

AnchorFace: An Anchor-based Facial Landmark Detector Across Large Poses

Zixuan Xu,^{1*} Banghuai Li,^{2*} Ye Yuan² and Miao Geng³

¹ Peking University

² Megvii Research

³ Beihang University

zixuanxu@pku.edu.cn, {libanghuai,yuanye}@megvii.com, geng_m@buaa.edu.cn

Abstract

Facial landmark localization aims to detect the predefined points of human faces, and the topic has been rapidly improved with the recent development of neural network based methods. However, it remains a challenging task when dealing with faces in unconstrained scenarios, especially with large pose variations. In this paper, we target the problem of facial landmark localization across large poses and address this task based on a split-and-aggregate strategy. To split the search space, we propose a set of anchor templates as references for regression, which well addresses the large variations of face poses. Based on the prediction of each anchor template, we propose to aggregate the results, which can reduce the landmark uncertainty due to the large poses. Overall, our proposed approach, named AnchorFace, obtains state-of-the-art results with extremely efficient inference speed on four challenging benchmarks, i.e. AFLW, 300W, Menpo, and WFLW dataset. Code will be available soon.

Introduction

Facial landmark localization, or face alignment, refers to detect a set of predefined landmarks on the human face. It is a fundamental step for many facial related applications, e.g. face verification/recognition, expression recognition, and facial attribute analysis.

With the recent development of convolutional neural network based methods (Ma et al. 2018; Xu et al. 2020), the performance for facial landmark localization in constrained scenarios has been greatly improved (Dong et al. 2018; Zhu et al. 2019; Wu et al. 2018). However, unconstrained scenarios, for example, faces with large pose, still limit the wide application of the existing landmark algorithms. In this paper, we target to address the problem of facial landmark localization across large poses.

There are two challenges for facial landmark detection across large poses. On one hand, faces with large poses will significantly increase the difficulty for landmark localization due to the large variations among different poses. As shown in Fig. 1, directly regressing the point coordinates may not be able to localize every landmark point precisely. On the

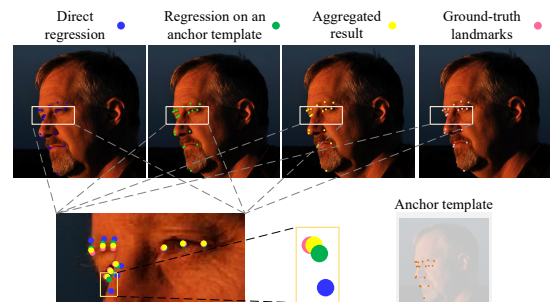


Figure 1: A comparison between direct regression and anchor-based regression (AnchorFace). Our AnchorFace includes two steps. The first step is to introduce the anchor templates and regress the offsets based on each anchor template (Second Column). The second step is to aggregate the prediction results from multiple anchor templates (Third Column)

other hand, there usually exists a large probability of uncertainty due to the self-occlusion and noisy annotations. For example, occlusion will usually lead to invisible landmarks, which will increase the uncertainty for the landmark prediction. Besides, the faces with a large pose will also cause difficulty during the data annotation process.

To address the above two challenges, we propose a novel pipeline for facial landmark localization based on an anchor-based design. The new pipeline includes two steps: split and aggregate. An overview of our pipeline can be found in Fig. 2. To deal with the first challenge with large pose variations, we adopt the divide-and-conquer way following an anchor-based design. We propose to use the anchor templates to split the search space, and each anchor will serve as a reference for regression. This can significantly reduce the pose variations for each anchor. To address the second issue with pose uncertainty, we propose to aggregate each anchor result weighted by the predicted confidence.

In summary, we propose AnchorFace to implement the split-and-aggregate strategy. There are three contributions in our paper.

- We propose a novel pipeline with a split-and-aggregate strategy which can well address the challenges for face

*These authors contributed equally

alignment across large poses.

- To implement the split-and-aggregate strategy, we introduce the anchor design into the facial landmark problem, which can simplify the search space for each anchor template and meanwhile improve the robustness for landmark uncertainty.
- Our proposed AnchorFace can achieve promising results on four challenging benchmarks with a realtime inference speed of ~ 45 FPS¹.

Related Work

Facial Landmark Localization. In the literature of facial landmark localization, a number of achievements have been developed including the classic ASMs (Milborrow and Nicolls 2008), AAMs (Ikeuchi, Hebert, and Delingette 1995; Kahraman et al. 2007), CLMs (Cristinacce and Cootes 2008, 2006), and Cascaded Regression Models (Cao et al. 2014; Zhu et al. 2015, 2016a). Nowadays, more and more deep learning-based methods have been applied in this area. These deep learning based methods could be divided into two categories, i.e. coordinate regression methods and heatmap regression methods.

