

Automatic Facial Paralysis Estimation with Facial Action Units

Xuri Ge, Joemon M. Jose, Pengcheng Wang, Arunachalam Iyer, Xiao Liu, Hu Han, *Member, IEEE*,

Abstract—Facial palsy is unilateral facial nerve weakness or paralysis of rapid onset with unknown cause. Automatically estimating facial palsy severeness can be helpful for the diagnosis and treatment of people suffering from it across the world. In this work, we develop and experiment with a novel model for estimating facial palsy severity. For this, an effective Facial Action Units (AU) detection technique is incorporated into our model, where AUs refer to a unique set of facial muscle movements used to describe almost every anatomically possible facial expression. In this paper, we propose a novel **Adaptive Local-Global Relational Network (ALGRNet)** for facial AU detection and use it to classify facial paralysis severity. ALGRNet mainly consists of three main novel structures: (i) an adaptive region learning module that learns the adaptive muscle regions based on the detected landmarks; (ii) a skip-BiLSTM that models the latent relationships among local AUs; and (iii) a feature fusion&refining module that investigates the complementary between the local and global face. Quantitative results on two AU benchmarks, *i.e.*, BP4D and DISFA, demonstrate our ALGRNet can achieve promising AU detection accuracy. We further demonstrate the effectiveness of its application to facial paralysis estimation by migrating ALGRNet to a facial paralysis dataset collected and annotated by medical professionals.

Index Terms—Facial action units, Facial paralysis estimation, Facial palsy, Skip-BiLSTM, Fusion&Refining

1 INTRODUCTION

RECENTLY, facial action units detection has attracted increasing research attention in computer vision due to its wide range of potential applications in facial state analysis, *i.e.*, diagnosing mental disease [1], face recognition [2], improving e-learning experiences [3], deception detection [4], *etc.* On a similar line to these applications, facial palsy is affecting a large number of population and automatic estimation of facial palsy severeness can be useful for both diagnosis and treatment. Facial palsy is the temporary or permanent weakness or lack of movement affecting one side of the face and is an acute, unilateral facial nerve weakness or paralysis of rapid onset (less than 72 hours) and unknown cause. It affects around 23 per 100,000 people per year, and current methods for facial palsy severity estimation is a relatively subjective process. However, in the literature it is rare to see such studies that extending the AU detection model to the assessment of facial paralysis grades. In fact,

the severity of facial paralysis can be estimated by the manifestations of muscle areas of the face together, which is really similar to the representation of individual expressions using the Facial Action Coding System (FACS) [5]. In this study, we propose a novel AU detection model and explore its ability to estimate facial paralysis severity.

From a biological perspective, the activation of AU corresponds to the movement of facial muscles; however, AU detection is challenging because of the subtle facial changes caused by AU. Hand-crafted features are used to represent the appearance of different local facial regions in early works [6], [7]. In fact, this also applies to some works on facial state analyses, such as facial paralysis estimation [8], [9] and patient pain detection [10]. However, due to their shallow nature, hand-crafted features are not discriminative enough to depict the morphological variations in the face. In order to improve the feature representation of AUs, deep learning-based approaches for AU detection have been developed in recent years.

To improve the AU representation, most existing facial AU detection methods combine local features from numerous independent AU branches, each corresponding to a separate AU patch. As shown in Fig. 1 (a), some grid-based deep learning studies [11], [12] combine regional (patch-based) Convolutional Neural Network (CNN) features from a face with equal grids. For example, LP-Net [13] using an LSTM model [14] to combine the local CNN features from equal partition grids. However, there are two significant issues with dividing the image into fixed grids; (i) It is hard to focus accurately on the muscle region related to each AU patch; and (ii) grid-based features may not adequately reflect the ROIs of irregularly shaped AU patches. Recent popular multi-branch combination-based methods [15], [16], [17] refine the AU-related features with irregular regions by fusing global or local features from independent

We thank Professor Brian O'reilly for sharing part of the facial paralysis dataset. This work was supported in part by China Scholarship Council (CSC) from the Ministry of Education of China (No. 202006310028).

- Xuri Ge is with the School of Computing Science, University of Glasgow, Scotland, UK (e-mail: x.ge.2@research.gla.ac.uk).
- Joemon M. Jose is with the School of Computing Science, University of Glasgow, Scotland, UK (e-mail: joemon.jose@glasgow.ac.uk).
- Pengcheng Wang is with Tomorrow Advancing Life Education Group (TAL), Beijing 100080, China (e-mail: wangpengcheng2@tal.com).
- Arunachalam Iyer is with the Department of Otolaryngology and Head and Neck Surgery, University Hospital Monklands, Airdrie, Scotland, UK, and also with University of Glasgow, Scotland, UK (aruniyerent@gmail.com).
- Xiao Liu is with the Online Media Business Unit at Tencent, Beijing 100080, China (e-mail: ender.liux@gmail.com).
- Hu Han is with the Key Laboratory of Intelligent Information Processing, Institute of Computing Technology, Chinese Academy of Sciences, Beijing 100190, China, and also with Peng Cheng Laboratory, Shenzhen 518055, China (e-mail: hanhu@ict.ac.cn).

Manuscript received April 19, 2005; revised August 26, 2015.

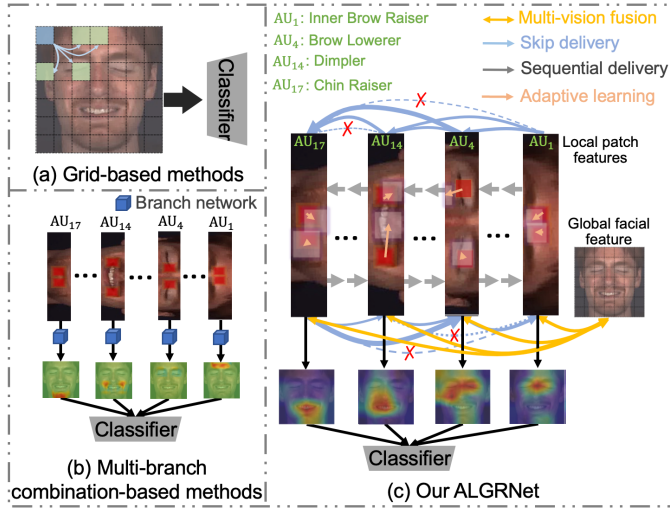


Fig. 1. Illustration of the different schemes: (a) the traditional grid-based feature extraction and classification, (b) the popular multi-branch combination-based detection methods, and (c) our ALGRNet method. Our ALGRNet, in comparison with (a) and (b), adaptively adjusts the AU areas in terms of different individuals based on detected landmarks, exploits mutual facilitation and inhibition of region-based multiple branches through a novel bidirectional structure with skipping gates and refines their irregular representations guided by the global facial feature.

AU branches based on the detected corresponding muscle region, as shown in Fig. 1 (b). For example, the scheme in [18] joints face alignment and AU detection in an end-to-end architecture, designed to extract corresponding AU features using multiple branches based on the detected and calculated AU centre coordinates. Furthermore, the latest approach [17] proposes a local AU recognition loss that refines the local attention map by the near-region pixel contributions for local regions in independent branches, rather than a pre-defined attention map; however, which also ignores the interrelationship between multiple AU areas of each face.

While AU detection methods based on multi-branch combinations have shown their success in the fusion of local AU features, limitations remain in establishing interrelationships between adaptive AU regions and in modelling local and global context. Firstly, unlike the localisation of corresponding AU regions based on fixed landmarks, adaptive learning of AU regions (shape and position) can improve the robustness of the model based on the diversity of expressions and individual characteristics, which however is usually neglected by the exiting literature. Secondly, according to the statistics of FACS [5], some patches corresponding to AUs are strongly correlated in some specific expressions (here we define positive correlation (mutual assistance) if multiple patches jointly influence the activation of the target AU, otherwise negative correlation (mutual exclusion)). For instance, the cheek area and the mouth corner of the face usually active simultaneously in a common facial behaviour called Duchenne smile, resulting in high correlations between AU6 (cheek raiser) and AU12 (lip corner puller). Furthermore, due to muscle linkage, adjacent AU2 (Outer Brow Raiser) and AU7 (Lid Tightener) will usually be activated simultaneously as the startle. Inspired by these

biological phenomena, we believe that it is important to capture the interactive relationship between patch-based branches, such as sequential/skipping information transfer of adjacent/non-adjacent related muscle regions, to enhance the AU features. Furthermore, the muscle activation areas of AU are often irregular due to individual and expression differences, and some non-AU areas are also commonly activated due to muscle linkage. We therefore argue that it is vital to use the information of global faces to complement the personalisation of each AU in terms of different individuals and expressions.

