# FAB: A Robust Facial Landmark Detection Framework for Motion-Blurred Videos

Keqiang Sun[1], Wenyan (Wayne) Wu[2], Tinghao Liu[3], Shuo Yang[4]
Quan Wang[3], Qiang Zhou[2], Zuochang Ye[1], Chen Qian[3]
[1]Institute of Microelectronics, Tsinghua University
[2]Beijing National Research Center for Information Science and Technology (BNRist), Tsinghua University
[3]SenseTime Research [4]Amazon Rekognition
{skq17,wwy15}@mails.tsinghua.edu.cn, {zuochang,zhouqiang}@tsinghua.edu.cn,
{liutinghao,wangquan,qianchen}@sensetime.com, shuoy@amazon.com

## Abstract

*Recently, facial landmark detection algorithms have achieved remarkable performance on static images. However, these algorithms are neither accurate nor stable in motion-blurred videos. The missing of structure information makes it difficult for state-of-the-art facial landmark detection algorithms to yield good results.*

*In this paper, we propose a framework named FAB that takes advantage of structure consistency in the temporal dimension for facial landmark detection in motion-blurred videos. A structure predictor is proposed to predict the missing face structural information temporally, which serves as a geometry prior. This allows our framework to work as a virtuous circle. On one hand, the geometry prior helps our structure-aware deblurring network generates high quality deblurred images which lead to better landmark detection results. On the other hand, better landmark detection results help structure predictor generate better geometry prior for the next frame. Moreover, it is a flexible video-based framework that can incorporate any static image-based methods to provide a performance boost on video datasets. Extensive experiments on Blurred-300VW, the proposed Real-world Motion Blur (RWMB) datasets and 300VW demonstrate the superior performance to the state-of-the-art methods. Datasets and models will be publicly available at https://keqiangsun.github.io/projects/FAB/FAB.html.*

## 1. Introduction

Facial landmark detection, or known as face alignment, serves as a key component for many face applications, *e.g.* face recognition, face verification and face augmented reality. Previous researches [41, 45, 46, 39, 8, 9, 38, 25] mainly
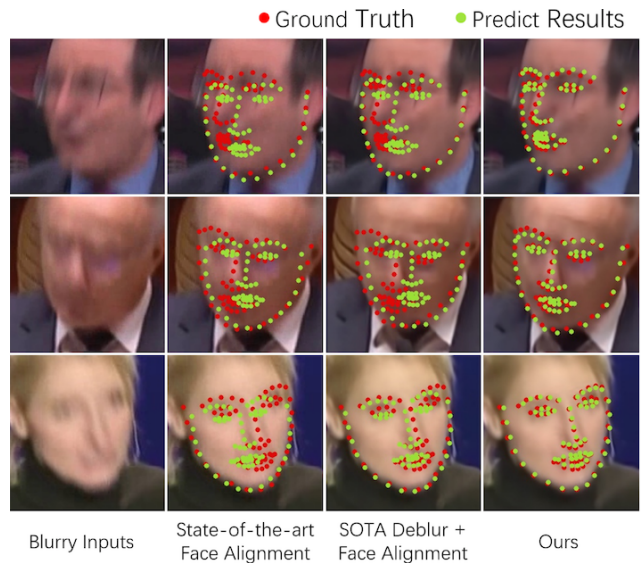


Figure 1: The first column is the frames of Blurred-300VW. In the second column, the results of state-of-the-art facial landmark detection algorithm are given. The third column corresponds to a naive combination of deblurring and facial landmark detection. In the fourth column are the results of our proposed algorithm.

focus on detecting facial landmarks in static images. One natural extension of image-based facial landmark detection is video facial landmark detection, which tries to locate facial landmarks in sequential frames. Different from static images, motion blur usually appears in videos, due to the mismatch of the motion speed and shutter closing speed. As shown in Figure 1, the missing of structure information, *i.e.* edges or boundaries, in the motion-blurred videos makes it difficult for the state-of-the-art facial landmark detection algorithms to capture facial structure. The objective of this paper is to devise an effective framework to handle facial landmark detection in motion-blurred videos.

An intuitive method is to employ a deblurring algorithm before facial landmark detection. It is straightforward that the deblurred image would promote facial landmark detection performance. State-of-the-art face deblurring algorithms [35, 6, 15] tend to rely on facial structure (*e.g.* facial landmarks and edges) from the input image as strong priors to restore the shape and details. Nevertheless, these face structure would be unavailable when the landmarks cannot be precisely predicted, *e.g.* in extremely blurry frames. In conclusion, the facial landmark detection in motion-blurred videos requires deblurred images, while face deblurring tends to rely on face structure like facial landmarks. They are mutually beneficial and interdependent, which makes it a "chicken or egg" dilemma in motion-blurred videos.

In a video clip, face structure keeps temporal continuity and consistency, and blurry frames usually intersperse along the time dimension. This motivates us to predict a reliable face structure according to previous structural information. Inspired by this idea, we proposed a Structure Predictor to predict current face structure. Specifically, given previous facial edges, the predictor figures out the optical flow and extends the motion to predict next face edges by assuming the optical flow is linear.

Based on this, we designed a framework, composed of three modules, *i.e.* structure predictor, structure-aware motion deblurring network, and replaceable facial landmark detection network, as is depicted in Figure 2. These three components work as an organic whole. The structure predictor predicts the current facial structure from previous temporal cues. The deblurring network removes the motion blur with the assistance of the structure prior. With the deblurred image, the landmark detector yields accurate landmarks, which are used to predict the next face structure.

Since our framework is proposed for Face Alignment in Blurred-videos, we call it "FAB". FAB is designed with the advantage of great flexibility. Each component of the framework can be replaced with faster or more accurate backbone, which allows our framework to be updated easily and meet more demands.

