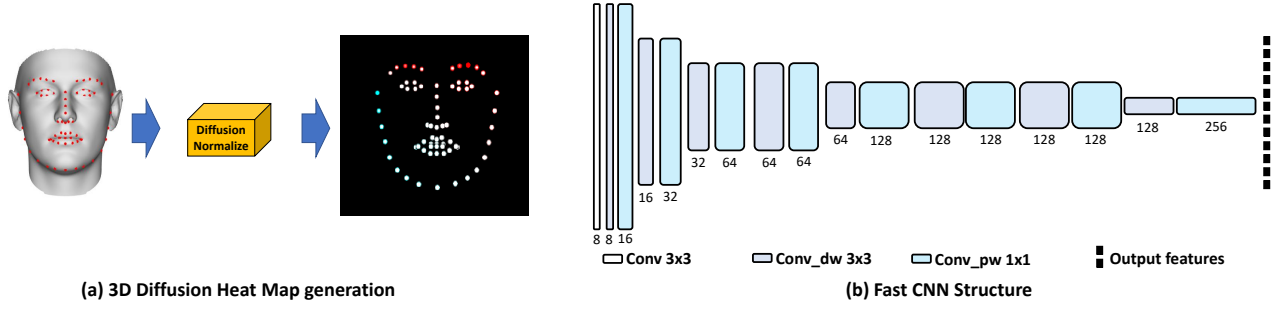
(a) 3D Diffusion Heat Map generation

(b) Fast CNN Structure

Fig. 2. Illustration of (a) generating 3D diffusion heat maps; (b) fast CNN module used as part of the new fast DHM framework.

### A. 3D Diffusion Heat Maps

To get feasible 3D landmarks in large poses, one of the reasonable ways is to build a 3D model of the face and simulate the details of a real face such as scale, expressions, and rotations, which can be formulated by the state-of-the-art 3DMM [14]. Typically, it represents factors of a 3D face by:

$$S = \overline{S} + E_{id}p_{id} + E_{exp}p_{exp}, \quad (1)$$

where $S$ is a predicted 3D face, $\overline{S}$ is the mean shape of the 3D face, $E_{id}$ is principle axes based on neutral 3D face, $p_{id}$ is shape parameters, $E_{exp}$ is principle axes based on the increment between expressional 3D face and neutral 3D face, and $p_{exp}$ is expression parameters. In our framework, the $E_{id}$ and $E_{exp}$ are calculated from a popular 3D face model named BFM [14]. Then we project the 3D face model by Weak Perspective Projection:

$$F(p) = f \times M \times R \times S + t_{2d}, \quad (2)$$

where $F(p)$ is the projected 3D face model, $f$ is the scale, $M$ is orthographic projected matrix $\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}$, $R$ is the rotation matrix written in $\begin{bmatrix} r_{pitch} & r_{yaw} & r_{roll} \end{bmatrix}$, $S$ is the 3D shape model calculated from Eq. (1), and $t_{2d}$ is the transition vector with the location coordinate $x$ and $y$. Since $F$ is the function of the parameter $p$, $p$ can be written as $p = \begin{bmatrix} f & R & t_{2d} & p_{id} & p_{exp} \end{bmatrix}$. We can generate the aligned 3D shape through Eq. (1) and Eq. (2). Afterwards, with the key point index provided by BFM, we will have precise locations of key points on the 3D model. To generate sparse 3D features, we normalize the coordinates in the 3D model around the key points. For a specific color channel $i$, the process could be described as:

$$\text{map}_i(k) = \frac{S_j(k)}{\max(S_j) - \min(S_j)}, \quad j \in \{x, y, z\}, \quad (3)$$

where "$\text{map}_i()$" is the 3D heat map with three channels R, G, and B, and the 3D coordinates triplet $\{x, y, z\}$ are mapped to the three channels. $S_j(k)$ means the value of the 3D shape at the location of the $k$-th landmarks. To incorporate the locality and increase the robustness of each landmark, we suggest to extend the 3D heat maps by a Gaussian diffusion map centered at each landmark's location, and thus, we obtain 3D DHM for robust representation. Specifically, we generate a set of heat maps centered at the landmarks and ranged by a 2D Gaussian

with standard deviation of 1 pixel (See Figure 2(a)). This strategy assures our framework can extract the features of the whole image with different weights instead of discarding the features located at a non-landmarks position. Note we conduct the normalization and diffusion on each single channel/axis, independently. A mean initialization DHM is generated before the training process.

**Discussions:** Recall the basic inputs of existing models usually include two separate parts: (1) an image; (2) initial mean landmarks. These methods usually depend on the initial landmarks for facial features for better performance, e.g., the features are usually extracted by cropping the image centered at initial landmarks. As discussed earlier, this strategy may degrade the performance in large pose situations, as the features centered at initial landmarks have significantly deviated from the ground truth. See the "Input" in Figure 1 (with face image and heap map). Thus, these features will misguide the learning model or regressor and probably converge to local trap. In contrast, we design a novel paradigm to address this issue. Namely, we keep the whole image as the input for robust facial feature learning and propose to employ sparse 3D shape information as the compliment. This avoids the issues of misguidance by the incorrect initial landmarks that propagate to the local features.

In the work 3DDFA [13], a diffusion features named PNCC was proposed. PNCC also fused 3D information with 2D image as the input. Compared with the PNCC features, our 3DDHM emphasize the features around the possible location of facial key points based on the mean shape. With such attention, our framework can localize the landmarks accurately using limited computation resource.

### B. Deep Evolutionary Diffusion Heat Maps

We offer an end-to-end trainable deep structure for face alignment in this section given the 3D heat maps and stacked image $I$. Since we have already generated DHM using 3D landmarks, we concentrate on: (1) discriminative and robust representation of image plus 3D DHM; (2) 3D DHM evolution; (3) recurrent HourGlass framework.

First, to formulate both discriminative and robust alignment features, we propose to use a plain CNN that absorbs both 2D and 3D information in a stacked structure. While handcraft features or direct use of 3D information may work, our practice reflects that they are less competitive than the well developed
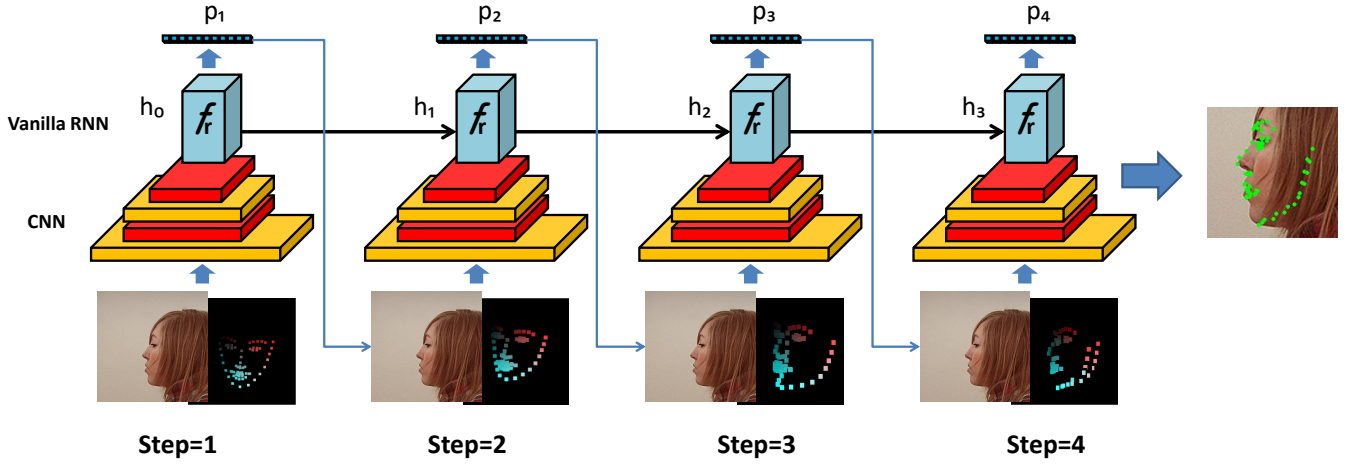
Fig. 3. Illustrations of the intermediate results of Deep Evolutionary Diffusion Heat Maps with state number $t = 4$ where changes of the heat maps can be seen. The heat map is scaled heavily in step 1 and deformed in step 2 to 4 with changed color. Finally, the heat map is well aligned on the face as shown in the output.