To this end, we propose a novel ALGRNet for AU detection and apply it to the facial paralysis estimation task to validate its robustness and transferability. Specifically, grid-based global features are extracted from a stem network consisting of multiple convolutional layers and local AU features are extracted from the calculated regions based on the detected AU centres. Different from previous methods, we calculate AU centres to suit different individuals and expressions by learning the landmarks and corresponding offsets. To catch the potential positive and negative relations among the branches, a skip-BiLSTM module is designed. To model the mutual support and exclusion information, the adjacent patches are transferred in BiLSTM [19] while the distant patches are connected via skipping-type gates. We model each branch as independent and equal and hence our kip connection manner can minimize the loss of information compared with traditional BiLSTM. Subsequently, in order to fuse global features to each local AU feature and also with even non-AU regional features, and in contrast to the previous approaches [17], a novel gated fusion architecture in the new feature fusion&refining module is proposed. This is important considering that the different AUs extracted from the face may focus on different information across the face regions. Finally, AU features combining other beneficial AU regions as well as non-AU regions are fed into a multi-branch classification network for AU detection or facial palsy class estimation.

Our contributions can be summarized as follows:

- We propose a skip-BiLSTM module to improve the robustness of local AU representations by modeling the mutual assistance and exclusion relationships of individual AUs based on the learned adaptive landmarks;
- We propose a feature fusion&refining module, filtering information that contributes to the target AU, even non-AU areas, to facilitate more discriminative local AU feature generation;
- The proposed ALGRNet has established new state-of-the-art performance for AU detection on two benchmark datasets, *i.e.*, BP4D and DISFA, without any external data or pre-trained models in additional data. Notably, we exploit a facial paralysis dataset, named FPara, to verify that the proposed AU detection model can be applied to facial paralysis estimation and achieves superior performance than the baseline methods.

In comparison to the earlier conference version [20] of this work, we propose a new adaptive region learning module in Section 3.2 in order to further improve the accuracy

of muscle regions and to better adapt to irregular muscle shapes. In particular, the adaptive region learning module contains learning of scaling factors to change the size of the corresponding muscle regions, as well as offset learning to slightly adjust landmark differences, with respect to different individuals. This suggests that adaptive region learning could better help the model to focus more accurately on the muscle region changes corresponding to each AU, and to obtain stronger robustness and generalisation ability. In addition, we did not evaluate the generalizability and transferability of the AU detection presented in the previous version, whereas we attempt it in the present study. Facial paralysis estimation is a time-consuming and subjective task for a traditional physician's diagnosis. Facial palsy is a condition characterised by motor dysfunction of the muscles of facial expression. It is usually qualified by observing the activation status of certain muscle areas and facial symmetry when the patient makes certain expressions, such as basic eyebrow raising, eye closing and mouth puckering, *etc.* In this study, we apply the proposed ALGRNet on facial paralysis estimation, which can improve the effectiveness of facial paralysis recognition and estimation by focusing on the activation of multiple muscle regions as well as global facial information. Specifically, we exploit a facial paralysis dataset which is annotated by medical professions to four grades of facial palsy degrees, *i.e.* normal, low, medium and high grade. For facial paralysis estimation we focus on the muscle areas that are more preferred in the facial paralysis ratings rather than the AU predefined muscle regions. Finally, we combine the multiple muscle region features enhanced by the interaction, as well as the useful global information, to obtain the final facial features for the facial palsy grade classification. To the best of our knowledge, there is no existing work in the literature on the estimation of facial paralysis using AU recognition methods. And we show the effectiveness and transferability of the proposed ALGRNet quantitatively through the Facial Palsy Assessment application, which was not present in our earlier conference version [20].

2 RELATED WORK

2.1 Facial Action Units Detection

Automatic AU detection is a task that detects the movement of a set of facial muscles. Recently, patch-learning based methods are the most popular paradigms for AU detection [21], [22], [23], [24], [25], [26], [27], [28]. For instance, [29] used the composition rules of different fixed patches for different AUs to recover facial expressions by sparse coding. [30] used a CNNs and BiLSTM to extract and model the image regions for AUs, which are pre-select by domain knowledge and facial geometry. However, all above methods need to pre-defined the patch location first. To address these issues, [18] proposed to jointly estimate the location of landmarks and the presence of action units in an end-to-end framework, where landmarks can be used to compute the attention map for each AU separately. Recent works [31], [32], [33], [34], [35], [36] explicitly take into consideration the linkage relationship between different AUs for AU detection, which rely on action unit relationship modeling to help improve recognition accuracy. Typically, [31], [37]

exploited the relationships between AU labels via a dynamic Bayesian network. [38] embedded the relations among AUs through a predefined graph convolutional network (GCN). [39] integrated the prior knowledge from FACS into an offline graph, which can construct a knowledge graph coding the AU correlations. However, these methods require prior connections by counting the co-occurrence probabilities in different datasets in advance. [36], [40], [41] applied a adaptive graph to model the relationships between AUs based on global feature, ignoring local-global interactions.

The most relevant previous studies to ours are [17], [18], which combine AU detection and face alignment into a multi-branch network. Different from these methods, our ALGRNet can adaptively adjust the target muscle region corresponding to each AU and utilizes the learned mutual assistance and exclusion relationships between the target muscle and other muscle regions to enhance the feature representation of the target AU. Doing so allows us to provide more robustness and interpretability than [17].

2.2 Facial Paralysis Estimation

Facial paralysis estimation has recently attracted extensive research attention [42], [43], due to the significant psychological and functional impairment to the patients. Nottingham system [44] is a widely accepted system for the clinical assessment of facial nerve function, which is similar with House-Brackmann (H-B) [45]. In addition to these, there are over twenty other methods of recognising and assessing facial palsy that are available in the literature. However, all of the above traditional methods are estimated by medical professionals and are both time consuming and subjective. More Recently, deep learning has been widely applied for facial representation learning, including using the deep representation for face recognition, face alignment, *etc.* [46] and [47] proposed two efficient quantitative assessment of facial paralysis based on the detected key points. [8] proposed to obtain the facial paralysis degree by calculating the changes of the surface areas of specific facial region. [48] considered both static facial asymmetry and dynamic transformation factors to evaluate the degree of facial paralysis. However, most of the exiting methods were only used the deep learning methods to pave the way for physical computation and do not directly model and predict the depth features of a face image. In addition, they apply some post-processing to obtain the final result.

In contrast to these existing methods, we employ an end-to-end deep framework ALGRNet to predict the grade of facial paralysis. We perform an automated estimation by combining depth features from high interest muscle regions with feature information from the global face, without any post-processing.

3 APPROACH

The framework of the proposed ALGRNet for AU detection is presented in Fig. 2. It is composed of four main modules, *i.e.*, adaptive region learning module (Subsection 3.2) for adaptive muscle region localisation, a skip-BiLSTM module (Subsection 3.3) for mutual facilitation and inhibition modeling, a feature fusion&refining module (Subsection