Coordinate regression methods directly map the discriminative features to the target landmark coordinates. The earliest work can be dated to (Sun, Wang, and Tang 2013). Sun et al. (Sun, Wang, and Tang 2013) used a three-level cascade CNN to do facial landmark localization in a coarse-to-fine manner, and achieved promising localization accuracy. MDM (Trigeorgis et al. 2016) was the first to apply a recurrent convolutional network model for facial landmark localization in an end-to-end manner. Zhang et al. (Zhang et al. 2014) utilized a multi-task learning framework to optimize facial landmark localization and correlated facial attributes analysis simultaneously. Recently, Wingloss (Feng et al. 2017) was proposed as a new loss function for landmark localization, which can obtain robust performance against widely used L_2 loss.

Heatmap regression methods generate a probability heatmap for each landmark, respectively. Benefit from FCN (Long, Shelhamer, and Darrell 2015) and Hourglass (Newell, Yang, and Deng 2016), heatmap regression methods have been successfully applied to landmark localization problems and have achieved state-of-the-art performance. JMFA (Deng et al. 2017) achieved high localization accuracy with a stacked hourglass network (Newell, Yang, and Deng 2016) for multi-view facial landmark localization in the Menpo (Zafeiriou et al. 2017) competition. Yang et al. (Yang, Liu, and Zhang 2017) adopted a supervised face transformation to normalize the faces, then employed an Hourglass network to regress it. Recently, LAB (Wu et al. 2018) proposed to use additional boundary lines as the geometric structure of a face image to help facial landmark localization.

Faces with Large Pose. Large pose is a challenging task for facial landmark localization, and different strate-

gies have been proposed to address the difficulty. Multi-view framework and 3D model are two popular ways. Multi-view framework uses different landmark configurations for different views. For example, TSPM (Zhu and Ramanan 2012) and CDM (Yu et al. 2013) employ DPM-like (Felzenszwalb et al. 2010) method to align faces with different shape models, and choose the highest possibility model as the final result. However, multi-view methods have to cover each view, making it impractical in the wild. 3D face models have been widely used in recent years, which fit a 3D morphable model (3DMM) (Banz and Vetter 2003) by minimizing the difference between face image and model appearance. Lost face information can be recovered to localize the invisible facial landmarks (Jourabloo et al. 2017; Liu et al. 2017; Zhu et al. 2016b; Bulat and Tzimiropoulos 2017). However, 3D face models are limited by their own database and the iterative label generation method. Besides, researchers have applied multi-task learning to address the difficulties resulting from pose variations. Other facial analysis tasks, such as pose estimation or facial attributes analysis, can be jointly trained with facial landmark localization (Xu and Kakadiaris 2017; Zhang et al. 2016). With joint training, multi-task learning can boost the performance of each subtask. The facial landmark localization task can achieve robust performance. But the multi-task framework is not specially designed for landmark localization, it contains much redundant information and contributes to large models.

In this paper, we propose an anchor-based model for facial landmark localization. Different from (Xiong et al. 2019), which utilized anchor points to predict the positions of a human 3D pose, our approach introduces a split-and-aggregate pipeline for the facial landmark localization. Anchor is utilized as a reference for regression in our approach. Overall, our model requires neither cascaded networks nor large backbones, leading to a great reduction in model parameters and computation complexity, while still achieving comparable or even better accuracy.

Proposed Method

In this paper, we propose a new split-and-aggregate strategy for facial landmark detector across large poses. An overview of our pipeline can be found in Fig. 2. To implement the split-and-aggregate strategy, we introduce the anchor-based design, and our approach is named AnchorFace. In the following section, we will discuss the **split** and **aggregate** steps separately, followed by the details on the network training.

Split Step

Due to the large pose variations among different poses, it is a challenging problem to directly regress the facial landmarks while maintaining high localization precision. In this paper, we propose to utilize the divide-and-conquer way to address the issue from large pose variations. More specifically, we propose to employ the anchor templates as regression references to split the search space. Different from the traditional methods which regress the landmarks with a uniform facial landmark detector, we propose to regress the offsets base on a set of anchor templates.

¹The computational speed of ~ 45 FPS is calculated on one Nvidia Titan Xp GPU with batchsize 1.