For a better evaluation, we proposed Blurred-300VW and RWMB dataset with severe artificial and real-world motion blur respectively. Extensive experiments demonstrate the effectiveness of our framework and the superior performance to state-of-the-art methods on both datasets. In conclusion, the contributions of this paper are:

(1) A framework, in which a face deblurring network and landmarks detector work as a *virtuous circle* and obtain state-of-the-art performance in motion-blurred videos.

(2) A novel component, structure predictor, which utilizes temporal information to provide reliable face structure.

(3) Two new datasets, (*i.e.*, Blurred-300VW and RWMB) which are more suitable benchmarks for video facial landmark detection tasks.

## 2. Related work

### 2.1. Facial Landmark Detection

**Facial landmark detection in static images.** The classic model-based methods ASMs [27], AAMs [7, 17, 23, 32], CLMs [19, 33], ESR [3], SDM [41], CFSS [46], and deep convolutional neural wetwork methods, *e.g.* TCDCN [45], FAN [1], DSRN [26], RFLD [25], SAN [8], LAB [38] have obtained increasingly excellent performance in static images under different poses, light conditions, expressions, etc. However, few works in the literature of face alignment pay attention to the motion blur. They cannot maintain outstanding performance in case of severe motion blur. Our paper looks into blurry scenario and propose a framework to remove the motion blur, and then promote all these state-of-the-art facial landmark detection algorithms.

**Facial landmark detection in videos.** To overcome challenging problems like large pose and occlusion, facial landmark detection is naturally extended to videos leveraging temporal information [4]. Xi *et al.* [30] propose a recurrent encoder-decoder network (RED), combined with spatial and temporal recurrent learning, to explicitly model the temporal dependency relationship on frames. Hao *et al.* [21] construct the two-stream transformer network (TSTN), where the temporal stream learns to capture the continuous consistency across multiple frames in video clips and the spatial stream is capable of locating landmarks. These [40, 34, 30, 11, 21, 10] are typical works employing temporal information in facial landmark detection. However, few works have ever noticed the motion blur, a common problem and challenge in most videos. Moreover, face structure information has not been paid attention to in these video-based methods. In our paper, we propose a framework that leverages the temporal and structural information to tackle the problem of motion blur.

### 2.2. Motion Deblurring

Motion blur, which usually happens for the mismatching of the motion speed and shutter closing speed, could be entailed with object movement, camera shake, etc. The unlimited nature of motion makes deblurring a complex problem. Recent years have witnessed great progress in general image deblurring [42, 43, 20, 29, 37, 44, 5]. Compared with image deblurring, video deblurring [16, 18, 36] can utilize temporal information to handle large motion blur with fewer network parameters. However, these methods are not specialized for face and therefore have not leveraged face structural information.

Structure information could efficiently assist deblurring. Prior knowledge, especially facial structure [15, 2, 35], has been proven to be an effective face prior in corresponding tasks such as super-resolution and deblurring. However, as mentioned in [35], these methods fail when the input face
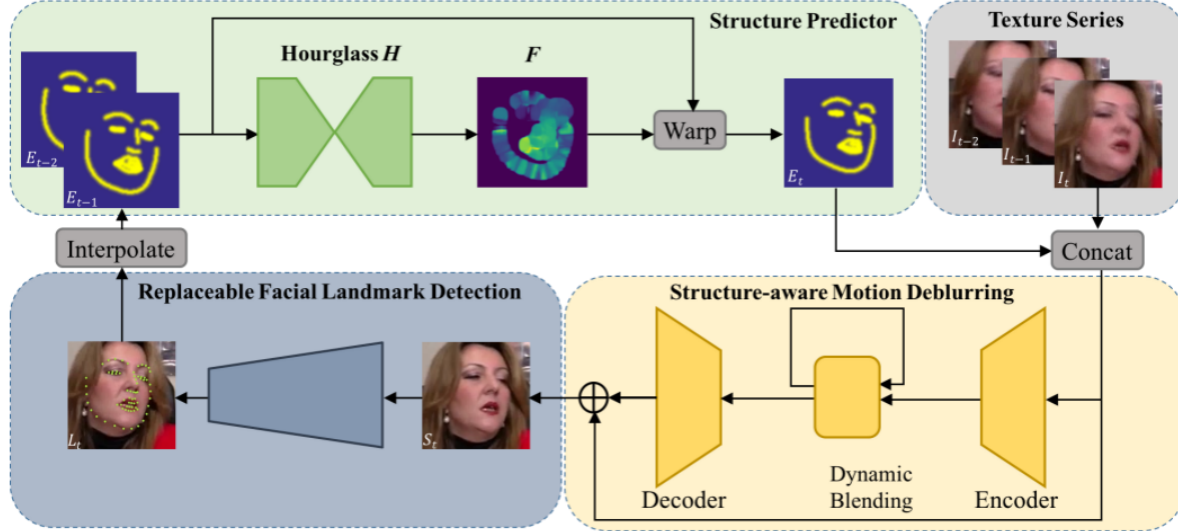
Figure 2: Framework. Given two previous face edges $E_{t-2}$ and $E_{t-1}$, hourglass $\boldsymbol{H}$ predicts the optical flow $F$ between the two boundary maps. Warping block $\boldsymbol{W}$ warps $E_{t-2}$ and $E_{t-1}$ into next boundary map $E_t$ according to the predicted optical flow $F$. Recent frames $I_{t-2}, I_{t-1}, I_t$, concatenated with the predicted boundary map, are feed to the Boundary-aware Deblur Network $\boldsymbol{D}$, which produces a sharp face $S_t$. Taking the deblurred sharp image $S_t$ as input, the replaceable facial landmark detection network predicts more accurate landmark location $L_t$, which is interpolated to face edges and provides the face structural information for next loop.

images are not well aligned, *e.g.* side faces or extremely large motion where semantic face parsing or landmark detection fails. Considering all these drawbacks, we designed the framework, which exploits temporal information to provide reliable structure information and furthermore removes the blur deeply.