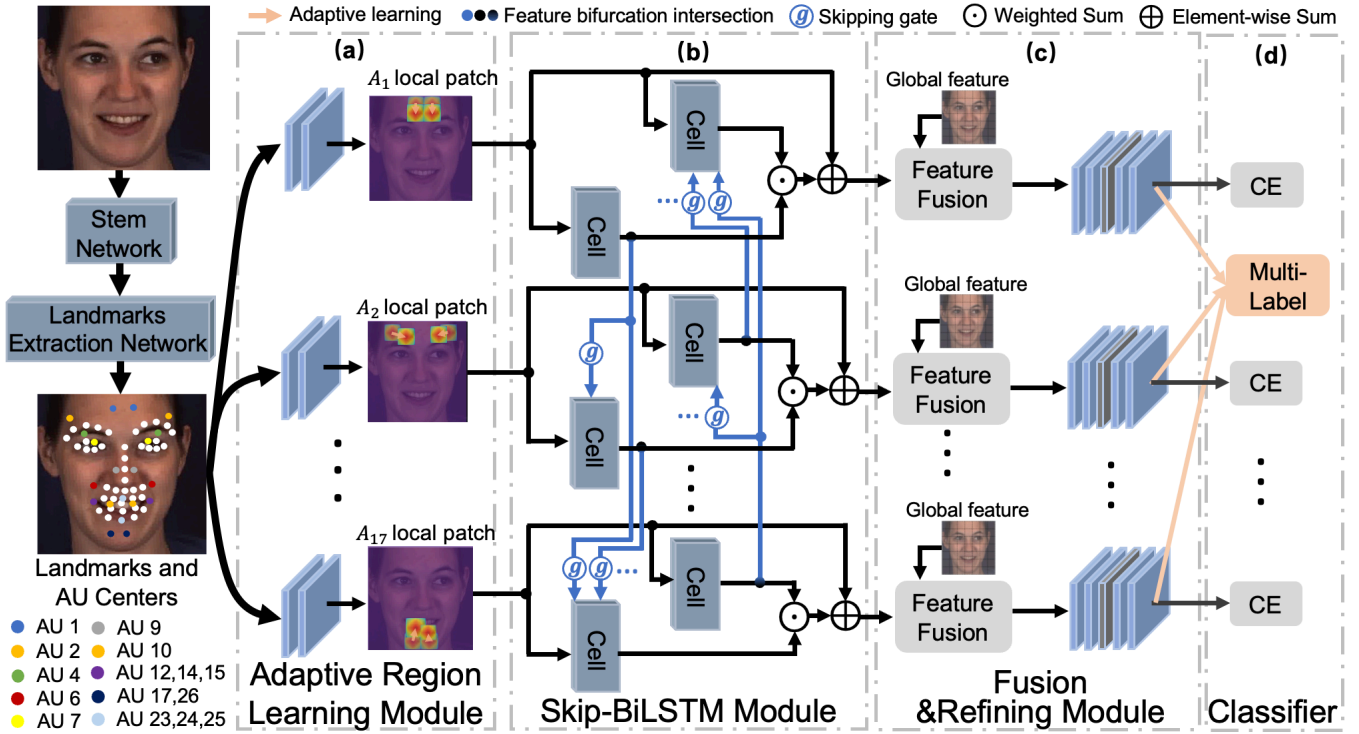


Fig. 2. The overall architecture of the proposed ALGRNet for facial AU detection. We utilize a simple landmark localization network to detect the landmarks and two linear-based network to learn the offsets and scaling factor of AU centers, which are used to compute local AU patches. We then feed the features into the novel multi-branch network with a skip-BiLSTM module and a feature fusion&refining module, with each branch corresponding to an AU. The skip-BiLSTM module mines positive and negative relations among different AU branches by different information delivery options. And the feature fusion&refining module in each branch helps the local AU region to fit irregular shape guided by the global grid-based feature. Finally, a multi-classifier is employed to predict individual AU activation probabilities.

3.4) for refining features of irregular muscle regions, and a multi-classifier module (Subsection 3.1) for predicting the AU activation probability. Note that, our ALGRNet, when applied to facial palsy detection, only modifies the definition of muscle regions according to the recommendations of the physician.

3.1 Overview of ALGRNet

Our method employs a multi-branch network [17], [18], [49] for both facial paralysis estimation and facial AU detection tasks. However, in contrast to previous methods, we believe that exploiting the relationship between multiple patches plays a crucial role in building a robust model for AU detection. In addition, due to the diversity of expressions and individual characteristics, we also attempted to learn adaptive muscle region offsets and scaling factors for each AU region. To this end, we design three modules (adaptive region learning module, skip-BiLSTM module and feature Fusion&Refining module) based on established multi-branch network that can fully exploit inter-regional and local-global interactions on the basis of adaptively adjusted muscle localization.

We first adapt a hierarchical and multi-scale region learning network from [18] as our stem network, which is used to extract the grid-based global feature and the local muscle region features. However, unlike the predefined muscle regions based on the detected landmark in [18], we add two simple networks combined with the previous

face alignment network, named adaptive region learning module (detailed in Section 3.2), to adaptively learn the offsets and scaling factors for each region. After that, local patches $A = \{A_1, A_2, \dots, A_n\}$ are computed from the learned locations and their features $V = \{v_1, v_2, \dots, v_n\}$ can be extracted through the stem network, where n is the numbers of selected patches. For the sake of simplicity, we do not repeat here the detailed structure of the stem network. The detailed structure of the stem network is not repeated here for the sake of simplicity.

In our ALGRNet, we design a novel skip-BiLSTM module (detailed in Section 3.3) to address the lack of sufficient delivery of local patch information between individual branches, which can transmit information in two ways (sequential delivery and skipping delivery) on both two directions (forward and backward), in contrast to the traditional sequence spreading of LSTM. The sequential delivery of information enables full exploration of the contextual relationships between adjacent patches. The skipping delivery highlight the interaction of information from non-adjacent related patches. After skip-BiLSTM, we get a set of local patch features $S = \{s_1, s_2, \dots, s_n\}$, which are expected to have all the useful information from adjacent and non-adjacent AU patches.

Furthermore, a novel feature Fusion&Refining module (detailed in Section 3.4) is developed to deal with irregular muscle areas, which can refine the local patches to obtain salient micro-level features for the global facial

feature G . Finally, the new patch-based representations $R = \{r_1, r_2, \dots, r_n\}$ for AUs are obtained by integrating local muscle features and global facial features.

In this work, we integrate face alignment and face AU detection (or facial paralysis estimation) into an end-to-end learning model. Our goal is to jointly learn all the parameters by minimizing both face alignment loss and facial AU detection loss (or facial paralysis estimation loss) over the training set. The face alignment loss is defined as:

$$\mathcal{L}_{align} = \frac{1}{2d_o^2} \sum_{i=1}^m [(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2], \quad (1)$$

where (x_i, y_i) and (\hat{x}_i, \hat{y}_i) denote the ground-truth (GT) coordinate and corresponding predicted coordinate of the i -th facial landmark, and d_o is the ground-truth inter-ocular distance for normalization [17].

In this paper, we also regard facial AU detection as a multi-label binary classification task following [17]. It can be formulated as a supervised classification training objective as follows,

$$\mathcal{L}_{rec} = -\frac{1}{n} \sum_{i=1}^n w_i [p_i \log \hat{p}_i + (1 - p_i) \log (1 - \hat{p}_i)], \quad (2)$$

where p_i denotes the GT probability of occurrence for the i -th AU, which is 1 if occurrence and 0 otherwise, and \hat{p}_i denotes the predicted probability of occurrence. w_i is the data balance weights, which are same as in [18]. Moreover, we also employ a weighted multi-label Dice coefficient loss [50] to overcome the sample imbalance problem, which is formulated as:

$$\mathcal{L}_{dice} = \frac{1}{n} \sum_{i=1}^n w_i \left(1 - \frac{2p_i \hat{p}_i + \tau}{p_i^2 \hat{p}_i^2 + \tau} \right), \quad (3)$$

where τ is the smoothness parameter. Finally, the facial AU detection loss is defined as:

$$\mathcal{L}_{au} = \mathcal{L}_{rec} + \mathcal{L}_{dice}, \quad (4)$$

Furthermore, we also minimize the loss of AU category classification \mathcal{L}_{int} by integrating all AUs information, including the refined AU features and the face alignment features, which is similar to the processing of \mathcal{L}_{au} . Finally, the joint loss of our ALGRNet for facial AU detection is defined as:

$$\mathcal{L} = (\mathcal{L}_{au} + \mathcal{L}_{int}) + \lambda \mathcal{L}_{align}. \quad (5)$$

where λ is a balancing parameter.

Similar with AU detection, we also joint the loss of facial paralysis estimation and face alignment, where the loss of facial paralysis estimation is formulated as:

$$\mathcal{L}_{par} = -w_i q \text{Log}(\hat{q}), \quad (6)$$

where q and \hat{q} are the label and predicted probability for the facial palsy grades, respectively. w_i is also the data balance weights, obtained by counting the different classes in the training set. Finally, we optimize the whole end-to-end network by minimizing the jointly loss function $\mathcal{L}_{par} + \lambda \mathcal{L}_{align}$ over the training set.

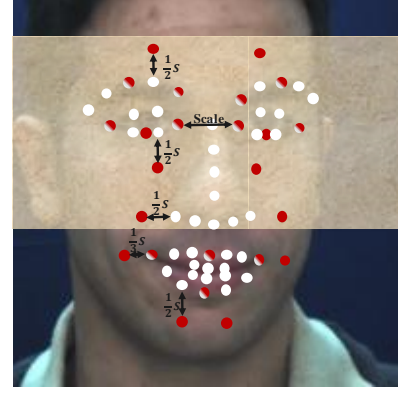


Fig. 3. New definitions for the 12 locations of muscle centers of facial paralysis estimation, which are marked in red or mixed red. The detected landmarks are marked in white or mixed red. "Scale" denotes the distance between two inner eye corners.

3.2 Adaptive Region Learning Module

Instead of the predefined muscle regions based on landmarks, we use two simple fully connected networks to adaptively learn the offsets and scaling factors for all AU regions respectively. Specially, we utilize an efficient landmark extraction network after the stem network to extract the landmarks $L = \{l_1, l_2, \dots, l_m\}$ (m is the numbers of landmarks) similar to [17], including three convolutional layers connected to a max-pooling layer. Simultaneously, two networks containing two fully-connected layers are used to get the adaptive offsets $O = \{o_1, o_2, \dots, o_{2n}\}$ and scaling factors $E = \{e_1, e_2, \dots, e_n\}$ respectively. According to the learned landmarks, offsets and scaling factors, local patches A are calculated. In particular, we first use the same rules in [17] to get the locations of AU centers based on the detected landmarks and then update the by adding the learned offsets. Please note that, we change the predefined muscle region centers, as shown in Fig. 3¹, based on the detected landmarks when we apply ALGRNet on facial paralysis estimation. Different from [17], we make the scaling factor E learnable rather than a fixed value, where e_i is the width ratio between the region of AU_i and whole feature map. After that, we generate an approximate Gaussian attention distribution for each AU region following [18]. Finally, based on the learned regions, local patch features V are extracted via the stem network.