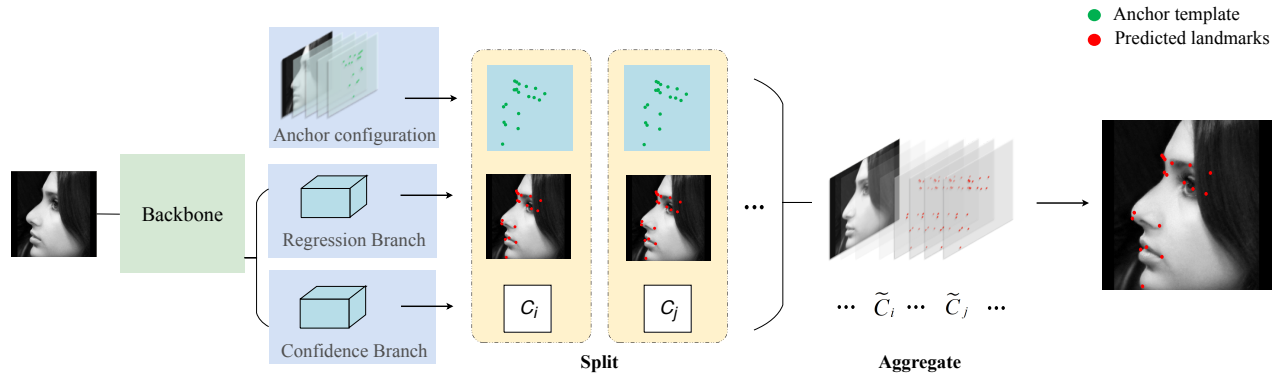


Figure 2: The pipeline of our proposed AnchorFace landmark detector. AnchorFace is based on a split-and-aggregate strategy, which consists of the backbone and two functional branches: the offset regression branch and the confidence branch. In the split step, we predict the landmark position based on each anchor template. During aggregate step, the predictions of multiple anchor templates are averaged by weighted confidence

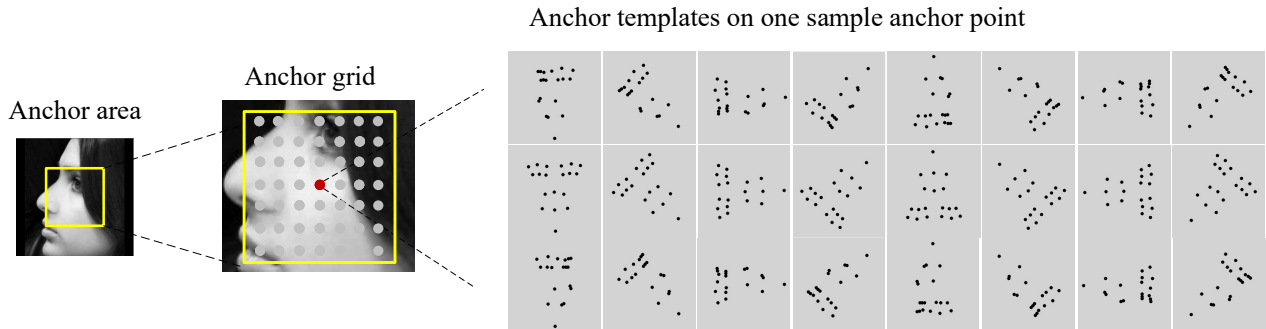


Figure 3: An illustration of our anchor configuration. Anchor area is a region centered at the image center with a spatial neighborhood. Based on the anchor area, we setup a grid of anchor points, where each anchor point contains a set of anchor templates to model various pose variations

Anchor Configuration. As shown in Fig. 3, there are three hyper-parameters for designing the anchor configuration: anchor area, anchor grid, and anchor templates.

Anchor area is denoted as the region to set the anchors. It is usually centered at the image center with a spatial neighborhood. The reason to define the anchor area is that the input image is cropped to put the face near the image center. Thus, we select a region near the image center, which is called anchor area, to set up the anchors. Based on the anchor area, we sample a set of anchor points in a grid, e.g. a 7×7 grid, as shown in Fig. 3. Each anchor point can be considered as the center of a set of anchor templates. The anchor templates are designed to address the challenges from large pose variations. Intuitively, these anchor templates are used to split the search space for regression and can serve as references for offsets prediction. Therefore, the sampling of anchor templates should be able to cover different variations of large poses and reduce the redundancy for the anchor sets.

To implement the anchor templates, we present two potential ways. The first one is to hand-design the anchors based on prior knowledge. The second one integrates the proposals generated by the data distribution.