## 3. Method

As shown in the figure 2, the proposed FAB can be divided into three conjoined components: structure predictor, structure-aware deblurring network and replaceable facial landmark detection network. The structure predictor predicts facial structure prior for the current frame from previous frames. With help from the predicted facial structure prior, the structure-aware deblurring network generates a clear image. Given a clear input image, the facial landmark detection network produces accurate facial landmarks for the current frame, which is fed back to the structure predictor for the next frame prediction. This makes the three networks an organic whole to perform facial landmarks localization and deblurring simultaneously and benefit each other in motion-blurred videos.

### 3.1. Structure Predictor

Structure predictor plays a key role in breaking the "chicken or egg" dilemma between deblur and facial landmark detection problems. It takes input facial structure of time $t-2$ and $t-1$ to predicts the facial structure of current frame $t$. The motion caused by camera shake or object motion can be modeled mathematically and are continuous

in a short time. Moreover, facial structure is a semantic meaningful, clear and well-defined representation of face regardless of the face texture. The motion between two facial structures is much easier to extract than motion between two face images. These properties make short time facial structure prediction a feasible problem.

A facial structure could be represented by landmarks, edges, part segments, 3D models *etc*. Landmarks are not semantically stable as they would drift along edges under different expressions and poses. The drift of landmarks adds noises to face motion which is harmful in structure prediction. On the other hand, complex annotation such as part segments and 3D models contains richer information but the large-scale dataset is hard to obtain. Thus we use face edges as facial structure information. The edges are interpolated from landmarks of each facial component.

Given two previous face edges $E_{t-2}$ and $E_{t-1}$, following previous paper [22], we use an hourglass $\boldsymbol{H}$ to predict the optical flow $F$ between the two boundary maps. Directly predicting face edges $E_t$ is also feasible but can hardly guarantee the sharpness property of the boundary map. However, this property is essential for following the deblur network. Then the warping block $\boldsymbol{W}$ would warp $E_{t-2}$ and $E_{t-1}$ into the next boundary map $E_t$ according to the predicted optical flow $F$:

$$\boldsymbol{H}(E_{t-2}, E_{t-1}) = \boldsymbol{F} \qquad (1)$$

$$\boldsymbol{W}(E_{t-2}, E_{t-1}, \boldsymbol{F}) = E_t \qquad (2)$$

where $t$ means the current time.

The facial structure predictor is pre-trained using facial landmark detection video dataset like 300VW, then fine-tuned together with other networks. For pretraining, the mean squared error (MSE) between the predicted and ground truth face boundary map. The loss is defined as

$$\mathcal{L}_{prd} = \frac{1}{N_{pixel}} \cdot \|E_t - E_{GT}\|_2 \qquad (3)$$

Where $N_{pixel}$ means the total pixel number in the generated image, and $E_{GT}$ is the ground truth of the current face edges. In the final experiment, we show that the motion of edges can be predicted accurately.

### 3.2. Structure-Aware Motion Deblurring

Face deblurring is difficult itself. Face motion is inconsistent with the motion blur degree because of the use of face detection/tracking. As is shown in figure 3, the face remains relatively static to the bounding box, but gradually getting blurry. However, face edges are exempt from these problems(see the second row of figure 3), and thus much easier to predict. With the reliably predicted structure prior, the deblurring network entangles the textural and structural information and reconstructed a deblurred face.

We disentangled a face into two parts, *i.e.* structure, and texture. Given accurate structure prior, boundary map, the deblur network $\boldsymbol{D}$ reconstructs a sharp image, which meets the boundary constraint and maintains the texture consistency simultaneously. Face structural information guides the motion deblurring of current frame. This process could also be viewed as filling the predicted structure with the texture from previous frames. Following state-of-the-art deblur network design [18, 37], an batch of three recent frames, $I_{t-2}, I_{t-1}, I_t$ are used as input. Then we concatenate the predicted boundary map, as structure prior, with the batch of frames as input to our deblur network. An encoder network extracts needed information from inputs. Then we use a dynamic temporal blending network, following [37] and [18], to combine information across different frames. Finally, a decoder is used to predicts the residual between blurry frames and groundtruth.

$$\boldsymbol{D}(E_t, I_{t-2}, I_{t-1}, I_t) = S_t \qquad (4)$$

Similar to facial structure predictor, we also pre-train structure-aware deblur network using 300VW. First, blurred videos are generated by using methods from [18]. Then the ground-truth edges are used as structure prior to pre-train this network. We employ the mean squared error $\mathcal{L}_{rec}$ between the deblurred image and the ground truth sharp image:

$$\mathcal{L}_{rec} = \frac{1}{N_{pixel}} \cdot \|S_t - S_{GT}\|_2 \qquad (5)$$
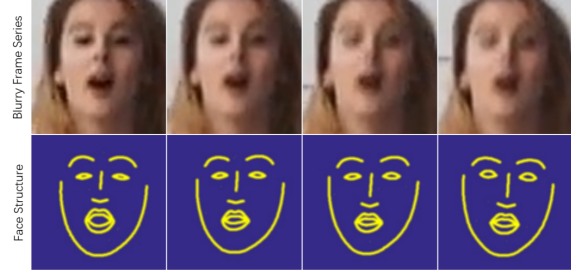


Figure 3: Face structural information is reliable. In the first row, the face remains relatively static to the bounding box, but gradually getting blurry. However, as is shown in the second row, face edges are exempt from these problems, and thus much easier to predict.