3.3 Skip-BiLSTM

Fig. 2 (b) shows the detailed structure of our skip-BiLSTM module for contextual and skipping relationship learning. Specifically, we extract a set of local patch features $V = \{v_1, v_2, \dots, v_n\}$ from the stem network, and feed them to skip-BiLSTM. Distinct from the prior works [13], we regard the multiple patches as a sequence structure from top to bottom, which can transfer information by a Bi-directional LSTM based model [19] with our skipping-type gate. Different from the traditional BiLSTM or tree-LSTM [51], [52], our skip-BiLSTM can directly calculate the correlation between

1. Due to patient confidentiality agreements, we cannot show real patients with facial palsy. This example image is from BP4D.

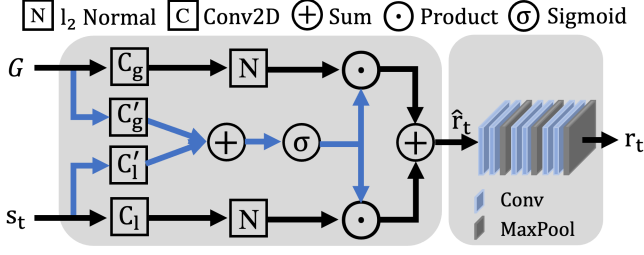


Fig. 4. The architecture of our feature fusion&refining module guided by global face feature.

a target AU and all other AUs. For the t -th patch ($t > 1$), the extracted feature v_t is used to learn the weights with forward hidden states $H = \{h_1, \dots, h_{t-1}\}$ by the skipping-type gates, which can determine the correlation coefficient between past AUs and current AU. And then the new states $\hat{H} = \{\hat{h}_1, \dots, \hat{h}_{t-1}\}$ and v_t are fed into the t -th forward cell in the skip-BiLSTM to learn the association weights, which can promote the transfer of relevant AUs information. The above process can be formulated as:

$$\vec{h}_t = \text{Cell}(\sum_{j=1}^{t-1} \vec{h}_j, v_t), \quad (7)$$

$$\vec{h}_j = \vec{h}_j f_j, \quad (8)$$

$$f_j = \sigma(\text{GAP}(W_j(\vec{h}_j v_t))), \quad (9)$$

where $\text{Cell}(\cdot)$ indicates the basic ConvLstm cell [53], and GAP denotes the global average pooling operation. W_j is the parameters of mapping function, in which we used Conv2D. σ denotes sigmoid function. We obtain the t -th patch feature for backward delivery, which follows the identical forward method as:

$$\overleftarrow{h}_t = \text{Cell}(\sum_{j=t+1}^n \overleftarrow{h}_j, v_t), \quad (10)$$

In order to fully promote the information interactive among individual AUs, the final representation for each patch is computed as the average of the hidden vectors in both directions, as well as the original patch feature:

$$s_t = v_t + (\vec{h}_t + \overleftarrow{h}_t)/2, \quad (11)$$

3.4 Feature Fusion&Refining Module

To exploit the useful global face feature, we design a gated fusion architecture and a refining architecture (F&R) that can selectively balance the relative importance of local patches and global face grids. We add these two architectures on each local AU branch because different AUs may focus on different global information. The grid-based global face feature G is extracted using a simple CNN with the same structure as the face alignment network [17]. As shown in Fig. 4, after obtaining the learned t -th local patch feature, it

TABLE 1
Overview information of our collected facial paralysis dataset.

Grade	Normal	Low	Medium	High
Num. of Video	20	29	20	20
Num. of Frame	9049	16970	11019	10547

is fused with the grid-based global feature G by the fusion architecture, which can be formulated as:

$$\alpha = \sigma(C'_g G + C'_l s_t), \quad (12)$$

$$\hat{r}_t = \alpha \odot \|C_g G\|_2 \oplus (1 - \alpha) \odot \|C_l s_t\|_2, \quad (13)$$

where σ is the sigmoid function, and $\|\cdot\|$ denotes the l_2 -normalization. C'_* and C_* denote the Conv2D operation. \oplus denotes the element-wise weighted sum of $\|C_g G\|_2$ and $\|C_l s_t\|_2$ according to the learned gate vector α .

The final local fusion feature s_t for t -th patch refined by our F&R module is shown in Fig. 4. F&R module contains three blocks. Each block consists of two convolutional layers and a maxpooling layer. Then multi-patch features R are sent to the multi-label binary classifier to calculate the occurrence probabilities of individual AUs.

4 EXPERIMENTS

4.1 Dataset

We evaluate the effectiveness of the proposed approach for facial AU detection on popular BP4D [56] and DISFA [57] datasets. **BP4D** consists of 328 facial videos from 41 participants (23 females and 18 males) who were involved in 8 sessions. Similar to [16], [17], [25], we consider 12 AUs and 140K valid frames with labels. **DISFA** consists of 27 participants (12 females and 15 males). Each participant has a video of 4,845 frames. We also limited the number of AUs to 8 similar to [17], [25]. In comparison to BP4D, the experimental protocol and lighting conditions deliver DISFA to be a more challenging dataset. Following the experiment setting of [17], we evaluated the model using the 3-fold subject-exclusive cross-validation protocol.

To evaluate the effectiveness of our ALGRNet for facial palsy severity estimation, we exploited a facial paralysis dataset from NHS, named **FPara** (the details in Table 1), which consists of 89 videos of facial palsy patients performing various types of facial palsy exercises inline with the House-Brackmann (H-B) scale [45]. Each of the videos consisted of facial palsy patients performing a set of exercises, such as raise eyebrows, close eyes gently, close eyes tightly, scrunch up face and smile, *etc.* They were part of a previous study on facial palsy with patient consent for research [58]. These videos are assigned a H-B scale from 1 to 6, and 1 being normal and 6 being severest with no body movements. We then further split into four grades, such as normal (H-B score=1), low (H-B score=2), medium ($3 \leq \text{H-B score} \leq 4$) and high ($5 \leq \text{H-B score} \leq 6$) grades. FPara data is summarised in Table 1. Similar to the settings of facial AU dataset, all facial paralysis grades are evaluated using subject exclusive 3-fold cross-validation, where two folds (about 80%) are used for training and the remaining one is used for testing (about 20%).

TABLE 2

Performance comparisons on F1-frame score of diverse AU detection for 12 AUs on BP4D. All values are in %. * means the method employed pretrained model on additional dataset, such as ImageNet and VGGFace2, *etc.*, so we do not compare. The first and second places are marked with the bold font and underline, respectively.

Method	AU Index												Avg.
	1	2	4	6	7	10	12	14	15	17	23	24	
DSIN [49]	<u>51.7</u>	40.4	56.0	76.1	73.5	79.9	85.4	62.7	37.3	62.8	38.8	41.6	58.9
MLCR [38]	42.4	36.9	48.1	77.5	77.6	83.6	85.8	61.0	43.7	63.2	42.1	55.6	59.8
CMS [54]	49.1	44.1	50.3	79.2	74.7	80.9	88.2	63.9	44.4	60.3	41.4	51.2	60.6
LP-Net [13]	46.9	45.3	55.6	77.1	76.7	83.8	87.2	63.3	45.3	60.5	48.1	<u>54.2</u>	61.0
JAA-Net [18]	47.2	44.0	54.9	77.5	74.6	<u>84.0</u>	86.9	61.9	43.6	60.3	42.7	41.9	60.0
ARL [16]	45.8	39.8	55.1	75.7	<u>77.2</u>	82.3	86.6	58.8	<u>47.6</u>	62.1	47.4	55.4	61.1
JAA-Net [17]	53.8	<u>47.8</u>	58.2	<u>78.5</u>	75.8	82.7	<u>88.2</u>	<u>63.7</u>	43.3	61.8	45.6	49.9	<u>62.4</u>
UGN-B* [41]	54.2	46.4	56.8	76.2	76.7	82.4	86.1	64.7	51.2	63.1	48.5	53.6	63.3
HMP-PS* [36]	53.1	46.1	56.0	76.5	76.9	82.1	86.4	64.8	51.5	63.0	49.9	54.5	63.4
DML* [55]	52.6	44.9	56.2	79.8	80.4	85.2	88.3	65.6	51.7	59.4	47.3	49.2	63.4
ALGRNet (Ours)	51.2	48.2	<u>57.3</u>	77.9	76.4	84.9	88.2	64.8	50.8	<u>62.8</u>	<u>47.6</u>	51.9	63.5

TABLE 3

Performance comparisons on F1-frame score of diverse AU detection for 8 AUs on DISFA. All values are in %. The first and second places are marked with the bold font and underline, respectively.