An overview of our hand-designed anchors can be found in Fig. 3. For each anchor point, we explore the 3D pose spaces (yaw, roll, pitch) and design the pose-level anchor set as follows. As unconstrained large-pose faces have large variations on the yaw direction, we first select N_{yaw} base anchors ($N_{yaw} = 3$ in our paper representing the anchors for the left, frontal, right faces). To generate the N_{yaw} base anchors, we utilize a heuristic approach to divide the training faces into three buckets and compute the average face landmarks for each bucket to obtain the anchor proposal. More specifically, we use the ratio of two eyes' width for bucket assignment. We define an indicator to estimate the yaw angle of each training face:

$$r = \frac{|p_{l_1} - p_{l_2}|_2}{|p_{r_1} - p_{r_2}|_2} - \frac{|p_{r_1} - p_{r_2}|_2}{|p_{l_1} - p_{l_2}|_2}, \quad (1)$$

where $p_{l_1}, p_{l_2}, p_{r_1}, p_{r_2}$ are the coordinates of left eye inner corner, left eye outer corner, right eye inner corner, and right eye outer corner respectively. With a threshold γ , we put the faces into the left or right bucket, when $r > \gamma$ or $r < -\gamma$. The other faces will be kept into the frontal bucket. We set $\gamma = 6$ in our experiments, as shown in Fig. 4.

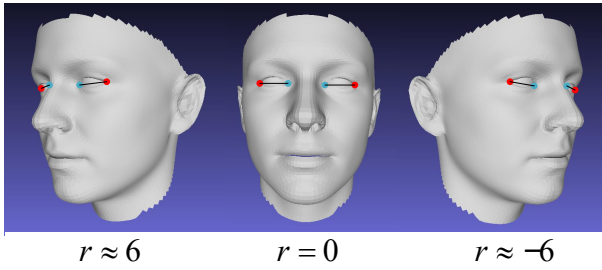


Figure 4: An illustration of the metric r for classifying the faces into three buckets along the yaw direction

Based on the N_{yaw} base anchors, to cover the roll variations, we rotate each anchor on the roll dimension. For example, we can get twenty-four templates by rotating the basic three anchors each 45° from 0° to 360° . Optionally, we can involve the pitch variations by directly projecting (rotating) along the pitch dimension. However, based on our experimental results, the anchors designed along the pitch view cannot further improve the performance but compromise the computational speed. Thus, in our final design, only anchors along the yaw and roll dimensions are utilized as shown in Fig. 3.

An alternative solution for the anchor design is based on the data distribution among the training faces. We first perform KMeans clustering, and we can generate a set of base anchors. One example is shown in Fig. 5. We can see that the clustered anchors among all the training faces obtain similar anchors along the yaw direction, as discussed in hand-designed anchors. Following the similar steps for the hand-designed anchors, we can rotate the generated prototypes along the roll and pitch direction to generate more anchors.

Regression and Confidence Branch. Based on the anchor proposals, we design a new head structure which involves two branches: regression branch and confidence branch. Regression branch aims to regress the landmark coordinate offsets based on each anchor. Confidence branch assigns each anchor with a confidence score. Among all the anchor templates, those anchors which are close to the pose of the ground-truth face should be given higher confidence.

As shown in Fig. 2, both the confidence branch and the regression branch are built upon the output feature map of the backbone network. While we set $h \cdot w$ anchors in the image, the output of the confidence and regression branch are $C_{con} \cdot h \cdot w$ and $C_{reg} \cdot h \cdot w$ respectively, where C_{con} and C_{reg} are denoted as the output channel number of the confidence branch and the regression branch respectively. Here $C_{con} = K$ and $C_{reg} = K \cdot 2L$, where K, L refer to the number of anchor templates on each anchor point and the number of facial landmarks respectively.

Aggregate Step

Large-pose faces will increase the uncertainty for the landmark prediction. To address this problem, we propose to aggregate the predictions from different anchor templates. More specifically, we first set a threshold C_{th} to pick up the

- Sample anchor templates by hand-design
- Sample anchor templates by KMeans

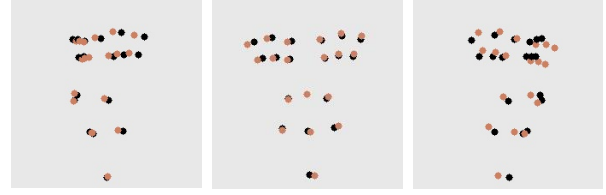


Figure 5: A comparison of three base anchors generated by hand-design approach and KMeans clustering based on AFLW (Köstinger et al. 2011) dataset