---

**Algorithm 1** Deep 3D Evolutionary Diffusion Heat Maps

**Inputs:** Image $I$, initial values: $p_0 = [p_{\exp}, p_{id}, f_0, R_0, t_{2d0}]$
**for** $i = 0$ to IterNum **do**
    generate $S_i$ in 3D
    $map = \text{zeros}()$
    **for** $j = 0$ to 3 **do**
        $S_i[j] -= \min(S_i[j])$
        $S_i[j] /= (\max(S_i[j]) - \min(S_i[j]))$
        $map[S_i[0], S_i[1], j] = S_i[j]$
    **end for**
    extract features $\phi^k(cat(I, map), C_{\text{conv2}})$
    $\hat{p}_k = RNN(\phi^k(cat(I, map)))$
**end for**
**Outputs:** $\hat{p}_k$

---

**CNN model.** The CNN model here can extract high-level features critical to alignment, and we are especially interested in global CNN features. In our experiments, we also find that an off-the-shelf CNN model such as VGG-net [45] works reasonably well in our case. We use reduced VGG network to fuse the 3D information and RGB information which proves useful in our experiments.

Second, we resort to evolutionary modeling for the alignment features. RNN has been widely applied to temporal data, as it is able to account for temporal dependencies. In training, RNN maintains the topology of feedforward networks while the feedback connections enable the representation of the current state of the system which encapsulates the information from the previous inputs, which can help update the parameter $p$ in the loops within the networks. Mathematically, the update rules can be written as:

$$h_{t+1} = \tanh(W_{\text{ih}} C_{\text{conv2}}([I, map]) + W_{\text{hh}} h_t), \quad (4)$$

where $h_t$ is the hidden state of the step $t$, $C_{\text{conv2}}([I, map])$ is the convolutional output features extracted from the input image $I$ and the heat maps "map" generated by Eq. 2. $W_{\text{ih}}$

is the weight from input to the hidden layer and the $W_{\text{hh}}$ is the weight from the hidden layer to the hidden layer. With the hidden features $h$, we can model the update rule for $\hat{p}$ using $\hat{p}_{t+1} = \hat{p}_t + W_{\text{ho}} h_t$, where $\hat{p}_t$ is the parameters in the step $t$, $W_{\text{ho}}$ is the weight from hidden layer to the output. Thus, from Eq. (4), we can prove that all the parameters in the step $t + 1$ are based on the state of the step $t$. Since the whole image has been engaged as the input in each step, we may rectify the errors caused by the previous steps. Besides, we have the generated heat maps to emphasize the change in the previous step so that the whole network can converge. An illustration of evolutions of 3D diffusion heat maps in four states can be found in Figure 3.

During the training, we define our loss function as

$$L = \min \|S_0 + \sum_{t=1}^{T} W_{\text{ho}} h_t(p) - S^*\|_F^2 \quad (5)$$

where $S^*$ is the ground truth shape, $F$ indicates the matrix Frobenius norm, $T$ is the total number of the steps. The complete algorithm is shown in Algorithm 1.

The optimized $p$ in the output layer will generate 3D landmarks for the test face which can be identified in the output of Figure 1. Here we suggest using Vanilla-RNN for simplicity. The evolutionary 3D DHM and intermediate results can be found in Figure 3. Note we update the parameters of the 3D model $p$ instead of the landmarks themselves as we would encourage the model to restrict the spatial relation of each landmark. The input is a $224 \times 224 \times 3$ image stacked by heat maps. The output of the RNN module is a 234-dimensional parameter, which will be cast to a 3D face model using Eq. (2). We can use specific vertices to find the landmark position. Last, to explore the generality of our framework, we upgrade the CNN by a popular stacked HourGlass module, and RNN by LSTM. In later experiments, we may compare with 3DFAN [9] that only uses the HourGlass module resulting in inferior performance. This also demonstrates our evolutionary learning strategy is more generic.

## C. Fast Deep Evolutionary Framework

An end-to-end trainable deep structure for face alignment given the 3D heat maps and stacked image $I$ is proposed before and its effectiveness has been proved [1]. However, such a framework is not efficient enough due to its iteration process. In this section, we concentrate on improving the efficiency and propose a novel framework named fast Deep Evolutionary Diffusion Heat Maps based on the Depthwise Separable Convolution module. Therefore, this subsection is divided into two parts: a brief introduction to Depthwise Separable Convolution [10], and the details of our fast Deep Evolutionary DHM framework.

*1) Depthwise Separable Convolution:* Depthwise Separable Convolution is one of the important factorization structures which plays key roles in building many lightweight architectures [12], [10], [11]. It consists of: (1) depthwise convolutional layer; (2) pointwise convolutional layer [10].

**Depthwise convolutional layer** is to apply a single convolutional filter to each input channel. This massively reduces the computational cost and the number of parameters. Particularly, the cost can be calculated as $S_F \times S_F \times S_k \times S_k \times C_{out}$, where $S_F$ is the size of the feature map, $S_k$ is the size of the kernel, $C_{out}$ is the number of the output features. The amount of the parameters without bias is computed by $S_k \times S_k \times C_{out}$, where $C_{out}$ is the channel number of the output.

**Pointwise convolutional layer** is to use $1 \times 1$ convolution to build the new features through computing the linear combinations of all input channels. Essentially, it is a traditional convolution layer with the kernel size being 1. Now, we show the computational cost of the traditional convolutional layer, which is calculated as $S_F \times S_F \times S_k \times S_k \times C_{in} \times C_{out}$. Since we assume there is no bias, the parameters for the traditional layer is computed by $S_k \times S_k \times C_{in} \times C_{out}$. The fact $S_k = 1$ in pointwise layer allows us to calculate its computational cost as $S_F \times S_F \times C_{in} \times C_{out}$, and the parameters as $C_{in} \times C_{out}$.

Because $C_{in}$ is usually much larger than $S_k^2$ in the depthwise layer, the depthwise convolution is more efficient than pointwise convolution in terms of computational cost and parameters.