### 3.3. Replaceable Facial Landmark Detection Network

By using the deblurred sharp image $S_t$ as input, the replaceable facial landmark detection network predicts the landmark location $L_t$. The facial landmark detection network can be any advanced network architecture. In this paper, we use the residual network with pre-activation [13, 12]. And L1 distance $\mathcal{L}_{align}$ is used as the loss for facial landmark detection. The network is first pre-trained on facial landmark detection dataset and fine-tuned together with the other two networks:

$$\mathcal{L}_{align} = \frac{1}{N_{point}} \cdot \|L_t - L_{GT}\|_1 \qquad (6)$$

where $L_{GT}$ means the ground truth facial landmark location. $N_{point}$ is the number of landmarks.

The deblurred frames would lead to better performance for facial landmark detection task regardless of what the facial landmark detection network is used in our framework. Almost every single landmark localization networks would have much better performance in motion-blurred videos. Moreover, the more accurate facial structure would further provide better information for structure predictor in the next loop, which forms a virtuous circle.

### 3.4. Alternate Fine-Tuning

Structure predictor, structure-aware motion deblurring network and facial landmark detection network mentioned above constitute our main architecture. The three networks are pretrained separately and finally alternately end-to-end fine-tuned together to minimize the total loss:

$$\mathcal{L}_{total} = \mathcal{L}_{str} + \mathcal{L}_{rec} + \mathcal{L}_{align} \qquad (7)$$

where the $\mathcal{L}_{str}$ $\mathcal{L}_{rec}$ $\mathcal{L}_{align}$ are calculated according to Equation 3, 5 and 6 respectively.

During the end-to-end fine-tuning, the three networks benefit each other to achieve better performance. The losses
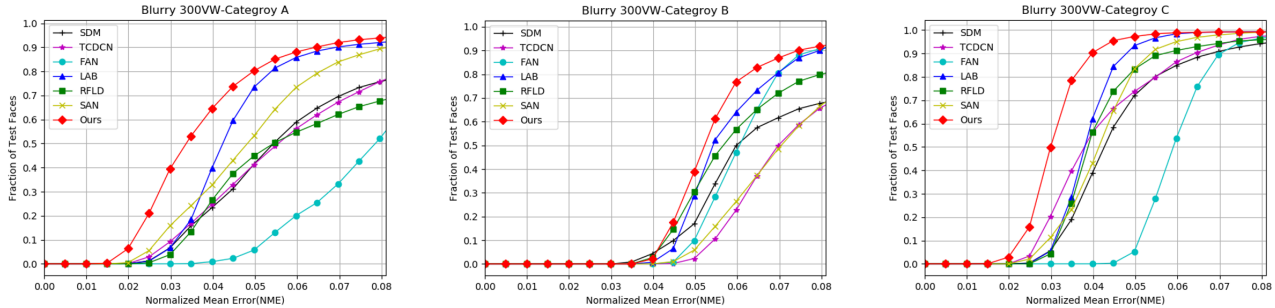
Figure 4: CED for Blurred-300VW Testset.

of other networks are also helpful. For example, with better deblur results, the alignment network can achieve higher accuracy. Thus the loss of facial landmark detector also helps optimizing face deblur network. The gradient of facial landmark detector would propagate back to encourage sharper input as well as pre-stage deblur network.

# 4. Implement Details

## 4.1. Network Structure

Our proposed algorithm is composed of three blocks, namely structure predictor, structure-aware motion deblurring network, and a replaceable facial landmark detection network. We used eight residual blocks in total to build the Hourglass [28] in the Structure Predictor. There are two convolutional layers and four residual blocks correspondingly in the Encoder and Decoder in the structure-aware motion deblurring network. Since the purpose of our work lies on the overall system framework instead of the design of the network or the loss function, we employed a simple network (*i.e.* pre-activated Resnet-18 [13]) as our replaceable facial landmark detection network. We present detailed information about the network structure in the first section of the supplementary material.

## 4.2. Training

All the three components are required to be pretrained respectively. For structure predictor, facial edges are implemented as inputs. We fit annotated landmarks face edges by cubic spline interpolation [24] and use them as input and ground truth. For structure-aware motion deblurring network, annotation results for the current frame are interpolated to face edge and used for structural information. For facial landmark detection network, we trained the model on 300W dataset [31], then finetune the network on the 300VW dataset [34]. Data augmentation such as translation, rotation, flipping, and zooming is also used in the total training stage.

During the end-to-end training, we designed an alternate training method to obtain better reconstruction and land-

mark localization results. Facial landmark detection network and the structure predictor take turns to be trained for one epoch. In this way, the landmark detection network would take deblurred images as input instead of original images.

# 5. Experiments

## 5.1. Datasets and Evaluation Metric

**300W [31] and 300VW [34].** 300W and 300VW are popular benchmarks for facial landmark localization methods. 300W contains 3,148 training images and 689 testing images. 300VW is designed as a benchmark for videos, containing 50 training videos and 64 testing videos. The testing set of 300VW is divided into three scenarios, *i.e.* category A: well-lit arbitrary expressions, category B: unconstrained illuminations and category C: arbitrary conditions. In this paper, these two datasets are combined in training.

**Blurred-300VW.** In order to generate more severe motion blur, larger motion in original videos is a need. We select subsets (31 videos in the training set and 9 videos in the testing set) from 300VW with large motion according to a face motion intensity index. The face motion intensity is defined by accumulating the movement of the left eye during a time unit and normalizing it with the inter-ocular distance. As is shown in Figure 5 (a), the selected videos are characterized with far more severe motion intensity.

Then, we blurred these picked subsets, following the method in [18]. For each adjacent three frames, we interpolated them with 20 subframes according to the optical flow [14]. The mean value of these 20 subframes is calculated to mimic motion blur. Annotation of each generated frame is taken from the middle-time subframe. The blurred-300VW dataset contains obvious motion blur especially when the face moves greatly, which is suitable for the illustration of this work.