Method	AU Index								Avg.
	1	2	4	6	9	12	25	26	
DSIN [49]	42.4	39.0	68.4	28.6	46.8	70.8	90.4	42.2	53.6
CMS [54]	40.2	44.3	53.2	57.1	<u>50.3</u>	73.5	81.1	59.7	57.4
LP-Net [13]	29.9	24.7	<u>72.7</u>	<u>46.8</u>	49.6	72.9	93.8	65.0	56.9
JAA-Net [18]	43.7	46.2	56.0	41.4	44.7	69.6	88.3	58.4	56.0
ARL [16]	43.9	42.1	63.6	41.8	40.0	76.2	95.2	66.8	58.7
JÂA-Net [17]	<u>62.4</u>	<u>60.7</u>	67.1	41.1	45.1	73.5	90.9	<u>67.4</u>	<u>63.5</u>
UGN-B* [41]	43.3	48.1	63.4	49.5	48.2	72.9	90.8	59.0	60.0
HMP-PS* [36]	21.8	48.5	53.6	56.0	58.7	57.4	55.9	56.9	61.0
DML* [55]	62.9	65.8	71.3	51.4	45.9	76.0	92.1	50.2	64.4
ALGRNet	63.8	65.4	73.6	44.5	54.1	<u>74.0</u>	<u>94.7</u>	69.9	67.5

4.2 Implementation Detail

Our model is trained on a single NVIDIA Tesla V100 GPU with 32 GB memory. The whole network is trained using PyTorch [59] with the stochastic gradient descent (SGD) solver, a Nesterov momentum [60] of 0.9 and a weight decay of 0.0005. The learning rate is set to 0.01 initially with a decay rate of 0.5 every 2 epochs. Maximum epoch number is set to 20. To enhance the diversity of training data, aligned faces are further randomly cropped into 176×176 and horizontally flipped. Regarding face alignment network and stem network, we set the value of the general parameters to be the same with [17]. The filters for the convolutional layers in refining architecture are used 3×3 convolutional filters with a stride 1 and a padding 1. In our paper, all of the mapping Conv2D operations are used 1×1 convolutional filters with a stride 1 and a padding 1. The dimensionality of hidden state in ConvLstm cell is set to 64. The filters for

the convolutional layers in ConvLstm cell are the same as refining architecture. λ is set to 0.5 for the jointly optimizing of AU detection and face alignment. The ground-truth annotations of 49 landmarks of training data is detected by SDM [61]. Different from JAA-Net [17], we averaged the predicted probability of the local information and the integrated information as the final predicted activation probability for each AU, rather than simply using the integrated information of all the AUs.

4.3 Performance Metric.

We evaluate the performance of all methods in terms of the F1 score (%) which has been widely used for classification. F1-frame score is the harmonic mean of the Precision P and Recall R and calculated by $F1 = 2PR/(P + R)$. For comparison, we calculate F1 score for all facial paralysis grades on FPara and for all the AUs on DISFA and BP4D

TABLE 4
Ablation study of ALGRNet for 8 AUs on DISFA. All values are in %.

Methods	Setting			AU Index								Avg.
	S-B	F&R	Ada	1	2	4	6	9	12	25	26	
w/o full				47.1	61.1	66.3	<u>44.7</u>	52.2	74.9	92.2	66.2	63.1
w/o F&R	✓			62.6	64.2	72.4	42.3	49.9	76.1	93.5	<u>72.6</u>	<u>66.7</u>
w/o S-B		✓		58.7	<u>65.2</u>	<u>73.5</u>	43.9	<u>53.5</u>	72.2	94.1	64.7	65.7
w/ Bi		✓		61.1	58.4	70.9	45.5	47.9	74.9	92.5	70.8	65.2
w/o Ada	✓	✓		<u>62.6</u>	64.4	72.5	46.6	48.8	<u>75.7</u>	<u>94.4</u>	73.0	67.3
ALGRNet	✓	✓	✓	63.8	65.4	73.6	44.5	54.1	74.0	94.7	69.9	67.5

TABLE 5
Ablation study of ALGRNet for 12 AUs on BP4D. All values are in %.

Methods	Setting			AU Index												Avg.
	S-B	F&R	Ada	1	2	4	6	7	10	12	14	15	17	23	24	
w/o full				50.1	47.1	54.3	77.3	75.1	82.5	88.1	61.7	44.9	62.7	45.2	49.9	61.6
w/o F&R	✓			50.4	46.9	53.4	79.0	<u>77.4</u>	<u>84.7</u>	87.4	63.0	45.3	63.3	47.0	55.7	62.8
w/o S-B		✓		51.3	47.6	56.3	<u>78.2</u>	76.2	83.7	88.1	64.4	49.1	61.9	46.1	49.8	62.7
w/ Bi		✓		50.7	50.0	55.2	77.0	75.7	84.1	<u>88.2</u>	63.4	49.1	62.3	<u>47.3</u>	52.0	62.9
w/o Ada	✓	✓		50.8	47.1	57.8	77.6	77.4	<u>84.7</u>	88.2	66.4	<u>49.8</u>	61.5	46.8	<u>52.3</u>	<u>63.4</u>
ALGRNet	✓	✓	✓	<u>51.2</u>	<u>48.2</u>	<u>57.3</u>	77.9	76.4	84.9	88.2	<u>64.8</u>	50.8	<u>62.8</u>	47.6	51.9	63.5

and then average them (denoted as **Avg.**) separately with “%” omitted.

4.4 Overall Performance of AU Detection

We compare the proposed ALGRNet for AU detection with several single-image based baselines in Table 2 and Table 3, including Deep Structure Inference Network (DSIN) [49], Joint AU Detection and Face Alignment (JAA-Net) [18], Multi-Label Co-Regularization (MLCR) [38], Cross Modality Supervision (CMS) [54], Local relationship learning with Person-specific shape regularization (LP-Net) [13], Attention and Relation Learning (ARL) [16], and Joint AU detection and face alignment via Adaptive Attention Network (JAA-Net) [17]. The performances of the baselines in Table 3 and 2 are their reported results. The first and second places are marked with the bold font and “_”, respectively.

For a more comprehensive display, we also show methods (marked with *) [36], [41], [55] that use additional data, such as ImageNet [62] and VGGFace2 [63], *etc.*, for pre-training. Due to the fact that our stem network only consists of a few simple convolutional layers, even if we pre-trained on additional datasets, it is unfair compared to pre-training on deeper feature extraction networks, such as ResNet50 [64]. In fact, our results are still excellent compared with them, which demonstrates the superiority and effectiveness of our proposed learning scheme. We omit the need for additional modal inputs and non-frame-based models [65], [66].

Quantitative comparison on BP4D: We report the performance comparisons between our ALGRNet and baselines on BP4D in Table 2. As it can be observed, our ALGRNet

significantly outperforms all the other methods in terms of F1-frame score and achieves the first and second places for most of the 12 AUs annotated in BP4D. JAA-Net is the latest state-of-the-art method which also joint AU detection and face alignment into an end-to-end multi-label multi-branch network. Our ALGRNet achieves 1.1% higher average F1-frame score compared with JAA-Net. The main reason lies in our ALGRNet overcomes the problem of non-transferable information between branches in the JAA-Net and adaptively adjusts the muscle regions corresponding to the AU.

Quantitative comparison on DISFA: We also report the performance of our proposed ALGRNet on DISFA. Table 3 shows the performance of our ALGRNet is the best in terms of average F1 score compared with all baselines. And our approach significantly outperforms all other methods for most of the 8 AUs annotated in DISFA. Compared with the existing end-to-end feature learning and multi-label classification methods DSIN [49] and ARL [16], the average F1-frame score of our proposed ALGRNet get 13.9% and 8.8% higher, respectively. Moreover, compared with the multi-branch combination-based state-of-the-art method JAA-Net [17], our ALGRNet achieves 4.0% improvements in terms of average F1-frame score.

The eventual experimental results of ALGRNet demonstrate its effectiveness in improving AU detection accuracy on DISFA and BP4D, as well as good robustness.

4.5 Ablative Analysis

To fully examine the impact of our proposed adaptive region learning module, skip-BiLSTM module and feature

TABLE 6
Mean error (lower is better) results of different face alignment models on BP4D, DISFA and FPara. All values are in %.

Methods	BP4D	DISFA	FPara
JAA-Net	3.80	3.87	5.15
ALGRNet	3.78	3.29	5.18

fusion&refining module, we conduct detailed ablative studies to compare different variants of ALGRNet for facial AU detection on DISFA and BP4D.