Symbol	Definition
A	A set of Anchor points on the spatial anchor grid
a	One anchor point $a \in A$
T	A set of anchor templates as in Fig. 3
$T(a, t)$	Anchor template $t \in T$ centering at anchor point a
$T_j(a, t)$	Landmark j on anchor template $T(a, t)$
$O(a, t)$	Output from the regression branch based on $T(a, t)$
$\bar{O}(a, t)$	Ground-truth (GT) offsets based on $T(a, t)$
$C(a, t)$	Output from the confidence branch based on $T(a, t)$
$\bar{C}(a, t)$	Confidence GT label based on $T(a, t)$

Table 1: The definition of symbols

reliable anchor predictions. The anchor predictions with low confidence scores are regarded as outliers and will be discarded. The remaining anchor predictions will be averaged by the weighted confidence for each prediction. As a result, the position of landmark j can be obtained as the weighted average of the outputs of all anchor faces as:

$$\tilde{S}_j = \frac{\sum_{a \in A, t \in T} \tilde{C}(a, t) \cdot (O_j(a, t) + T_j(a, t))}{\sum_{a \in A, t \in T} \tilde{C}(a, t)}, \quad (2)$$

where

$$\tilde{C}(a, t) = \begin{cases} 0, & C(a, t) < C_{th} \\ C(a, t), & \text{others} \end{cases} \quad (3)$$

The definition of the symbols can be found from Table 1, and the threshold C_{th} is set to 0.6 in our experiments.

Network Training

In this subsection, we will discuss the ground-truth setting for the regression and confidence branch as well as the related losses. For the regression branch, the target is to regress the offsets against each of the predefined anchor. The regression loss L_{reg} can be defined as:

$$L_{reg} = \sum_{a \in A, t \in T} C(a, t) \sum_j |O_j(a, t) - \bar{O}_j(a, t)|, \quad (4)$$

where $O_j(a, t)$ and $\bar{O}_j(a, t)$ refer to the prediction offsets and the ground-truth offsets for j th landmark. $C(a, t)$ is denoted as the confidence weight for the anchor template

$T(a, t)$. The detailed symbol definitions can be found in Table 1.

For the confidence branch, we set the targeted confidence output $\bar{C}(a, t)$ as the inverse L2 distance between the anchor pose \mathbf{v}_1 and the ground-truth pose \mathbf{v}_2 as $\|\mathbf{v}_1 - \mathbf{v}_2\|_2$, where $\mathbf{v}_1, \mathbf{v}_2$ refer to flatten landmark coordinates. To normalize the pose difference, we perform a tanh operation as:

$$\bar{C} = \tanh\left(\left(\frac{\|\mathbf{v}_1 - \mathbf{v}_2\|_2}{\beta \cdot 2L}\right)^{-1}\right), \quad (5)$$

where β is a hyperparameter and L refers to the count of facial landmarks. The confidence loss is then defined as:

$$L_{con} = \sum_{a \in A, t \in T} (-C(a, t) \cdot \log \bar{C}(a, t) - (1 - C(a, t)) \cdot \log(1 - \bar{C}(a, t))). \quad (6)$$

The network is jointly supervised by the two loss functions above with end-to-end training. The final training loss is then defined as:

$$L_{total} = L_{reg} + \lambda \cdot L_{con} \quad (7)$$

where λ is a hyperparameter in our method, and it is insensitive to the localization accuracy in our experiments.

Experiment

Experiment Settings

Datasets. The experiments are evaluated on four challenging datasets, i.e. AFLW (Köstinger et al. 2011), 300W (Sagonas et al. 2013), Menpo (Deng et al. 2019; Zafeiriou et al. 2017), and WFLW (Wu et al. 2018). They are all widely used benchmarks in the facial landmark research area. More details about these datasets can be found in our supplementary materials.

Evaluation Metric. We adopt the normalized mean error (NME) for evaluation. The normalized mean error is defined as the average Euclidean distance between the predicted facial landmark locations $O_{i,j}$ and their corresponding ground-truth facial landmark annotations $\bar{O}_{i,j}$:

$$NME = \frac{1}{N} \sum_{i=1}^N \frac{\frac{1}{L} \sum_{j=1}^L |O_{i,j} - \bar{O}_{i,j}|_2}{d} \quad (8)$$

where N is the number of images in the testing set, L is the number of landmarks, and d is the normalization factor. On AFLW dataset, we follow (Zhu et al. 2015) to use face size as the normalization factor. On Menpo dataset, we use the distance between left-top corner and right-bottom corner as the normalization factor. On 300W and WFLW dataset, we follow MDM (Trigeorgis et al. 2016) and (Sagonas et al. 2013) to use the ‘‘inter-ocular’’ normalization factor, i.e. the distance between the outer eye corners.