Moreover, to evaluate the limit of recent facial landmark detection methods, we further form a Challenge Subset, which contains hundreds of images from 6 test videos with extreme motion blur. The comparisons of the blurry de-
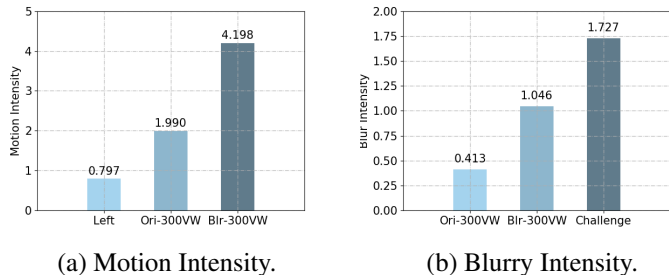
(a) Motion Intensity.      (b) Blurry Intensity.

Figure 5: (a) Comparison of the motion intensity among the left videos of original 300VW, images of original 300VW and images of blurred-300VW. (b) Comparison of the blur intensity among the original 300VW, Blurred-300VW and Challenge Subset.



Figure 6: Comparison of random sampled corresponding images on original 300VW and Blurred-300VW.

gree for these three sets are shown in Figure 5 (b). Also, we compare images on 300VW and Blurred-300VW at the same time in Figure 6.

**RWMB.** Real-World Motion Blur (RWMB) is our newly proposed in-the-wild benchmark for facial landmark detection task in realistic motion-blurred videos. It currently contains 20 videos with obvious real-world motion blur picked from YouTube, which include dancing, boxing, jumping, *etc*. There are $35,540$ frames, which are all annotated with 98 landmarks following the protocol of WFLW [38]. Moreover, RWMB dataset will be further enlarged to hundreds of videos, including millions of frames. Please notice the news in the project page for more information.

Even though human beings could generally recognize the position of facial features of blurry faces, it is challenging to determine the specific location of each landmark. The annotation of the previous frame is presented to the annotator as a reference. Each frame is annotated by three expert annotators and checked by two quality inspector. As motion blur exists widely in regular videos, research in this scenario is challenging and meaningful. This benchmark aims at promoting the development of facial landmark detection in motion-blurred videos.

**Annotation protocol.** We provided two versions of annotation of RWMB and 300VW, *i.e*. 98 landmarks following [38] and 68 landmarks following [34]. In this way, we

| Method | CateA | CateB | CateC | Year |
|--------|-------|-------|-------|------|
| SDM [41] | 10.4 | 8.52 | 4.83 | 2013 |
| TCDCN [45] | 6.81 | 7.90 | 4.36 | 2016 |
| FAN [1] | 10.84 | 6.80 | 6.11 | 2017 |
| RFLD [25] | 10.95 | 8.19 | 4.63 | 2018 |
| SAN [8] | 5.62 | 7.71 | 4.35 | 2018 |
| LAB [38] | 5.28 | 6.07 | 3.96 | 2018 |
| **Ours** | **4.24** | **5.67** | **3.16** | - |

Table 1: NME(%) value of our method and several state-of-the-arts on Testset of Blurred-300VW.

unify these two video facial landmark datasets and hope to facilitate the cross-dataset evaluation in future works.

**Evaluation metric.** To evaluate the performance of facial landmark detection, Normalized Mean Error (NME), Cumulative Errors Distribution (CED) curve, Area Under the Curve (AUC) and Failure Rate are employed in this paper. Mean Error is normalized by inter-ocular distance (namely outer eye corner distance) throughout this paper.

## 5.2. Results

### 5.2.1 Evaluation on Blurred-300VW

This experiment is conducted to demonstrate that our algorithm is capable to handle *artificial* motion blur, and accurately detect facial landmarks.

We re-implemented several typical state-of-the-art algorithms for comparison, including RFLD [25], SAN [8] and LAB [38]. For a fair comparison, we made great efforts to reproduce the comparable results reported in their papers. The results on Blurred-300VW are shown in Table 1. The Experiment shows that our algorithm significantly outperforms previous methods by a large margin, and has improved the state-of-the-art performance (NME) by relatively 18.1% on average on the blurred-300VW dataset from 5.28 to 4.24. Figure 4 shows the CED curve of all these excellent algorithms. Note that SAN [8] and LAB [38] are rather robust in image-based facial landmark detection. However, when confronting the severe motion blur, they fail to yield reliable and robust facial landmarks due to the destruction of the structure information. Our method benefits from the reliable predicted facial structure and the deblurred input and obtains state-of-the-art performance.

### 5.2.2 Evaluation on RWMB

Although the experiments in Blurred-300VW have demonstrated that the proposed method could handle *artifical* motion blur, we still wonder how our method performs when facing *real-world* motion blur. Therefore, we propose this dataset, Real-World Motion Blur, as well as the testing result of several state-of-the-art algorithms to provide a

| Method | NME | NME threshold: 0.2 | | NME threshold: 0.1 | | NME threshold: 0.08 | |
|---|---|---|---|---|---|---|---|
| | | Failure Rate(%) | AUC | Failure Rate(%) | AUC | Failure Rate(%) | AUC |
| RFLD [25] | 16.34 | 25.67 | 0.59 | 44.84 | 0.38 | 54.39 | 0.29 |
| FAN [1] | 13.73 | 10.42 | 0.56 | 34.18 | 0.27 | 53.26 | 0.18 |
| TCDCN [45] | 10.48 | 11.27 | 0.60 | 34.82 | 0.41 | 45.79 | 0.34 |
| SAN [8] | 10.40 | 11.79 | 0.62 | 31.37 | 0.43 | 40.93 | 0.36 |
| LAB [38] | 9.47 | 8.28 | 0.61 | 32.20 | 0.43 | 41.12 | 0.34 |
| **Ours** | **8.43** | **5.25** | **0.63** | **28.77** | **0.45** | **38.54** | **0.39** |