*2) Fast DHM framework:* In brief, we accelerate both CNN and RNN models in our DHM framework. In the work of DHM [1], a plain CNN structure is used to extract both 2D and 3D information in a stacked structure. In our fast Deep Evolutionary framework, however, such CNN structure is replaced by the depthwise separable convolution block to reduce the computation complexity and parameters. The structure is shown in Figure 2 (b).

A novel fast Recurrent Network is also proposed to further accelerate the whole framework by utilizing the depthwise convolution and pointwise convolution in the RNN structure. The structure is shown in Figure 4. Inspired by the work [20], [25], we use a recurrent regression to learn the optimal of the non-linear problem $\min\|S - S^*\|_F^2$, in which $S$ is the predicted shape and $S^*$ is the ground truth. Different from the work [1], the output of each iteration is the feature maps for final landmarks instead of the parameter $p$ [1]. In this way, our framework can avoid the time-consuming computation of BFM in the recurrent process. However, the recurrent process

is still inefficient due to the convolution blocks in its iteration procedure. To improve efficiency, the factorization convolution block is implemented in the RNN module. Mathematically, the update rules can be written as:

$$F_{t+1} = F_t + \tanh(W_{ih}F_t + W_{hh}h_t), \qquad (6)$$

where $F_{t+1}$ is the output features from stage $t + 1$, $F_t$ is the output features from stage $t$, $W_{ih}$ is the weight from input to the hidden layer, $h_t$ is the hidden state of the stage $t$, and the $W_{hh}$ is the weight from hidden layer to hidden layer. In the initial stage, the input features are the feature maps extracted from CNN structure, and the output of the pointwise convolution layer is considered as the initial hidden state $h_0$. The output features can be considered as the increment of the landmarks in each stage. To keep the consistency during the evolutionary process, the input features in each stage (except initial stage) is the combination of previous features and the incremental features.

To avoid the over interference of previous features, a depthwise convolution layer is implemented on the input features of the current stage. Then a depthwise separable convolution structure is adopted to replace the traditional hidden layer with the input from the previous hidden state. Following the current hidden layers, a sum function and a $tanh$ activate function is applied to generate the final incremental features. The process of Eq. (6) can be written as:

$$F_{t+1} = F_t + \tanh(D_i \otimes F_t + D_h W_{hh} \otimes h_t), \qquad (7)$$

where $D_i$ is the weight of the depthwise convolution for the input, $D_h$ is the weight of the depthwise convolution for the hidden state, $W_{hh}$ is the weight of the $1 \times 1$ hidden layer, and $\otimes$ is the convolution operation.

During the training, the loss function is defined as $\min\|S - S^*\|_F^2$ where $S$ is the final predicted shapes after the fully connected layer as shown in Figure 1, $S^*$ is the ground truth shape, $F$ indicates the matrix Frobenius norm.

## IV. EXPERIMENTS

In this section, we will first detail the evaluation datasets for large pose face alignment and show the evaluation metrics. Then, baseline methods and parameter settings will be briefly introduced. Third, we will conduct an analysis of the proposed method and evaluate different modules. Last, we will compare with the state-of-the-art face alignment methods and analyze the results from accuracy, time complexity, and space complexity.

### A. Dataset and Evaluation Metrics

We use 68-point landmarks to conduct fair comparisons with the state-of-the-art methods, though our method can adapt to any numbers of landmarks. Note in the training process, we may need 3D landmarks or parameters which are not always available. Thus, we estimate the 3D information through [9] in this situation. Evaluation datasets are detailed below:

- 300W-LP: The dataset has four parts, a total of 61,225 samples across large poses (1,786 from IBUG, 5,207 from
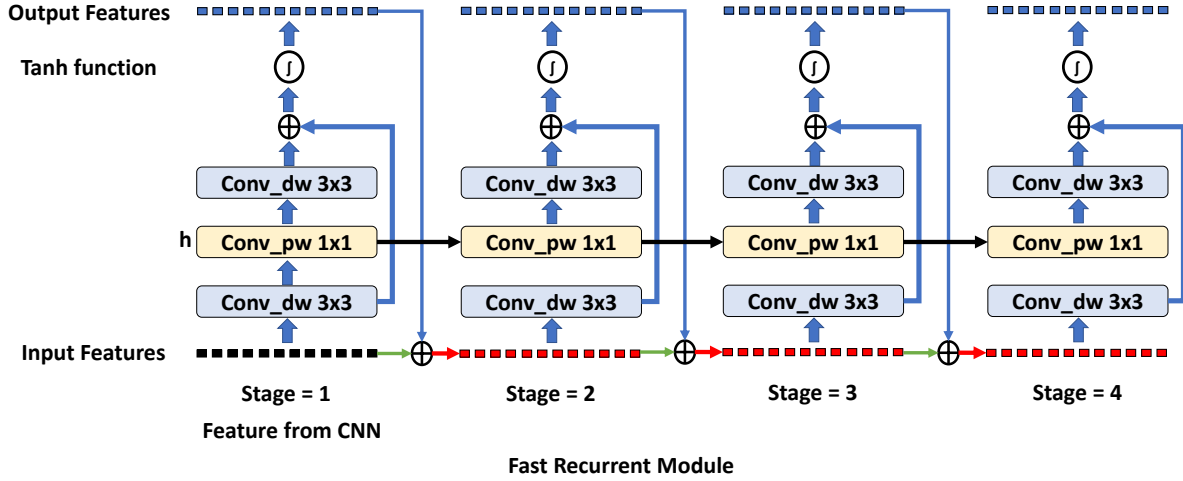
Fig. 4. The demonstration of our fast Recurrent module. Different from traditional RNN, we use one depthwise convolution layer to calculate the input weight and the output weight and use one pointwise convolution layer as the hidden layer. With our fast Recurrent module, the speed of our framework is improved from 17 FPS (with accelerated CNN module only) to 80 FPS. And as a result, the size of the model has been reduced to 976 KB.

AFW, 16,556 from LFPW and 37,676 from HELEN) [13]. Note we used 58,164 images for training and 3,061 as the validation.

- AFLW2000-3D: The dataset is essentially a reconstruction by Zhu et al. [13] given 2D landmarks. Note we use it for testing with 2000 images in total.
- Re-annotated AFLW2000-3D: The dataset is relabeled by Bulat et al. [9] from AFLW2000-3D given 2D landmarks. We use it for testing with 2000 images in total.
- LS3D-W: The dataset is also a re-annotated by Bulat et al. [9]. We use it for training and testing to make a fair comparison. We use 218,595 images for training, and use its sub-datasets Menpo-3D, 300W-3D, 300VW-3D (A), 300VW-3D (B), and 300VW-3D (C) for testing. Note this dataset only has 2D landmarks projected from 3D space.