In addition, on WFLW dataset, two further statistics i.e. the area-under-the-curve (AUC) (Yang et al. 2015) and the failure rate (which is defined as the proportion of failed detected faces) are measured for further analysis. Especially, any normalized error above 0.1 is considered as a failure (Wu et al. 2018).

Implementation Details. In our method, the original images are cropped and resized to a fixed resolution, i.e. 224×224 , according to the provided bounding boxes. Anchor templates are generated based on KMeans clustering following, while anchor area and anchor grid are set as 56×56 and 7×7 respectively. Random rotation and translation are applied for data augmentation. We apply the Adam optimizer with the weight decay of 1×10^{-5} and train the network for 50 epochs in total. The learning rate is set to 1×10^{-3} and divided by ten at 20-th, 30-th, 40-th epoch. $\beta = 0.05$ and $\lambda = 0.5$ are applied to all models across four benchmarks. If there are no special instructions, ShuffleNet-V2 (Ma et al. 2018) is utilized as the backbone for our algorithm in this paper.

Comparison with State-of-the-art Methods

For a fair comparison, we only compare the methods following the standard settings, as discussed in Section . Therefore, those methods which are trained from external datasets or combined multiple datasets are not compared. Due to the space limit, comparisons about Menpo are reported in our supplementary materials.

AFLW dataset: We first evaluate our algorithm on the AFLW dataset. The performance comparisons are given in Table 2. It can be observed that, on this large dataset, our network outperforms the other approaches. As mentioned in Section , AFLW contains lots of faces with large poses. Note that our method has a significant improvement on AFLW-Full set against AFLW-Frontal, which means that we achieve more robust localization performance on faces in unconstrained scenarios, including large pose. This essentially validates the superiority of our approach.

300W dataset: We compare our approach against several state-of-the-art methods on 300W Fullset. The results are shown in Table 3. Since there are fewer large pose variations across the whole dataset and the cropped faces normally center near the image center point, 300W dataset is not very challenging compared with the other three benchmarks. However, our algorithm still can achieve promising localization performance with an efficient speed at ~ 45 fps with batch size 1. Compared with LAB (Wu et al. 2018), which is slightly better than our method, our approach is much faster (45 fps vs 17 fps).

WFLW dataset: A comparison of the performance from our proposed approach as well as state-of-the-art methods on WFLW dataset is shown in Table 4. Considering the difficulty of WFLW, we attempt to adopt HRNet-18 (Sun et al. 2019) as the backbone to implement our AnchorFace. As indicated in Table 4, AnchorFace equipped with HRNet-18 (Sun et al. 2019) achieves the best performance on NME & AUC metrics and achieves the second best performance on Failure Rate metric. Furthermore, the efficient AnchorFace with ShuffleNet-V2 (Ma et al. 2018) also achieves impressive results considering the huge computation cost of other methods. This also directly verifies the versatility of our method. Detailed results about each subset of WFLW are reported in our supplementary materials.

Methods	AFLW-Full	AFLW-Frontal
LBF (Ren et al. 2016)	4.24	2.74
CFSS (Zhu et al. 2015)	3.92	2.69
CCL (Zhu et al. 2016a)	2.72	2.17
TSR (Lv et al. 2017)	2.17	-
SAN (Dong et al. 2018)	1.91	1.85
Wing (Feng et al. 2017)	1.65	-
SA (Liu et al. 2019)	1.62	-
ODN (Zhu et al. 2019)	1.63	1.38
AnchorFace	1.56	1.38

Table 2: Normalized mean error (%) on AFLW-Full and AFLW-Frontal set.

Methods	Common	Challenge	Full
MDM (Trigeorgis et al. 2016)	4.83	10.14	5.88
Two-Stage (Lv et al. 2017)	4.36	7.42	4.96
RDR (Xiao et al. 2017)	5.03	8.95	5.80
Pose-Invariant (Jourabloo et al. 2017)	5.43	9.88	6.30
SBR (Dong et al. 2018)	3.28	7.58	4.10
PCD-CNN (Kumar and Chellappa 2018)	3.67	7.62	4.44
LAB (Wu et al. 2018)	2.98	5.19	3.49
SAN (Dong et al. 2018)	3.34	6.60	3.98
ODN (Zhu et al. 2019)	3.56	6.67	4.17
AnchorFace	3.12	6.19	3.72

Table 3: Normalized mean error (%) on 300W Common subset, Challenging subset, and Full set.