Table 2: Failure rate and Area under curve(AUC), normalized by inter-ocular distance, on RWMB dataset
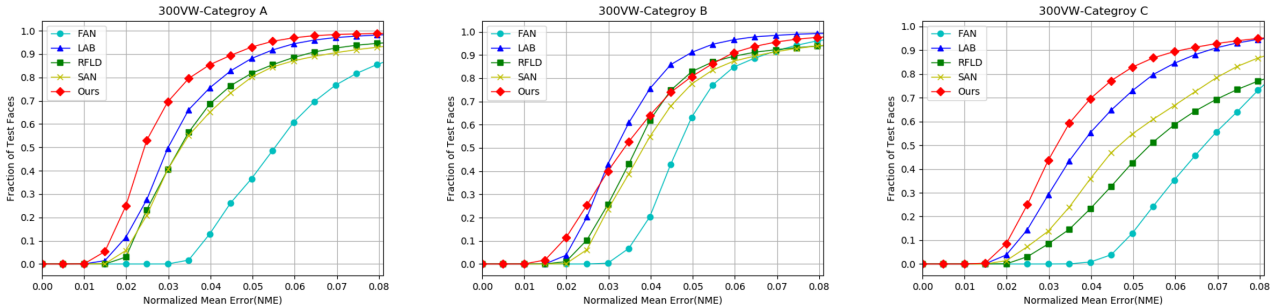


Figure 7: CED for original 300VW Testset.

benchmark.

We retrained several state-of-the-art algorithms as is introduced in Sec. 5.2.1 and tested them on the 68-points version of the RWMB. Comparison results on the RWMB against these methods are presented in Table 2. The experiment shows that our algorithm has improved the state-of-the-art performance (NME) by relatively 10.61% on average on the RWMB dataset from 9.47 to 8.43, significantly outperforms previous methods by a large margin.

Note that the RWMB is full of real-world motion blur, and thus rather difficult to handle. Most state-of-the-art methods fail to perform well on this dataset for the absence of reliable facial structure information. However, our method manages to remove the motion blur and obtain accurate landmarks. This experiment demonstrates that our algorithm is robust under real-world motion blur circumstance.

### 5.2.3 Evaluation on 300VW

Since most of the video facial landmark detection literature provide testing results on the 300VW dataset, we also evaluate our method on the 300VW dataset to make this paper complete.

As is shown in the Figure 7. There actually exists some motion-blurred frames in category A and C (even though much slight than Blurred-300VW dataset), in which our method largely obtains the state-of-the-art results. Since category B is designed to evaluate the robustness under different illuminations, dark rooms and overexposed shots

| Method | CateA | CateB | CateC | Challset |
|---|---|---|---|---|
| REDN [30] | - | - | - | 6.25 |
| TSTN [21] | 5.36 | 4.51 | 12.84 | 5.59 |
| Ours | **3.56** | **3.88** | **5.02** | **3.96** |

Table 3: Comparison with temporal methods.

(which is not the main scenario considered in our work), our method obtains comparable results. Note that comparing to other sophisticated methods, only a vanilla Resnet is used in our model. Generally, this experiment demonstrates that our method does not sacrifice accuracy on regular dataset when obtaining state-of-the-art performance on blurry datasets.

To compare with other methods using temporal information, we report the NME in the Table 3. Our method works well in the original 300VW dataset and outperforms the existing methods considering temporal information.

### 5.3. Ablation Study

In this section, we investigate the effectiveness of each pivotal component in the proposed framework on the Blurred-300VW Challenging Subset. We still employ the Resnet-18 as the basic facial landmark detector (FA) in our framework. Based on FA, we analyze each proposed component or potential designing, *i.e.* the Structure Predictor (SP), and the Structure-aware Motion Deblurring (SMD), by comparing the NME score. The overall result is presented in Figure 8 (a). We also analyze the different way to obtain facial edges, *i.e.* by state-of-the-art Structure Detec-
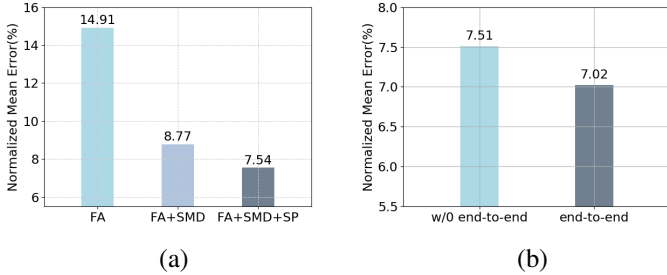
Figure 8: (a) Normalized Mean Error to analyze the effect of the pivotal components. (b) Normalized Mean Error to evaluate the effectiveness of alternate fine-tuning strategy.
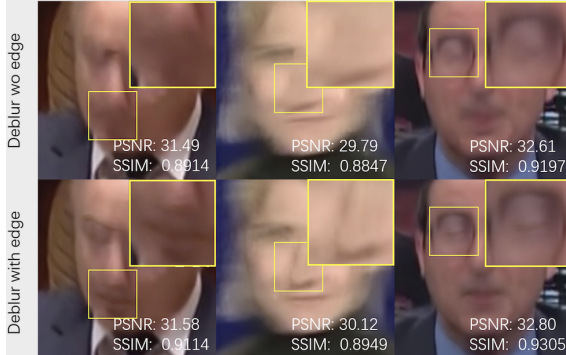


Figure 9: Deblurring w/wo structure information. images in the first row is deblurred by state-of-the-art algorithms, and images in the second row is produced by our algorithm.

tor(SD) or by our Structure Predictor (SP). The comparing results is presented in the Table 4.