As our focus is face alignment, we should reduce the negative effects of face detection. Thus, the detected bounding box of each face is computed by ground-truth landmarks. To compare with other methods, we use the same metric "Normalized Mean Error (NME)" defined as NME $= \frac{1}{N} \sum_{i=1}^{N} \frac{\left\| \hat{X}_i - X_i^* \right\|_2^2}{d}$ where the $\hat{X}$ and $X^*$ is predicted and ground truth landmarks, respectively, $N$ is the number of the landmarks, $d$ is normalized distance computed by the width and height of the bounding box using $d = \sqrt{w_{\text{bbox}} \times h_{\text{bbox}}}$. The lower NME means higher accuracy. We also show the curve of cumulative errors distribution (CED) and set the failure threshold as $6\%$. The CED curve indicates the percentage of successful cases in the test dataset. The speed of all methods is evaluated on Intel-i7 CPU with one core only. The storage size in this paper is calculated from the compressed model generated from the source code.

*B. Baselines and Parameter Settings*

We conduct comprehensive evaluations with state-of-the-art methods. In this paper, a comparison is made with deep state-of-the-art methods PCD-CNN [46], 3D-FAN [9], Hyperface [47], 3DSTN [7], 3DDFA [13], and MDM [25]. Among these baselines, the results of 3DSTN and PCD-CNN are cited from their original papers. We also compare the results and speed on CPU with some traditional methods running on CPU, including SDM [20], ERT [18], and ESR [48]. The brief introduction of each method is:

- DAMDNet [49]: DAMDNet is the state-of-the-art methods for large pose face alignment proposed in 2019. It implemented dual attention into the densenet framework and achieved impressive accuracy and speed.
- PRNet [46]: PRNet is the methods for 3D face reconstruction and alignment proposed in 2018. With the shape constraint imposed by UVMap, it improved the performance significantly.
- Hyperface [47]: Hyperface is a framework for multitask, including face alignment. Its main idea is to fuse different features by concatenation operation and usefully connected layer to the decoder the features for different tasks.
- 3DSTN [7]: 3D Spatial Transformer Networks utilize both the true 3D model of the subject in the image and the properties of the camera used to capture the image to model how the face changes from one viewpoint to another.
- 3DFAN [9]: 3DFAN is another state-of-the-art algorithm for 3D face alignment. It uses the Hour Glass (HG) Network to do both 2D and 3D alignment. In their structure, it has four HGs, and all the bottleneck blocks in HGs were replaced with a hierarchical, parallel and multi-scale block. In this way, it has achieved an impressive performance.
- 3DDFA [13]: 3DDFA is one of the state-of-the-art algorithms for 3D face alignment, using cascaded CNN as its main structure. It learns a set of CNN structure with the dense 3D features. The performance is further improved
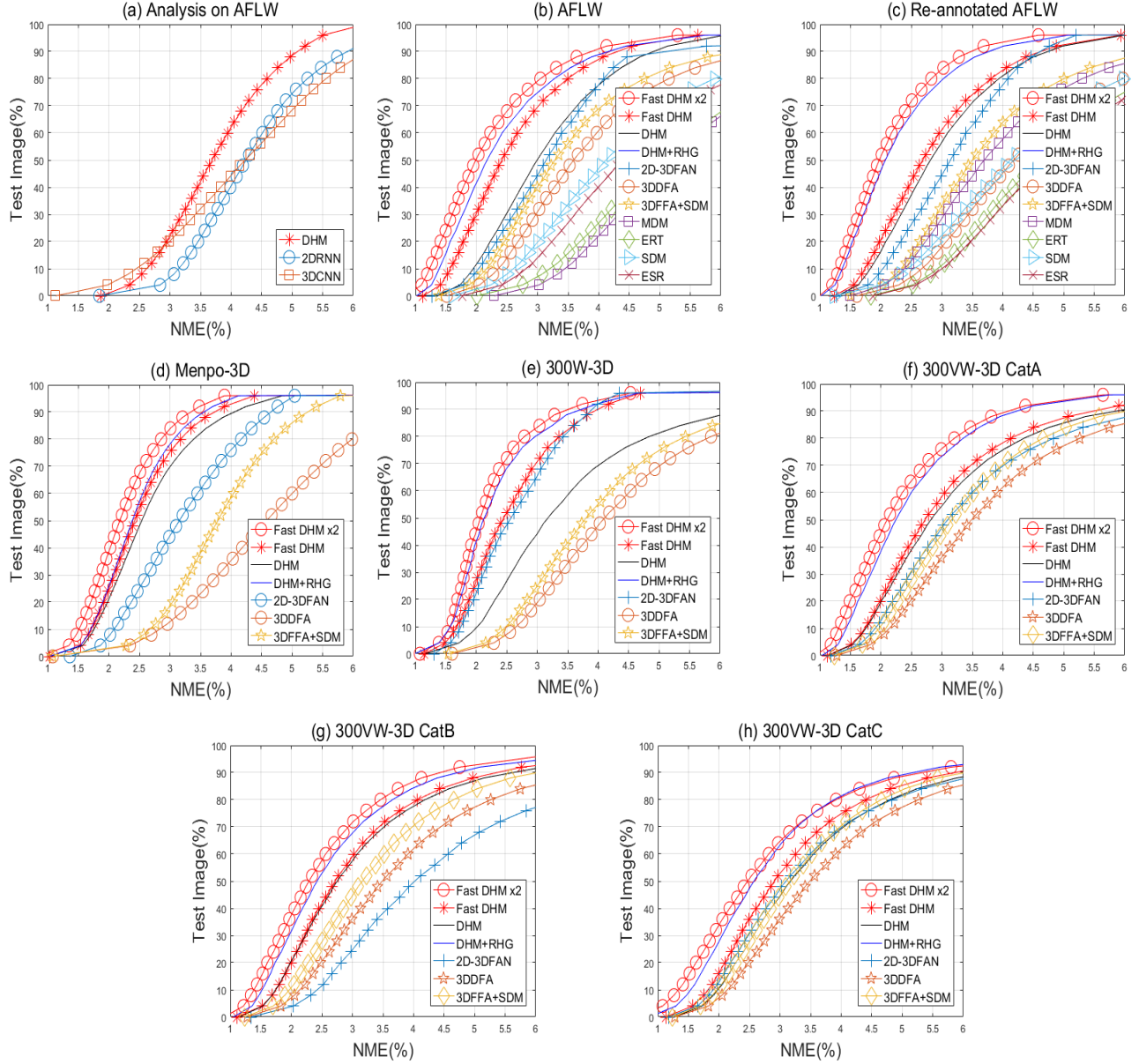
Fig. 5. (a) shows analysis of replacing 3D diffusion heat map and RNN model. (b) to (h) show the comparisons results on different test datasets. The CED area of our fast DHM ×2 is the largest among all baselines in different test datasets. The DHM with Recurrent Hourglass [1] has similar results with the Fast DHM×2, but it is much slower and larger.

recently [50].

- MDM [25]: MDM is the state-of-the-art algorithm for 2D alignment. The MDM regresses the patch features with the CNN structure and uses RNN to improve the cascade learning process.
- ERT [18]: ERT is a popular algorithm using regression trees using cascade to learn a set of weak regressors. It uses the pixel-wise feature and has achieved outstanding performance on 2D alignment.
- SDM [20]: SDM is a method calculating the descent in feature space to minimize the error. It also uses cascade to learn a set of linear regressors using the features and the offset of the landmarks.