4.5.1 Effects of adaptive region learning module

To cancel out the adaptive region learning (indicated w/o Ada), we follow the same experiment setting as [17] (It means each scaling factor e is set to 0.14.) to predefined muscle region based on the detected landmarks for each AU. In Table 4 and 5, ALGRNet decreases its F1 score to 67.3% and 63.4% on DISFA and BP4D respectively. It has been shown that the adaptive region learning module can contribute the AU detection greatly.

4.5.2 Effects of skip-BiLSTM

In Table 4 and 5, when the skip-BiLSTM module is removed (indicated by w/o S-B), ALGRNet (without adaptive region learning module) shows an absolute decrease of 1.6% and 0.7% in the average F1-frame score for DISFA and BP4D, respectively. In addition, in order to explicitly validate the effectiveness of our skipping operation, we use the basic BiLSTM [19] (indicated by w/ Bi) instead of skip-BiLSTM for information sequential transfer across different branches in the ALGRNet (also with Fusion&Refining module), ALGRNet obtains lower average F1 scores of 65.2% and 62.9% on DISFA and BP4D, respectively. The performance reduction clearly verifies that roughly defining the relationships between AU-related branches from top to bottom may not be the best way to model the real relationships between AUs. Notably, skipping operation can significantly improve performance, suggesting that our skip-type gates play an important role in our model. Furthermore, the eventual experimental results demonstrate the effectiveness of our skip-BiLSTM in modelling mutual assistance and exclusion relationship between different patch-based branches for AU detection.

4.5.3 Effects of feature fusion&refining module

In order to illustrate the effectiveness of feature fusion&refining module, we directly conduct the classification over the output of skip-BiLSTM without fusion&refining module (indicated by w/o F&R in Table 4 and 5). When the fusion&refining module is not used, the average F1-frame score drops significantly from 67.3% to 66.7% on DISFA and from 63.4% to 62.8% on BP4D, due to the lack of supplementary information from the global face for each patch. This suggests that the refined local AU features from the proposed fusion&refining module, guided by the grid-based global features, make a significant contribution in our model.

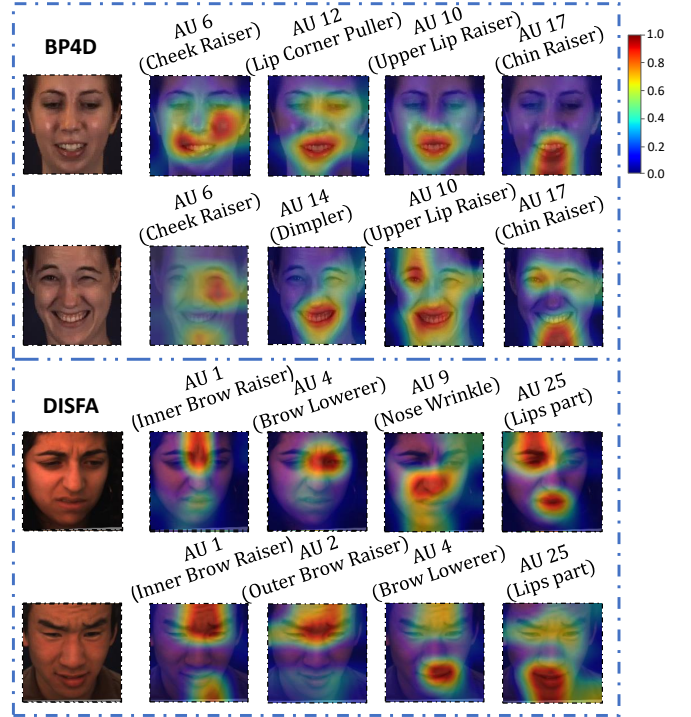


Fig. 5. Class activation maps that show the discriminative regions for different AUs in terms of different expressions and individuals on DISFA and BP4D datasets.

Finally, after simultaneously removing all the proposed adaptive region learning module, skip-BiLSTM and fusion&refining module (marked by w/o full in Table 4 and 5), a significant performance degradation in AU detection can be observed, *i.e.*, a 4.4% drop on DISFA and a 1.9% drop on BP4D in terms of average F1-frame score. This sufficiently demonstrates that the potential mutual assistance and exclusion relationships between the adaptive AU patches, complemented by the global facial features, can significantly improve the performance of facial AU detection.

4.6 Results for Face Alignment

We integrate face alignment and face AU detection into our end-to-end ALGRNet, which can be beneficial for each other as they are coherently related to each other. For example, the detected landmarks can help the model focus on the exact muscle areas of the AU patches. As shown in Table 6, compared with baseline method JAA-Net [17], our ALGRNet performs comparably to baseline on FPara and better on BP4D and DISFA. The robustness of the adaptive region learning module allows our ALGRNet to outperform JAA-Net in AU detection and facial paralysis estimation, even if sometimes with slightly lower landmark detection accuracy.

4.7 Visualization of Results

Fig. 5 shows some examples of the learned class activation maps of ALGRNet (the outputs of F&R module) corresponding to different AUs. For an adequate display, we show two individuals from BP4D and DISFA respectively, containing visualisations of different genders with different AU categories. Through the learning of ALGRNet, not only

TABLE 7
Performance comparisons on F1-frame score (in %) of diverse facial paralysis estimation for 14 grades on FPara.

Method	Facial Paralysis Grades				Avg.
	Normal	Low	Medium	High	
ResNet18 [64]	99.8	50.7	47.7	67.9	66.5
ResNet50 [64]	99.9	53.9	54.7	71.4	70.0
Transformer-based [67]	100	63.0	58.6	68.7	72.6
JAA-Net [17]	<u>100</u>	<u>58.8</u>	<u>64.3</u>	<u>72.9</u>	<u>72.8</u>
ALGRNet(Ours)	100	55.9	72.1	73.2	75.4

the concerned AU regions can be accurately located, but also the positive (in red) or negative (in blue) correlation with other AU areas can be established and other details of the global face can be supplemented. This obviously improves the flaws of the excessive localisation of JAA-Net [17] and the negative influence of unrelated regions of ARL [18]. In addition, it also adapts well to irregular muscle areas for different AUs. The heatmaps for the same AU category in the different examples are broadly consistent but also vary slightly by individual, demonstrating that our ALGRNet can learn certain rules across different datasets and adaptively adjust to different samples.

4.8 Facial Palsy Severity Estimation

In this section we evaluate the effectiveness of facial palsy severity estimation.

4.8.1 Facial paralysis estimation

Different from facial AU detection, the exiting deep-learning-based facial paralysis estimation methods are rare, so we apply currently popular deep learning classification methods, such as the ResNet [64] and Transformer [68], on our collected facial palsy dataset (FPara). Besides, we also compare it with the state-of-the-art AU detection approach, JAA-Net [17]. Specially, we evaluate the following methods:

- ResNet18 and ResNet50 [64]: These methods use different depth layers based on ResNet to model the input face images, which are similar to [69].
- Transformer-based method [67]: This baseline is motivated from self-attention and uses the Transformer [68] architecture. The output of the Transformer-based encoder [67] is treated as the latent representation for the input of the multi-label AU classifier.
- JAA-Net [17]: This is a recently proposed multi-branch combination-based AU detection method, which can extract precise local muscle features thanks to a joint facial alignment network.

4.8.2 Quantitative comparison on the collected FPara

Facial paralysis estimation results by different methods on our FPara are shown in Table 7. It has been shown that our ALGRNet outperforms all its competitors with impressive margins. Specifically, compared with the state-of-the-art AU detection method JAA-Net [17], our ALGRNet achieves 2.6% improvements in terms of average F1 score. Moreover,

the average F1 score of our ALGRNet get 2.8% higher compared to the currently popular Transformer-based approach [67].

The eventual experimental results of our ALGRNet demonstrate that it is successful in boosting AU detection accuracy on BP4D and DISFA, as well as having high generalisation ability on our new facial paralysis dataset.

5 CONCLUSION

In this paper, we present ALGRNet, a novel adaptive local-global relational network for both facial action units detection and facial paralysis severity estimation, which can exploit the mutual assistance and exclusion relationships of adaptive muscle regions as well as interactions with global information. Specifically, ALGRNet adaptively represents the corresponding muscle areas in terms of different expressions and individual characteristics. It then models the potential mutual assistance and exclusion relationships between AU branches and enables efficient information transfer via a novel skip-BiLSTM to get local features. Finally, a novel feature fusion&refining module is proposed and equipped in each branch, facilitating complementarity between local feature and grid-based global feature as well as adaptation to irregular muscle regions. Experimental results on two widely used AU detection benchmarks show the effectiveness of the AU detection and the superiority of the Facial Palsy severity estimation on a facial paralysis estimation benchmark. This not only witnesses the distinct performance gains over the state-of-the-arts, but also demonstrates effective migration capability from AU detection to the facial paralysis estimation task.