**Effectiveness of pivotal component.** We first evaluate the effectiveness of the Structure-aware Motion Deblurring (SMD) for facial landmark detection tasks in blurred-videos. The first bar of the Figure 8 (a) shows the NME score of the baseline facial landmark detection method fed with blurry images without the assistance of SP or SMD. By adding the state-of-the-art deblurring network [18] before facial landmark detector (FA+SMD), we reduce the NME from 14.91% to 8.77%, as is shown in the first (FA) and the second bar (FA+SMD) in the Figure 8 (a). Then, we evaluate the effectiveness of Structure Predictor (SP). By leveraging the predicted facial edges with SP, we optimize the NME from 8.77 to 7.54 as is shown in the second (FA+SMD) and the third bar (FA+SMD+SP) of Figure 8 (a). As shown in Figure 9, though motion deblurring is a component rather than the main objective of this work, we still report the comparison of PSNR and SSIM of images deblurred by a state-of-the-art deblurring method [18] and our methods with SP. Superior results demonstrate the necessity and the efficiency of our Structure Predictor (SP).

**Choice of Structure Predictor.** One intuitive way to obtain the facial structure is to use a pre-trained state-of-the-art landmark detector. For comparison to our temporal-based

| Row | Method | NME (%) |
|-----|--------|---------|
| 1 | FA + SMD + SP (FAN [1]) | 8.52 |
| 2 | FA + SMD + SP (ResNet [13]) | 8.33 |
| 3 | FA + SMD + SP (SAN [8]) | 8.18 |
| 4 | FA + SMD + SP (LAB [38]) | 7.92 |
| 5 | **FA + SMD + SP(Ours)** | **7.54** |
| 6 | FA + SMD + GT | 6.72 |

Table 4: NME of different choices of Structure Predictor as the source of structure prior.

| Method | w/o Ours | w/ Ours | Improvement |
|--------|----------|---------|-------------|
| FAN [1] | 45.42 | 10.95 | 75.89% ↑ |
| ResNet [13] | 14.91 | 7.54 | 49.43% ↑ |
| SAN [8] | 9.21 | 7.02 | 23.78% ↑ |
| RFLD [25] | 41.24 | 20.07 | 51.33% ↑ |
| LAB [38] | 8.94 | 4.91 | 45.53% ↑ |

Table 5: NME(%) value of state-of-the-art methods, working with or without our framework. Our method enhance these methods efficiently.

structure predictor (SP), we report NME in the Table 4 of different choices of structure predictor (*i.e.*, FAN[1], ResNet [13], SAN [8], LAB [38]). Besides, GT is reported as the upper-bound performance of structure prior usage, in which ground-truth landmarks are used. Our proposed temporal-based structure predictor (SP) obtains the best NME score of 7.54%.

**Replaceable facial landmark detector.** To verify the feasibility of the replaceable facial landmark detector, we report the NME results in Table 5, which demonstrates that our framework could significantly enhance the performance of single networks by replacing Resnet with them in our framework. It is exciting and encouraging to see that our method enhances most state-of-the-arts significantly.

**Alternate fine-tuning strategy.** After pretraining each component, we end-to-end finetuned them as a whole. As is introduced in Figure 8 (b), the baseline (Resnet [13]) obtains obvious improvement after the finetuning.

## 6. Discussion and Conclusion

In this paper, we introduced an important but omitted issue, the facial landmark detection in motion-blurred videos. Since motion blur is ubiquitous in practical work, we believe further research on this field is of great significance. We proposed a novel framework, in which the three components work as an organic whole. The face boundary yielded by the structure predictor assists the deblurring module to relieve the motion blur. And the deblurred face leads to more accurate facial landmarks. There is still a lot of room for improvement for the extreme motion blur of the proposed RWMB dataset. We hope to see further development of this work in the future.

# References

[1] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem?(and a dataset of 230,000 3d facial landmarks). In *ICCV*, 2017. 2, 6, 7, 8

[2] Adrian Bulat and Georgios Tzimiropoulos. Super-fan: Integrated facial landmark localization and super-resolution of real-world low resolution faces in arbitrary poses with gans. 2018. 2

[3] Xudong Cao, Yichen Wei, Fang Wen, and Jian Sun. Face alignment by explicit shape regression. *International Journal of Computer Vision*, 2014. 2

[4] Grigorios G Chrysos, Epameinondas Antonakos, Patrick Snape, Akshay Asthana, and Stefanos Zafeiriou. A comprehensive performance evaluation of deformable face tracking. *International Journal of Computer Vision*, 2018. 2

[5] Grigorios G. Chrysos, Paolo Favaro, and Stefanos Zafeiriou. Motion deblurring of faces. *IJCV*, 2018. 2

[6] Grigorios G Chrysos and Stefanos Zafeiriou. Deep face deblurring. In *CVPR*, 2017. 2

[7] Timothy F Cootes, Gareth J Edwards, and Christopher J Taylor. Active appearance models. *TPAMI*, 2001. 2

[8] Xuanyi Dong, Yan Yan, Wanli Ouyang, and Yi Yang. Style aggregated network for facial landmark detection. In *CVPR*, 2018. 1, 2, 6, 7, 8

[9] Xuanyi Dong, Shoou-I Yu, Xinshuo Weng, Shih-En Wei, Yi Yang, and Yaser Sheikh. Supervision-by-registration: An unsupervised approach to improve the precision of facial landmark detectors. In *CVPR*, 2018. 1

[10] Grigorios G. Pd2t: Person-specific detection, deformable tracking. In *TPAMI*, 2018. 2

[11] Jinwei Gu, Xiaodong Yang, Shalini De Mello, and Jan Kautz. Dynamic facial analysis: From bayesian filtering to recurrent neural network. In *CVPR*, 2017. 2

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 4

[13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *ECCV*, 2016. 4, 5, 8

[14] Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial intelligence*, 1981. 5

[15] Yinghao Huang, Hongxun Yao, Sicheng Zhao, and Yanhao Zhang. Efficient face image deblurring via robust face salient landmark detection. In *Pacific Rim Conference on Multimedia*, 2015. 2