To learn the weights of the network, we use Adam stochastic optimization [52] with default hyperparameters. The initial learning rate is 0.001 for 300W-LP with exponential decay of 0.95 every 2000 iterations, and an initial learning rate of 0.0001 is employed in our training process. The batch size is set to 50. The Recurrent HourGlass (HG) network starts with a $7\times7$ convolutional layer with stride 2. A residual module and a round of max pooling are added after it to bring the resolution down from 256 to 64. We use 3 stacked HG modules to extract features and an LSTM [53] as our recurrent module. The initial learning rate is 0.001, and we set weight decay at epoch 5, 15, 30. The total number of epochs is 40. We use RMSprop [54] as our optimizer. The training batch is 32, and validation
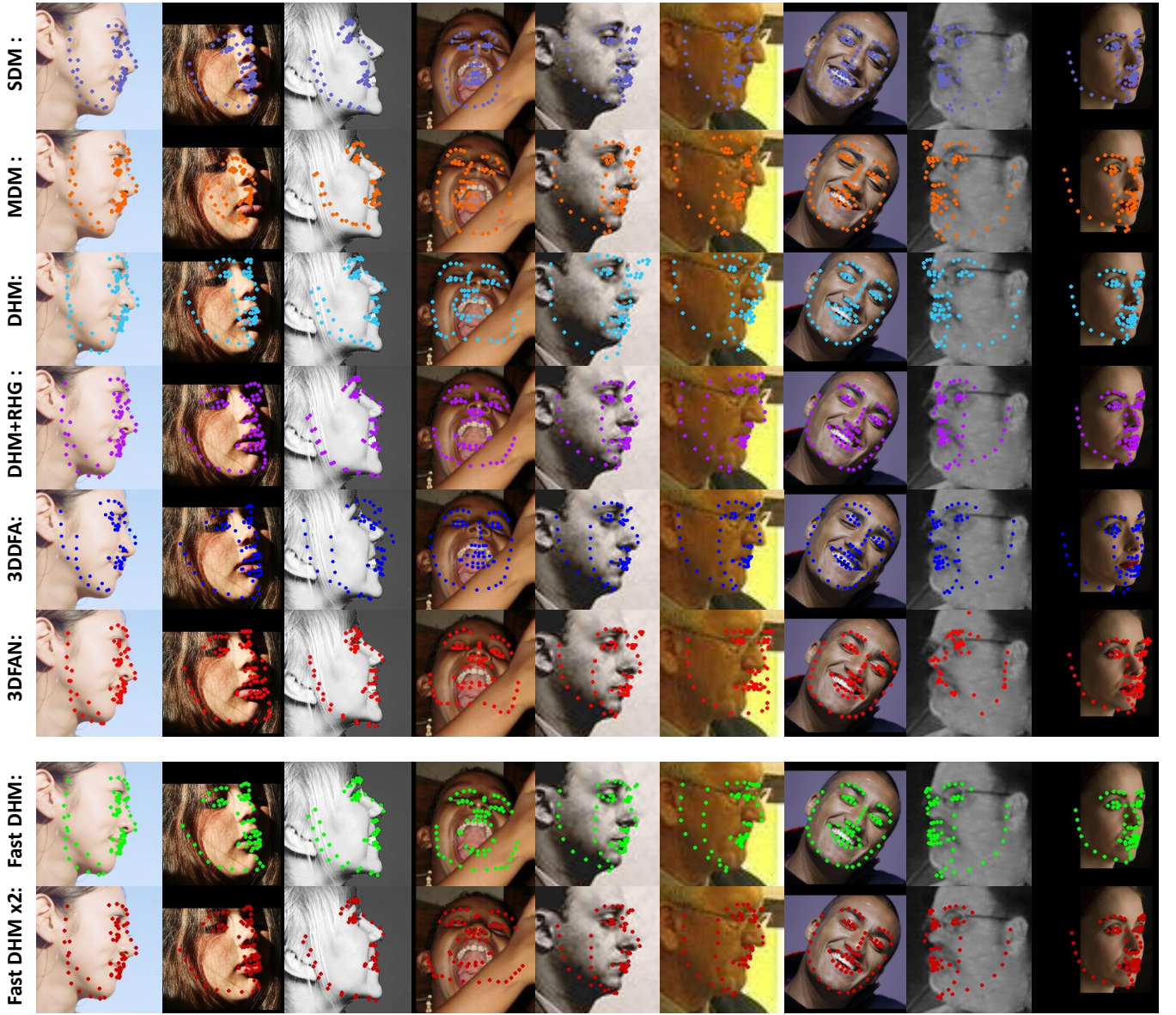
Fig. 6. Illustrations of alignment results by different methods. The bottom two rows of the figure show the results of Fast DHM×2 and Fast DHM.

batch is 16.

Fast DHM framework is trained with 0.005 as its initial learning rate. The optimization method is also Adam stochastic optimization. The training batch size is set to 50. Training epoch is set to 120. After 60 epochs, the learning rate is reset to 0.00001.

### C. Performance Analysis of Our Model

In this section, we will demonstrate the advantage and necessity of two modules: (1) 3D diffusion heat maps; (2) RNN for deep evolution.

First, we evaluate the importance of evolutionary 3D maps. Instead of using RNN, we use a plain 3D-CNN structure. That is being said, we use the same 3D heat maps but replace the RNN by a conventional CNN structure. The rest parts remain the same. Note we use 300W-LP dataset for training and

AFLW2000-3D dataset for evaluation. Results of this setting (3D-CNN) can be found in Figure 5(a).

Second, we keep the RNN structure and test the importance of 3D heat maps in our framework. We replace the 3D module by initial 2D landmarks. Accordingly, we change the output of the framework from a $234 \times 1$ vector to a $204 \times 1$ vector, which is the increment of locations of predicted landmarks. The rest of the framework remains the same. Note we also use the same training and testing datasets. The result of this setting (2D-RNN) can be found in Figure 5(a). It can be found that the combination of both modules performs best.

### D. Comparisons with Existing Methods

In this section, we conduct comprehensive evaluations with state-of-the-art methods. Especially, all methods are trained on the 300W-LP dataset including both ours and others. All of the

TABLE I
COMPARISONS WITH STATE-OF-THE-ART METHODS. WE HIGHLIGHT THE TOP THREE RESULTS IN EACH SETTING. FAST DHM×2 MEANS THE CHANNELS IN THE FAST CNN MODULE ARE DOUBLED.