REFERENCES

- [1] D. R. Rubinow and R. M. Post, "Impaired recognition of affect in facial expression in depressed patients," *Biological Psychiatry*, vol. 31, no. 9, pp. 947–953, 1992.
- [2] B. Klare and A. K. Jain, "Heterogeneous face recognition: Matching nir to visible light images," in *International Conference on Pattern Recognition*, 2010, pp. 1513–1516.
- [3] X. Niu, H. Han, J. Zeng, X. Sun, S. Shan, Y. Huang, S. Yang, and X. Chen, "Automatic engagement prediction with gap feature," in *ACM International Conference on Multimodal Interaction*, 2018, pp. 599–603.
- [4] R. S. Feldman, L. Jenkins, and O. Popoola, "Detection of deception in adults and children via facial expressions," *Child Development*, pp. 350–355, 1979.
- [5] P. Ekman and E. L. Rosenberg, *What the face reveals: Basic and applied studies of spontaneous expression using the Facial Action Coding System (FACS)*. Oxford University Press, USA, 1997.
- [6] Y. Tong and Q. Ji, "Learning bayesian networks with qualitative constraints," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2008, pp. 1–8.
- [7] Y. Li, S. Wang, Y. Zhao, and Q. Ji, "Simultaneous facial feature tracking and facial expression recognition," *IEEE Transactions on Image Processing*, vol. 22, no. 7, pp. 2559–2573, 2013.
- [8] G. Cheng, J. Dong, S. Wang, H. Qu *et al.*, "Evaluation of facial paralysis degree based on regions," in *Third International Conference on Knowledge Discovery and Data Mining*. IEEE, 2010, pp. 514–517.
- [9] G. Barrios Dell'Olio and M. Sra, "Farapy: An augmented reality feedback system for facial paralysis using action unit intensity estimation," in *The 34th Annual ACM Symposium on User Interface Software and Technology*, 2021, pp. 1027–1038.
- [10] Z. Chen, R. Ansari, and D. Wilkie, "Learning pain from action unit combinations: a weakly supervised approach via multiple instance learning," *IEEE Transactions on Affective Computing*, 2019.

- [11] P. Liu, S. Han, Z. Meng, and Y. Tong, "Facial expression recognition via a boosted deep belief network," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1805–1812.
- [12] P. Liu, J. T. Zhou, I. W.-H. Tsang, Z. Meng, S. Han, and Y. Tong, "Feature disentangling machine—a novel approach of feature selection and disentangling in facial expression analysis," in *European Conference on Computer Vision*, 2014, pp. 151–166.
- [13] X. Niu, H. Han, S. Yang, Y. Huang, and S. Shan, "Local relationship learning with person-specific shape regularization for facial action unit detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 917–11 926.
- [14] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [15] K. Zhao, W.-S. Chu, F. De la Torre, J. F. Cohn, and H. Zhang, "Joint patch and multi-label learning for facial action unit and holistic expression recognition," *IEEE Transactions on Image Processing*, vol. 25, no. 8, pp. 3931–3946, 2016.
- [16] Z. Shao, Z. Liu, J. Cai, Y. Wu, and L. Ma, "Facial action unit detection using attention and relation learning," *IEEE transactions on Affective Computing*, 2019.
- [17] Z. Shao, Z. Liu, J. Cai, and L. Ma, "Jaa-net: joint facial action unit detection and face alignment via adaptive attention," *International Journal of Computer Vision*, vol. 129, no. 2, pp. 321–340, 2021.
- [18] —, "Deep adaptive attention for joint facial action unit detection and face alignment," in *European Conference on Computer Vision*, 2018, pp. 705–720.
- [19] A. Graves and J. Schmidhuber, "Framewise phoneme classification with bidirectional lstm and other neural network architectures," *Neural Networks*, vol. 18, no. 5–6, pp. 602–610, 2005.
- [20] X. Ge, P. Wan, H. Han, J. M. Jose, Z. Ji, Z. Wu, and X. Liu, "Local global relational network for facial action units recognition," in *IEEE International Conference on Automatic Face and Gesture Recognition*, 2021, pp. 01–08.
- [21] Y. Zhu, F. De la Torre, J. F. Cohn, and Y.-J. Zhang, "Dynamic cascades with bidirectional bootstrapping for action unit detection in spontaneous facial behavior," *IEEE Transactions on Affective Computing*, vol. 2, no. 2, pp. 79–91, 2011.
- [22] Z. Wang, Y. Li, S. Wang, and Q. Ji, "Capturing global semantic relationships for facial action unit recognition," in *IEEE International Conference on Computer Vision*, 2013, pp. 3304–3311.
- [23] Y. Li, J. Chen, Y. Zhao, and Q. Ji, "Data-free prior model for facial action unit recognition," *IEEE Transactions on Affective Computing*, vol. 4, no. 2, pp. 127–141, 2013.
- [24] K. Zhao, W.-S. Chu, F. De la Torre, J. F. Cohn, and H. Zhang, "Joint patch and multi-label learning for facial action unit detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 2207–2216.
- [25] W. Li, F. Abtahi, and Z. Zhu, "Action unit detection with region adaptation, multi-labeling learning and optimal temporal fusing," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1841–1850.
- [26] C. Tang, W. Zheng, J. Yan, Q. Li, Y. Li, T. Zhang, and Z. Cui, "View-independent facial action unit detection," in *IEEE International Conference on Automatic Face & Gesture Recognition*. IEEE, 2017, pp. 878–882.
- [27] C. Ma, L. Chen, and J. Yong, "Au r-cnn: Encoding expert prior knowledge into r-cnn for action unit detection," *Neurocomputing*, vol. 355, pp. 35–47, 2019.
- [28] I. Ntinou, E. Sanchez, A. Bulat, M. Valstar, and Y. Tzimiropoulos, "A transfer learning approach to heatmap regression for action unit intensity estimation," *IEEE Transactions on Affective Computing*, 2021.
- [29] S. Taheri, Q. Qiu, and R. Chellappa, "Structure-preserving sparse decomposition for facial expression analysis," *IEEE Transactions on Image Processing*, vol. 23, no. 8, pp. 3590–3603, 2014.
- [30] S. Jaiswal and M. Valstar, "Deep learning the dynamic appearance and shape of facial action units," in *IEEE Winter Conference on Applications of Computer Vision*, 2016, pp. 1–8.
- [31] Y. Tong, W. Liao, and Q. Ji, "Facial action unit recognition by exploiting their dynamic and semantic relationships," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 29, no. 10, pp. 1683–1699, 2007.
- [32] S. Wu, S. Wang, B. Pan, and Q. Ji, "Deep facial action unit recognition from partially labeled data," in *IEEE International Conference on Computer Vision*, 2017, pp. 3951–3959.
- [33] S. Wang, S. Wu, G. Peng, and Q. Ji, "Capturing feature and label relations simultaneously for multiple facial action unit recognition," *IEEE Transactions on Affective Computing*, vol. 10, no. 3, pp. 348–359, 2017.
- [34] T. Song, W. Zheng, P. Song, and Z. Cui, "Eeg emotion recognition using dynamical graph convolutional neural networks," *IEEE Transactions on Affective Computing*, vol. 11, no. 3, pp. 532–541, 2018.
- [35] S. Wang, G. Peng, and Q. Ji, "Exploring domain knowledge for facial expression-assisted action unit activation recognition," *IEEE Transactions on Affective Computing*, vol. 11, no. 4, pp. 640–652, 2018.
- [36] T. Song, Z. Cui, W. Zheng, and Q. Ji, "Hybrid message passing with performance-driven structures for facial action unit detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 6267–6276.
- [37] K. Zhao, W.-S. Chu, and H. Zhang, "Deep region and multi-label learning for facial action unit detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 3391–3399.
- [38] X. Niu, H. Han, S. Shan, and X. Chen, "Multi-label co-regularization for semi-supervised facial action unit recognition," in *Advances in Neural Information Processing Systems*, 2019, pp. 909–919.
- [39] G. Li, X. Zhu, Y. Zeng, Q. Wang, and L. Lin, "Semantic relationships guided representation learning for facial action unit recognition," in *AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 8594–8601.
- [40] T. Song, Z. Cui, Y. Wang, W. Zheng, and Q. Ji, "Dynamic probabilistic graph convolution for facial action unit intensity estimation," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4845–4854.
- [41] T. Song, L. Chen, W. Zheng, and Q. Ji, "Uncertain graph neural networks for facial action unit detection," in *AAAI Conference on Artificial Intelligence*, 2021, p. 5993–6001.
- [42] M. J. Fields and N. S. Peckitt, "Facial nerve function index: a clinical measurement of facial nerve activity in patients with facial nerve palsies," *Oral surgery, oral medicine, oral pathology*, vol. 69, no. 6, pp. 681–682, 1990.
- [43] E. A. Chu, T. Y. Farrag, L. E. Ishii, and P. J. Byrne, "Threshold of visual perception of facial asymmetry in a facial paralysis model," *Archives of Facial Plastic Surgery*, vol. 13, no. 1, pp. 14–19, 2011.
- [44] G. E. Murty, J. P. Diver, P. J. Kelly, G. O'Donoghue, and P. J. Bradley, "The nottingham system: objective assessment of facial nerve function in the clinic," *Otolaryngology—Head and Neck Surgery*, vol. 110, no. 2, pp. 156–161, 1994.
- [45] W. House, "Facial nerve grading system," *Otolaryngol Head Neck Surg*, vol. 93, pp. 184–193, 1985.
- [46] J. Dong, L. Ma, Q. Li, S. Wang, L.-a. Liu, Y. Lin, and M. Jian, "An approach for quantitative evaluation of the degree of facial paralysis based on salient point detection," in *International Symposium on Intelligent Information Technology Application Workshops*. IEEE, 2008, pp. 483–486.
- [47] J. Barbosa, K. Lee, S. Lee, B. Lodhi, J.-G. Cho, W.-K. Seo, and J. Kang, "Efficient quantitative assessment of facial paralysis using iris segmentation and active contour-based key points detection with hybrid classifier," *BMC medical imaging*, vol. 16, no. 1, pp. 1–18, 2016.
- [48] T. Wang, S. Zhang, J. Dong, L. Liu, and H. Yu, "Automatic evaluation of the degree of facial nerve paralysis," *Multimedia Tools and Applications*, vol. 75, no. 19, pp. 11 893–11 908, 2016.
- [49] C. Corneanu, M. Madadi, and S. Escalera, "Deep structure inference network for facial action unit recognition," in *European Conference on Computer Vision*, 2018, pp. 298–313.
- [50] F. Milletari, N. Navab, and S.-A. Ahmadi, "V-net: Fully convolutional neural networks for volumetric medical image segmentation," in *The Fourth International Conference on 3D Vision*, 2016, pp. 565–571.
- [51] K. S. Tai, R. Socher, and C. D. Manning, "Improved semantic representations from tree-structured long short-term memory networks," in *The Association for Computer Linguistics*, 2015, pp. 1556–1566.
- [52] X. Ge, F. Chen, J. M. Jose, Z. Ji, Z. Wu, and X. Liu, "Structured multi-modal feature embedding and alignment for image-sentence retrieval," in *ACM International Conference on Multimedia*, 2021, pp. 5185–5193.
- [53] X. Shi, Z. Chen, H. Wang, D.-Y. Yeung, W.-K. Wong, and W.-c. Woo, "Convolutional lstm network: A machine learning approach for precipitation nowcasting," *Advances in Neural Information Processing Systems*, pp. 802–810, 2015.
- [54] N. Sankaran, D. D. Mohan, S. Setlur, V. Govindaraju, and D. Fedorishin, "Representation learning through cross-modality super-