[16] Tae Hyun Kim and Kyoung Mu Lee. Generalized video deblurring for dynamic scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015. 2

[17] Fatih Kahraman, Muhittin Gokmen, Sune Darkner, and Rasmus Larsen. An active illumination and appearance (aia) model for face alignment. In *CVPR*, 2007. 2

[18] Tae Hyun Kim, Kyoung Mu Lee, Bernhard Schölkopf, and Michael Hirsch. Online video deblurring via dynamic temporal blending network. In *ICCV*, 2017. 2, 4, 5, 8

[19] Neeraj Kumar, Peter Belhumeur, and Shree Nayar. Facetracer: A search engine for large collections of images with faces. In *European conference on computer vision*, pages 340–353. Springer, 2008. 2

[20] Orest Kupyn, Volodymyr Budzan, Mykola Mykhailych, Dmytro Mishkin, and Jiri Matas. Deblurgan: Blind motion deblurring using conditional adversarial networks. 2018. 2

[21] Hao Liu, Jiwen Lu, Jianjiang Feng, and Jie Zhou. Twostream transformer networks for video-based face alignment. *IEEE transactions on pattern analysis and machine intelligence*, 2018. 2, 7

[22] Ziwei Liu, Raymond A Yeh, Xiaoou Tang, Yiming Liu, and Aseem Agarwala. Video frame synthesis using deep voxel flow. In *ICCV*, 2017. 3

[23] Iain Matthews and Simon Baker. Active appearance models revisited. *International journal of computer vision*, 2004. 2

[24] Sky McKinley and Megan Levine. Cubic spline interpolation. *College of the Redwoods*, 45(1):1049–1060, 1998. 5

[25] Daniel Merget, Matthias Rock, and Gerhard Rigoll. Robust facial landmark detection via a fully-convolutional local-global context network. In *CVPR*, 2018. 1, 2, 6, 7, 8

[26] Xin Miao, Xiantong Zhen, Xianglong Liu, Cheng Deng, Vassilis Athitsos, and Heng Huang. Direct shape regression networks for end-to-end face alignment. In *CVPR*, 2018. 2

[27] Stephen Milborrow and Fred Nicolls. Locating facial features with an extended active shape model. In *ECCV*, 2008. 2

[28] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. *CoRR*, 2016. 5

[29] Haesol Park and Kyoung Mu Lee. Joint estimation of camera pose, depth, deblurring, and super-resolution from a blurred image sequence. In *ICCV*, 2017. 2

[30] Xi Peng, Rogerio S Feris, Xiaoyu Wang, and Dimitris N Metaxas. A recurrent encoder-decoder network for sequential face alignment. In *ECCV*, 2016. 2, 7

[31] Christos Sagonas, Georgios Tzimiropoulos, Stefanos Zafeiriou, and Maja Pantic. 300 faces in-the-wild challenge: The first facial landmark localization challenge. In *ICCV-W*, 2013. 5

[32] Jason Saragih and Roland Goecke. A nonlinear discriminative approach to aam fitting. In *Computer Vision, 2007. ICCV 2007. IEEE 11th International Conference on*, 2007. 2

[33] Jason M Saragih, Simon Lucey, and Jeffrey F Cohn. Deformable model fitting by regularized landmark mean-shift. *International Journal of Computer Vision*, 2011. 2

[34] Jie Shen, Stefanos Zafeiriou, Grigoris G Chrysos, Jean Kossaifi, Georgios Tzimiropoulos, and Maja Pantic. The first facial landmark tracking in-the-wild challenge: Benchmark and results. In *ICCV-W*, 2015. 2, 5, 6

[35] Ziyi Shen, Wei-Sheng Lai, Tingfa Xu, Jan Kautz, and Ming-Hsuan Yang. Deep semantic face deblurring. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 2

[36] Shuochen Su, Mauricio Delbracio, Jue Wang, Guillermo Sapiro, Wolfgang Heidrich, and Oliver Wang. Deep video deblurring for hand-held cameras. In *CVPR*, 2017. 2

[37] Patrick Wieschollek, Michael Hirsch, Bernhard Schölkopf, and Hendrik PA Lensch. Learning blind motion deblurring. In *ICCV*, 2017. 2, 4

[38] Wayne Wu, Chen Qian, Shuo Yang, Quan Wang, Yici Cai, and Qiang Zhou. Look at boundary: A boundary-aware face alignment algorithm. In *CVPR*, 2018. 1, 2, 6, 7, 8

[39] Wayne Wu and Shuo Yang. Leveraging intra and inter-dataset variations for robust face alignment. In *CVPR*, 2017. 1

[40] P. Perona X. P. Burgos-Artizzu, D. Hall and P. Dollr. Merging pose estimates across space and time. 2013. 2

[41] Xuehan Xiong and Fernando De la Torre. Supervised descent method and its applications to face alignment. In *CVPR*, 2013. 1, 2, 6

[42] Li Xu and Jiaya Jia. Two-phase kernel estimation for robust motion deblurring. In *European conference on computer vision*, 2010. 2

[43] Li Xu, Shicheng Zheng, and Jiaya Jia. Unnatural l0 sparse representation for natural image deblurring. In *CVPR*, 2013. 2

[44] Jiawei Zhang, Jinshan Pan, Jimmy Ren, Yibing Song, Linchao Bao, Rynson WH Lau, and Ming-Hsuan Yang. Dynamic scene deblurring using spatially variant recurrent neural networks. In *CVPR*, 2018. 2

[45] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. Facial landmark detection by deep multi-task learning. In *ECCV*, 2014. 1, 2, 6, 7

[46] Shizhan Zhu, Cheng Li, Chen Change Loy, and Xiaoou Tang. Face alignment by coarse-to-fine shape searching. In *CVPR*, 2015. 1, 2