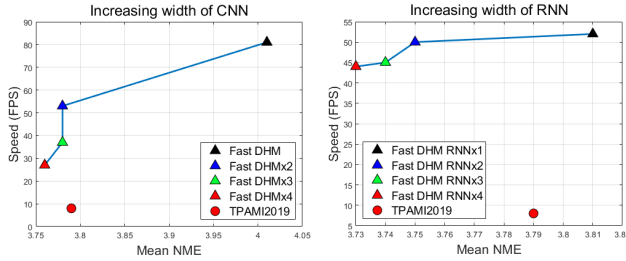| | Normalized Mean Error on AFLW2000-3D | | | | Speed (FPS) | | Memory |
| Method Name | $[0°30°]$ | $[30°60°]$ | $[60°90°]$ | Mean | GPU | CPU | params (Bytes) |
| --- | --- | --- | --- | --- | --- | --- | --- |
| **Fast DHM×2** | **2.43** | **3.51** | 5.41 | **3.78** | **740** | 53 | **2.2M** |
| **Fast DHM** | **2.57** | 3.60 | 5.88 | 4.01 | >900 | **81** | **0.9M** |
| DAMDNet (2019ICCVW) [49] | 2.91 | 3.83 | **4.95** | 3.89 | 50 | 12.5 | 11.3M |
| Improved 3DDFA (2019TPAMI) [50] | 2.84 | 3.57 | **4.96** | **3.79** | 45 | 8 | 12.6M |
| DHM+RHG (2018BMVC) [1] | **2.52** | **3.21** | 5.48 | 3.85 | 15 | < 1 | 221.9M |
| DHM (2018BMVC) [1] | 2.75 | 4.21 | 6.91 | 4.11 | 27 | < 1 | 193.7M |
| PRNet (2018ECCV) [51] | 2.75 | **3.51** | **4.61** | **3.62** | 100 | 5 | 153M |
| Hyperface (2017TPAMI) [47] | 3.93 | 4.14 | 6.71 | 4.26 | - | - | 119.7M |
| 3DSTN (2017ICCV) [7] | 3.15 | 4.33 | 5.98 | 4.49 | 52 | < 1 | - |
| 3DFAN (2017ICCV) [9] | 2.75 | 3.76 | 5.72 | 4.07 | 6 | < 1 | 183M |
| 3DDFA (2016CVPR) [13] | 3.78 | 4.54 | 7.93 | 5.42 | 15 | 8 | 111M |
| 3DDFA+SDM (2016CVPR) [13] | 3.43 | 4.24 | 7.17 | 4.94 | 10 | 7 | 121M |
| MDM (2016CVPR) [25] | 3.67 | 5.94 | 10.76 | 6.45 | 5 | < 1 | 307M |
| ERT (2014CVPR) [18] | 5.40 | 7.12 | 16.01 | 10.55 | - | **300** | 95M |
| ESR (2014IJCV) [48] | 4.60 | 6.70 | 12.67 | 7.99 | - | **83** | 248M |
| SDM (2013CVPR) [20] | 3.67 | 4.94 | 9.76 | 6.12 | - | 80 | **10M** |



Fig. 7. Relationships of speed and mean NME among different sizes of fast CNN module (left) and fast RNN module (right) in our fast DHM framework. $\times n$ means the channels of fast CNN module is extended by $n$. We also show the result of TPAMI 2019 [50] for clear comparison of the baseline. Note that the experiments of fast RNN module is implemented based on Fast DHM CNN ×2.
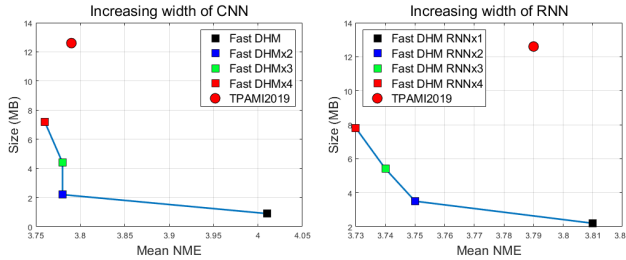


Fig. 8. Relationships of storage size and mean NME among different sizes of fast CNN module (left) and fast RNN module (right) in our fast DHM framework. $\times n$ means the channels of fast CNN module is extended by $n$. We also show the result of TPAMI 2019 [50] for clear comparison of the baseline. Note that the experiments of fast RNN module is implemented based on Fast DHM CNN ×2.

input faces are cropped by the bounding box calculated from landmarks. The methods with released codes are trained on the 300W-LP and the comparison CED curves can be found in Figure 5 (b) to (h). The quantitative results can be found in Table I.

**Accuracy:** In the Table I, we highlight the top three methods in each columns. Our proposed fast deep evolutionary network

has the second lowest mean NME among all the methods. Although our NME in $[60°90°]$ is 0.8 higher than PRNet and is not in the top three, our performance ranks top three in $[0°30°]$, $[30°60°]$ and mean error evaluation, which means our fast deep evolutionary framework with DHM is competitive with the state-of-the-art methods on the accuracy. From the experiments, we can prove that our approach is very robust and reliable on different datasets. The visualization of eight methods are shown in Figure 6.

**Time Complexity:** Since face alignment algorithms are often used on mobile devices, real-time implementation of the algorithms without GPU support is important. Compared with those deep learning methods [49], [50], [47], [7], [13], [9], [41], [25], our fast deep evolutionary framework has much better speed on both one core CPU and GPU. Our speed is ×6 of DAMDNet on CPU and ×14 on GPU. Besides, our Fast DHM×2 is ×10 faster than PRNet on CPU and ×7 faster on GPU. The results are shown in Table I. There are two reasons for the impressive acceleration. First, the lightweight depthwise convolution is used in the CNN structure and the dimension of each layer is small. Second, the recurrent block is replaced by only one pointwise layer and two depthwise layers. In the table, we notice that the SDM [20], ERT [18] and ESR [48] have very impressive speed on CPU. The reason is all of these methods use hand-craft features which are computationally efficient. However, the accuracy based on these features is inferior and fail to provide great representation. We also explore the relationship between the speed and mean NME among different sizes of fast CNN module and RNN module relatively. The results are shown in Figure 7. From the figure, it is clear to see that there is a saturate point when the mean NME is around 3.8, which can be reached by just double the channels of fast CNN module. After the saturate point, the loss of the speed is far beyond of the benefit from the improvement of the accuracy. However, the trade-off between the speed and the increment of the accuracy brought by RNN module seems acceptable.

**Space Complexity:** For the applications on mobile devices, the face alignment model may use limited memory for
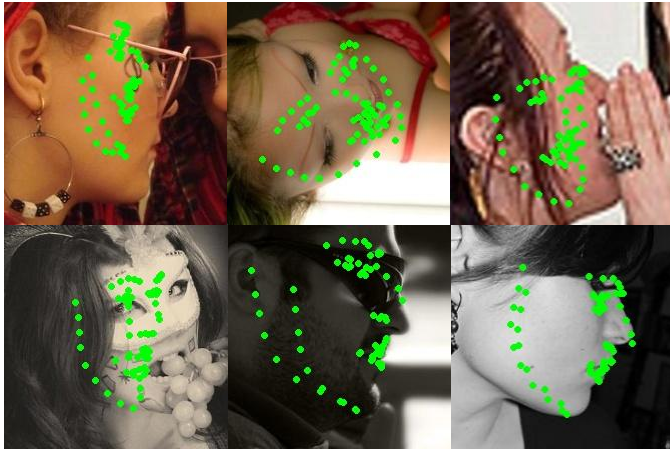
Fig. 9. Illustration of some failure cases by our models. The errors mainly occur when the large-pose faces are occluded.

functionality. Thus, it is important to measure the memory consumption of each model. From Table I, it can be observed that the DHM with the fast evolutionary framework is $\times 12.5$ smaller than the smallest model in baseline deep learning methods. Besides, it is $\times 10$ smaller than the smallest model in all baselines. Note that, since the DHM is pre-defined before the testing stage, we do not take the size of BFM model into consideration. We also show the relationship between the size and mean NME among different sizes of fast CNN module in Figure 8. From the figure, we can observe that our Fast DHM reaches the saturate point at the Fast DHM $\times 2$, which means the channels of the fast CNN module are doubled. For RNN module, the reduction of mean NME is only $0.03$ while the storage size increases almost $50\%$ every time the RNN module is increased, which is not cost-effective. Thus, the Fast DHM $\times 2$ is utilized to achieve the best performance on mean NME, speed, and model size.