vision," in *IEEE International Conference on Automatic Face & Gesture Recognition*, 2019, pp. 1–8.

- [55] S. Wang, Y. Chang, and C. Wang, "Dual learning for joint facial landmark detection and action unit recognition," *IEEE Transactions on Affective Computing*, 2021.
- [56] X. Zhang, L. Yin, J. F. Cohn, S. Canavan, M. Reale, A. Horowitz, P. Liu, and J. M. Girard, "Bp4d-spontaneous: a high-resolution spontaneous 3d dynamic facial expression database," *Image and Vision Computing*, vol. 32, no. 10, pp. 692–706, 2014.
- [57] S. M. Mavadati, M. H. Mahoor, K. Bartlett, P. Trinh, and J. F. Cohn, "Disfa: A spontaneous facial action intensity database," *IEEE Transactions on Affective Computing*, vol. 4, no. 2, pp. 151–160, 2013.
- [58] B. F. O'Reilly, J. J. Soraghan, S. McGrenary, and S. He, "Objective method of assessing and presenting the house-brackmann and regional grades of facial palsy by production of a facogram," *Otology & Neurotology*, vol. 31, no. 3, pp. 486–491, 2010.
- [59] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems*, 2019, pp. 8026–8037.
- [60] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, "On the importance of initialization and momentum in deep learning," in *International Conference on Machine Learning*, 2013, pp. 1139–1147.
- [61] X. Xiong and F. De la Torre, "Supervised descent method and its applications to face alignment," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 532–539.
- [62] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.
- [63] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, "Vg-face2: A dataset for recognising faces across pose and age," in *IEEE International Conference on Automatic Face & Gesture Recognition*, 2018, pp. 67–74.
- [64] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [65] H. Yang, T. Wang, and L. Yin, "Adaptive multimodal fusion for facial action units recognition," in *ACM International Conference on Multimedia*, 2020, pp. 2982–2990.
- [66] P. Liu, Z. Zhang, H. Yang, and L. Yin, "Multi-modality empowered network for facial action unit detection," in *IEEE Winter Conference on Applications of Computer Vision*, 2019, pp. 2175–2184.
- [67] J. M. J. Valanarasu, P. Oza, I. Hacihaliloglu, and V. M. Patel, "Medical transformer: Gated axial-attention for medical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2021, pp. 36–46.
- [68] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [69] A. Song, Z. Wu, X. Ding, Q. Hu, and X. Di, "Neurologist standard classification of facial nerve paralysis with deep neural networks," *Future Internet*, vol. 10, no. 11, p. 111, 2018.

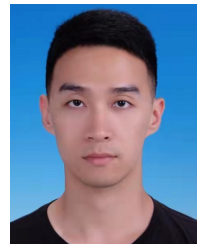


Xuri Ge received the M.S. degree in computer science from the School of Information Science and Engineering, Xiamen University, China, in 2020. He is currently pursuing the Ph.D. degree with the school of computer science, University of Glasgow, Scotland, UK. His current research interests include computer vision, medical image processing, and multimedia processing.



related to the above themes.

Joemon M. Jose is a Professor at the School of Computing Science, University of Glasgow, Scotland and a member of the Information Retrieval group. His research focuses around the following three themes: (i) Social Media Analytics; (ii) Multi-modal interaction for information retrieval; (iii) Multimedia mining and search. He has published over 300 papers with more than 8000 Google Scholar citations, and an h-index of 47. He leads the Multimedia Information Retrieval group which investigates research issues



Pengcheng Wang is currently a Researcher with Tomorrow Advancing Life Education Group (TAL), Beijing, China. His research interests include the applied AI, such as intelligent multimedia processing, and computer vision. As a Key Team Member, he achieved the best performance in various competitions, such as the EmotioNet facial expression recognition challenge.



Arunachalam Iyer is a Honorary Associate Clinical Professor of Surgery, University of Glasgow, Scotland, UK. He is also a consultant ENT surgeon, University Hospital Monklands, UK. He is currently working as a Consultant ENT surgeon & Otologist at Lanarkshire and a council member of the British society of Otology and the section of Otology, Royal society of Medicine. He also teaches at the Temporal bone course at University of Glasgow and is involved in training junior doctors.



Xiao Liu received the Ph.D. degree in computer science from Zhejiang University in 2015. He worked at Baidu from 2015 to 2019 and at Tomorrow Advancing Life Education Group (TAL) from 2019 to 2021. He is currently a Researcher with the Online Media Business Unit at Tencent, Beijing, China. His research interests include the applied AI, such as intelligent multimedia processing, computer vision, and learning systems. His research results have expounded in more than 40 publications at journals and conferences, such as IEEE TRANSACTIONS ON IMAGE PROCESSING, IEEE TRANSACTIONS ON NEURAL NETWORKS AND LEARNING SYSTEMS, CVPR, ICCV, ECCV, AAAI, and MM. As a Key Team Member, he achieved the best performance in various competitions, such as the ActivityNet challenges, NTIRE super resolution challenge, and EmotioNet facial expression recognition challenge.



Hu Han (Member, IEEE) is a Professor at the Institute of Computing Technology (ICT), Chinese Academy of Sciences (CAS). His research interests include computer vision, pattern recognition, biometrics, and medical image analysis. He has published more than 70 papers with more than 4700 Google Scholar citations. He was a recipient of the IEEE Signal Processing Society Best Paper Award (2020), ICCV2021 Human-centric Trustworthy Computer Vision Best Paper Award, IEEE FG2019 Best Poster Presentation

Award, and 2016/2018 CCBR Best Student/Poster Award. He is/was the Associate Editor of Pattern Recognition, IJCAI2021 SPC, ICPR2020 Area Chair, ISBI2022 Session Chair, and VALSE LAC. He has organized a number of special sessions and workshops of "vision based vital sign analysis and health monitoring" in ICCV2021, CVPR2020, FG2019/20/21.