Methods	NME(%)	Failure Rate(%)	AUC	Flops(G)
CFSS (Zhu et al. 2015)	9.07	29.40	0.3659	-
DVLN (Wu and Yang 2017)	6.08	10.84	0.4551	-
LAB (Wu et al. 2018)	5.27	7.56	0.5323	18.85
SAN (Dong et al. 2018)	5.22	6.32	0.5355	-
Wing (Feng et al. 2017)	5.11	6.00	0.5504	5.40
AVS (Qian et al. 2019)	5.25	7.44	0.5034	-
AWing (Wang, Bo, and Fuxin 2019)	4.36	2.84	0.5719	26.79
HRNET (Sun et al. 2019)	4.60	4.64	0.5237	-
AnchorFace	4.62	4.20	0.5516	1.71
AnchorFace*	4.32	2.96	0.5769	5.30

Table 4: Evaluation results about WFLW test set for NME(%), Failure Rate@0.1(%) and AUC@0.1. * means we adopt HRNet-18 (Sun et al. 2019) as the backbone network to implement AnchorFace.

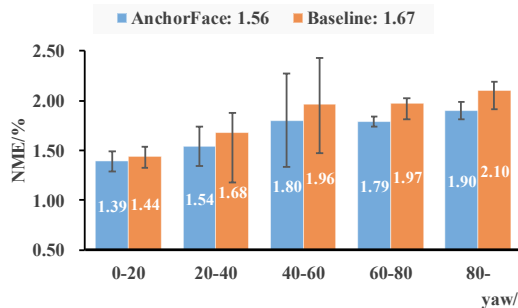


Figure 6: A comparison of baseline and AnchorFace across yaw dimension

Model Analysis

Our proposed AnchorFace introduces a novel split-and-aggregate strategy based on anchor design to address the face alignment across large poses. In this section, we perform further analysis of its mechanism.

Anchor design. Anchor templates serve as regression references to split the search space in our proposed approach. In comparison with directly regressing target landmark coordinates in whole 2D space, regress offsets based on anchor templates can simplify the search space and boost the robustness of localization accuracy. We conduct several experiments on AFLW dataset and make statistics across yaw dimension which is shown in Fig. 6. It is quite clear that AnchorFace significantly outperforms the baseline with lower NME and smaller variances in each subinterval especially for large pose, which can well verify our assumptions.

Split-and-aggregate strategy. In our proposed algorithm, we follow the divide-and-conquer way to address the challenges for face alignment across large poses. To verify its effectiveness, we adopt *Pearson correlation coefficient* to measure the correlation between **confidence scores** $C(a, t)$ and **prediction errors** $|O(a, t) - \bar{O}(a, t)|_2$:

$$P = \frac{1}{N} \sum_{i=1}^N [r_{a,t}^i(|O(a, t) - \bar{O}(a, t)|_2, C(a, t))] \quad (9)$$

Where r represents the calculation function of *Pearson correlation coefficient*. We conduct experiments on AFLW dataset and get $P = -0.82$, which means a strong negative correlation between them. In other words, anchor template with larger confidence score can achieve more accurate predicted landmarks. It can help filter prediction outliers and aggregate remaining predictions to mitigate the uncertainty of the localization result on a single anchor face. Due to the confidence score is defined as mathematical modeling of the distance between the anchor pose and the ground-truth pose, we can come to another conclusion that *closer* anchors tend to achieve more accurate localization, which also directly proves our search space split strategy based on anchors. Comparison details can be found in our ablation studies and we show several intuitive samples in our supplementary materials.

Ablation Study

In this section, we perform the ablation study for our proposed algorithm on the AFLW dataset, which is a challenging benchmark with large pose variations. More specifically, we divide the test set into four subsets according to the yaw dimension, i.e. Light ($0^\circ \sim 30^\circ$), Medium ($30^\circ \sim 60^\circ$), Large ($60^\circ \sim 90^\circ$), and Heavy ($90^\circ \sim$). Normalized mean error is utilized to evaluate the performance of our algorithm. Without explicitly specified, we use anchor templates as KMeans-24 (KMeans clustering to generate 24 anchor templates), anchor area as 56×56 , anchor grid as 7×7 , and the aggregating strategy is weighted average for ablation.

Comparison with Regression Baseline. Table 5 compares the performance of our proposed approach with the baseline of direct regression on AFLW dataset. “Baseline”

Method	Full	Light	Medium	Large	Heavy
Baseline	1.67	1.46	1.92	1.99	2.13
AnchorFace	1.56	1.40	1.74	1.80	1.96

Table 5: A comparison of direct regression and anchor-based regression.