We also illustrate some failure cases of our model in Figure 9 which indicate that our model may be fragile given large-angle rolled faces with heavy occlusions. The primary reason is the lack of faces in similar cases in the training dataset, which is a common issue for all other methods. Possible solutions include adding corresponding training data to approach the special case, employing multiple initializations to avoid bad local minimal, and utilizing advanced loss to fit the data.
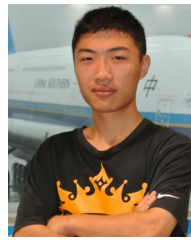
## V. CONCLUSIONS

In this paper, we focused on improving the face alignment algorithms with the sparse 3D landmarks to approach the challenge of large poses. We presented a deep evolutionary framework to progressively update the 3D heat maps to generated target face landmarks. First, we proposed to use as a robust representation. Second, we demonstrated that an RNN based evolutionary learning paradigm was able to model the dynamics of least square problems and optimize the facial landmarks. Third, we proposed a fast framework for deep evolutionary learning and a lightweight recurrent block. The results show the efficiency of our framework on large pose

face alignment task, and also show the possibility of further improvement with complicated structures in our evolutionary DHM strategy. At last, we believe the fast recurrent block can be further implemented on other temporal tasks.

## REFERENCES

[1] B. Sun, M. Shao, S.-Y. Xia, and Y. Fu, "Deep evolutionary 3d diffusion heat maps for large-pose face alignment," in *BMVC*, 2018.
[2] W. Zhao, R. Chellappa, P. J. Phillips, and A. Rosenfeld, "Face recognition: A literature survey," *ACM computing surveys (CSUR)*, vol. 35, no. 4, pp. 399–458, 2003.
[3] D. W. Hansen and Q. Ji, "In the eye of the beholder: A survey of models for eyes and gaze," *TPAMI*, vol. 32, no. 3, pp. 478–500, 2010.
[4] D. Guo and T. Sim, "Digital face makeup by example," in *CVPR*. IEEE, 2009, pp. 73–79.
[5] Y. Fu, G. Guo, and T. S. Huang, "Age synthesis and estimation via faces: A survey," *TPAMI*, vol. 32, no. 11, pp. 1955–1976, 2010.
[6] S. Bazrafkan, H. Javidnia, and P. Corcoran, "Face synthesis with landmark points from generative adversarial networks and inverse latent space mapping," *arXiv preprint arXiv:1802.00390*, 2018.
[7] C. Bhagavatula, C. Zhu, K. Luu, and M. Savvides, "Faster than real-time facial alignment: A 3d spatial transformer network approach in unconstrained poses," *Proc. ICCV, page to appear*, vol. 2, 2017.
[8] F. Liu, D. Zeng, Q. Zhao, and X. Liu, "Joint face alignment and 3d face reconstruction," in *ECCV*. Springer, 2016, pp. 545–560.
[9] A. Bulat and G. Tzimiropoulos, "How far are we from solving the 2d & 3d face alignment problem?(and a dataset of 230,000 3d facial landmarks)," *arXiv preprint arXiv:1703.07332*, 2017.
[10] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
[11] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Inverted residuals and linear bottlenecks: Mobile networks for classification, detection and segmentation," *arXiv preprint arXiv:1801.04381*, 2018.
[12] X. Zhang, X. Zhou, M. Lin, and J. Sun, "Shufflenet: An extremely efficient convolutional neural network for mobile devices," *arXiv preprint arXiv:1707.01083*, 2017.
[13] X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li, "Face alignment across large poses: A 3d solution," in *Proc CVPR*, 2016, pp. 146–155.
[14] V. Blanz and T. Vetter, "Face recognition based on fitting a 3d morphable model," *TPAMI*, vol. 25, no. 9, pp. 1063–1074, 2003.
[15] T. F. Cootes, C. J. Taylor, D. H. Cooper, and J. Graham, "Active shape models-their training and application," *CVIU*, vol. 61, no. 1, pp. 38–59, 1995.
[16] T. F. Cootes, G. J. Edwards, and C. J. Taylor, "Active appearance models," *TPAMI*, vol. 23, no. 6, pp. 681–685, 2001.
[17] T. Baltrušaitis, P. Robinson, and L.-P. Morency, "3d constrained local model for rigid and non-rigid facial tracking," in *CVPR*. IEEE, 2012, pp. 2610–2617.
[18] V. Kazemi and J. Sullivan, "One millisecond face alignment with an ensemble of regression trees," in *Proc CVPR*, 2014, pp. 1867–1874.
[19] S. Ren, X. Cao, Y. Wei, and J. Sun, "Face alignment at 3000 fps via regressing local binary features," in *Proc CVPR*, 2014, pp. 1685–1692.
[20] X. Xiong and F. De la Torre, "Supervised descent method and its applications to face alignment," in *Proc CVPR*, 2013, pp. 532–539.
[21] J. Zhang, S. Shan, M. Kan, and X. Chen, "Coarse-to-fine auto-encoder networks (cfan) for real-time face alignment," in *ECCV*. Springer, 2014, pp. 1–16.
[22] S. Zhu, C. Li, C. Change Loy, and X. Tang, "Face alignment by coarse-to-fine shape searching," in *Proc CVPR*, 2015, pp. 4998–5006.
[23] Y. Sun, X. Wang, and X. Tang, "Deep convolutional network cascade for facial point detection," in *Proc CVPR*, 2013, pp. 3476–3483.
[24] X. Xiong and F. De la Torre, "Global supervised descent method," in *Proc CVPR*, 2015, pp. 2664–2673.
[25] G. Trigeorgis, P. Snape, M. A. Nicolaou, E. Antonakos, and S. Zafeiriou, "Mnemonic descent method: A recurrent process applied for end-to-end face alignment," in *Proc CVPR*, 2016, pp. 4177–4187.
[26] X. Dong, Y. Yan, W. Ouyang, and Y. Yang, "Style aggregated network for facial landmark detection," in *CVPR*, 2018, pp. 379–388.
[27] X. Miao, X. Zhen, X. Liu, C. Deng, V. Athitsos, and H. Huang, "Direct shape regression networks for end-to-end face alignment," in *CVPR*, 2018, pp. 5040–5049.

[28] P. Dollár, P. Welinder, and P. Perona, "Cascaded pose regression," in *CVPR*. IEEE, 2010, pp. 1078–1085.

[29] A. Jourabloo and X. Liu, "Large-pose face alignment via cnn-based dense 3d model fitting," in *Proc CVPR*, 2016, pp. 4188–4196.

[30] Y. Liu, A. Jourabloo, W. Ren, and X. Liu, "Dense face alignment," *arXiv preprint arXiv:1709.01442*, 2017.

[31] A. Jourabloo and X. Liu, "Pose-invariant face alignment via cnn-based dense 3d model fitting," *IJCV*, vol. 124, no. 2, pp. 187–203, 2017.

[32] J. Yang, Q. Liu, and K. Zhang, "Stacked hourglass network for robust facial landmark localisation," in *CVPRW*. IEEE, 2017, pp. 2025–2033.

[33] A. Bulat and G. Tzimiropoulos, "Binarized convolutional landmark localizers for human pose estimation and face alignment with limited resources," in *ICCV*, 2017.

[34] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, 2011, pp. 315–323.

[35] L. Sifre and P. Mallat, "Rigid-motion scattering for image classification," Ph.D. dissertation, Citeseer, 2014.

[36] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.

[37] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," *arXiv preprint*, 2016.

[38] J. Jin, A. Dundar, and E. Culurciello, "Flattened convolutional neural networks for feedforward acceleration," *arXiv preprint arXiv:1412.5474*, 2014.

[39] ——, "Flattened convolutional neural networks for feedforward acceleration," *CoRR*, vol. abs/1412.5474, 2014. [Online]. Available: http://arxiv.org/abs/1412.5474

[40] J. Ba and R. Caruana, "Do deep nets really need to be deep?" in *Advances in neural information processing systems*, 2014, pp. 2654–2662.

[41] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.

[42] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and¡ 0.5 mb model size," *arXiv preprint arXiv:1602.07360*, 2016.

[43] M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1," *arXiv preprint arXiv:1602.02830*, 2016.

[44] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "Xnor-net: Imagenet classification using binary convolutional neural networks," in *ECCV*. Springer, 2016, pp. 525–542.

[45] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[46] A. Kumar and R. Chellappa, "Disentangling 3d pose in a dendritic cnn for unconstrained 2d face alignment," *arXiv preprint arXiv:1802.06713*, 2018.

[47] R. Ranjan, V. M. Patel, and R. Chellappa, "Hyperface: A deep multi-task learning framework for face detection, landmark localization, pose estimation, and gender recognition," *IEEE TPAMI*, 2017.

[48] X. Cao, Y. Wei, F. Wen, and J. Sun, "Face alignment by explicit shape regression," *IJCV*, vol. 107, no. 2, pp. 177–190, 2014.

[49] L. Jiang, X.-J. Wu, and J. Kittler, "Dual attention mobdensenet (damd-net) for robust 3d face alignment," in *ICCVW*, 2019, pp. 0–0.

[50] X. Zhu, X. Liu, Z. Lei, and S. Z. Li, "Face alignment in full pose range: A 3d total solution," *IEEE TPAMI*, vol. 41, no. 1, pp. 78–92, 2019.

[51] Y. Feng, F. Wu, X. Shao, Y. Wang, and X. Zhou, "Joint 3d face reconstruction and dense alignment with position map regression network," in *ECCV*, 2018, pp. 534–551.

[52] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[53] A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *ICASSP*. IEEE, 2013, pp. 6645–6649.

[54] T. Tieleman and G. Hinton, "Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude," *COURSERA: Neural networks for machine learning*, vol. 4, no. 2, pp. 26–31, 2012.

**Bin Sun** received the B.Eng. degree in photoelectric information engineering from Beijing Institute of Technology (BIT), China, in 2015. He is pursuing the Ph.D. degree with the Department of Electrical and Computer Engineering, Northeastern University, Boston, MA, USA. His research interests include Face Synthesis and light weight network design. He has served as PC member for AAAI, and published paper in BMVC 2018.
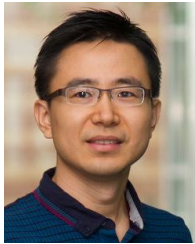
**Ming Shao** (S11-M16) received the B.E. degree in computer science, the B.S. degree in applied mathematics, and the M.E. degree in computer science from Beihang University, Beijing, China, in 2006, 2007, and 2010, respectively. He received the Ph.D. degree in computer engineering from Northeastern University, Boston MA, 2016. He is a tenure-track Assistant Professor affiliated with College of Engineering at the University of Massachusetts Dartmouth since 2016 Fall. His current research interests include sparse modeling, low-rank matrix analysis, deep learning, and applied machine learning on social media analytics. He was the recipient of the Presidential Fellowship of State University of New York at Buffalo from 2010 to 2012, and the best paper award winner/candidate of IEEE ICDM 2011 Workshop on Large Scale Visual Analytics, and ICME 2014. He has served as the reviewers for many IEEE Transactions journals including TPAMI, TKDE, TNNLS, TIP, and TMM. He has also served on the program committee for the conferences including AAAI, IJCAI, CVPR, ICCV, ICLR, etc. He is the Associate Editor of SPIE Journal of Electronic Imaging, and IEEE Computational Intelligence Magazine. He is a member of IEEE.

**Siyu Xia** received his BE and MS degrees in automation engineering from Nanjing University of Aeronautics and Astronautics, Nanjing, China, in 2000 and 2003, respectively, and the PhD degree in pattern recognition and intelligence system from Southeast University, Nanjing, China, in 2006. He is currently working as an associate professor in the School of Automation at Southeast University, Nanjing, China. His research interests include object detection, applied machine learning, social media analysis, and intelligent vision systems. He was the recipient of the Science Research Famous Achievement Award in Higher Institution of China in 2015. He has served as the reviewer of many journals including TIP, T-SMC-B, T-IFS, T-MM, IJPRAI, and Neurocomputing. He received Outstanding Reviewer Award for Journal of Neurocomputing in 2016. He has also served on the PC/SPC for the conferences including AAAI, ACM-MM, ICME, and ICMLA. He is a member of IEEE and ACM.

**Yun Fu** (S'07-M'08-SM'11-F'19) received the
B.Eng. degree in information engineering and the
M.Eng. degree in pattern recognition and intelli-
gence systems from Xian Jiaotong University, China,
respectively, and the M.S. degree in statistics and
the Ph.D. degree in electrical and computer engi-
neering from the University of Illinois at Urbana-
Champaign, respectively. He is an interdisciplinary
faculty member affiliated with College of Engineer-
ing and the College of Computer and Information
Science at Northeastern University since 2012. His
research interests are Machine Learning, Computational Intelligence, Big Data
Mining, Computer Vision, Pattern Recognition, and Cyber-Physical Systems.
He has extensive publications in leading journals, books/book chapters and
international conferences/workshops. He serves as associate editor, chairs,
PC member and reviewer of many top journals and international confer-
ences/workshops. He received seven Prestigious Young Investigator Awards
from NAE, ONR, ARO, IEEE, INNS, UIUC, Grainger Foundation; nine
Best Paper Awards from IEEE, IAPR, SPIE, SIAM; many major Industrial
Research Awards from Google, Samsung, and Adobe, etc. He is currently an
Associate Editor of the IEEE Transactions on Neural Networks and Leaning
Systems (TNNLS). He is fellow of IEEE, IAPR, OSA and SPIE, a Lifetime
Distinguished Member of ACM, Lifetime Member of AAAI and Institute
of Mathematical Statistics, member of ACM Future of Computing Academy,
Global Young Academy, AAAS, INNS and Beckman Graduate Fellow during
2007-2008.