Template	Full	Light	Medium	Large	Heavy
Kmeans-3	1.60	1.43	1.80	1.83	2.10
Kmeans-24	1.56	1.40	1.74	1.80	1.96
Kmeans-48	1.58	1.41	1.79	1.82	2.06
HandDesign-24	1.58	1.42	1.76	1.84	2.00

Table 6: A comparison of different template settings.

directly maps the discriminative features to the target landmark coordinates with ShuffleNet-V2 backbone. A fully connected layer with length $2L$ is used as the output of the baseline network. As shown in Table 5, our proposed anchor-based method significantly outperforms the baseline by a large margin across yaw variations. The improvements are attributed to two reasons. First, the anchor design can significantly reduce the search space and simplify the regression problem. Second, the aggregating of different anchors can further improve model robustness.

Comparison of Various Split Configurations. Due to the challenges from the large-pose faces, we propose a set of anchor templates as references for regression to split search space. Split strategy consists of three hyper-parameters: anchor templates, anchor area, and anchor grid, as shown in Fig. 3.

Anchor template plays a voting role in our method as a reference for regression. As mentioned in Section , we get three basic template faces from the training dataset by hand design or KMeans clustering. Then we do some transformations to get more templates, corresponding to the pose variations in yaw, roll, and pitch dimension. By comparing KMeans-24 against HandDesign-24 in Table 6, KMeans is better than hand-design approach based on the same anchor number (24). The potential reason is that KMeans utilizes more data features to generate the base anchors, which should be more general compared with hand-designed based anchors. Besides, as shown in Table 6, 24 may be a good option for the number of anchor templates compared with 3 or 48 in our algorithm.

Anchor area is the area where we set anchors in the image for the spatial domain. As the input image is cropped and resized to 224×224 and the face is around the image

Anchor area	Full	Light	Medium	Large	Heavy
112×112	1.58	1.40	1.79	1.83	2.06
56×56	1.56	1.40	1.74	1.80	1.96
28×28	1.58	1.42	1.79	1.82	2.04
14×14	1.58	1.41	1.81	1.83	2.01

Table 7: A comparison of different anchor area settings.

Anchor count	Full	Light	Medium	Large	Heavy
3×3	1.58	1.40	1.78	1.82	2.17
5×5	1.57	1.40	1.77	1.82	2.05
7×7	1.56	1.40	1.74	1.80	1.96
13×13	1.56	1.40	1.75	1.81	2.07

Table 8: A comparison of different anchor grid settings.

Aggregate	Full	Light	Medium	Large	Heavy
Argmax	1.58	1.42	1.76	1.83	2.00
Weighted	1.56	1.40	1.74	1.80	1.96
Mean	1.61	1.43	1.80	1.89	2.22

Table 9: A comparison of different aggregation strategies.

center, we set anchor points at a center area with size 14×14 , 28×28 , 56×56 , 128×128 , respectively. As shown in Table 7, 56×56 around the image center would be a good choice for putting anchors in the spatial domain.

Anchor grid defines how many anchors we set in the anchor area. For example, 7×7 means we sample 49 spatial points in a 7×7 grid from the anchor area to generate the anchor templates. As shown in Table 8, it is a good choice to set at 7×7 .

Comparison of Various Aggregate Strategies. To mitigate the uncertainty of the localization result on a single anchor face, we aggregate the predictions from different anchor templates. We introduce three aggregate strategies: Mean, Argmax, confidence weighted voting (Weighted). As shown in Table 9, aggregating the predictions with weighted confidences can obtain superior results compared with the argmax choice without aggregating. Besides, the confidence generated by the confidence branch is important if we compare the strategy of “Weighted” and “Mean”.

Conclusions

In this paper, a novel split-and-aggregate strategy is proposed for large-pose faces. By introducing an anchor-based design, our proposed approach can simplify the regression problem by splitting the search space. Moreover, aggregating the prediction results contributes to reducing uncertainty and improving the localization performance. As validated on four challenging benchmarks, our proposed AnchorFace obtains state-of-the-art results with extremely fast inference speed.

References

- Blanz, V.; and Vetter, T. 2003. Face recognition based on fitting a 3D morphable model. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 25(9): 1063–1074. doi:10.1109/TPAMI.2003.1227983.
- Bulat, A.; and Tzimiropoulos, G. 2017. How Far are We from Solving the 2D & 3D Face Alignment Problem? (and a Dataset of 230,000 3D Facial Landmarks). *Proceedings of the IEEE International Conference on Computer Vision 2017-October*: 1021–1030. ISSN 15505499. doi:10.1109/ICCV.2017.116.
- Cao, X.; Wei, Y.; Wen, F.; and Sun, J. 2014. Face Alignment by Explicit Shape Regression. *International Journal of Computer Vision* 107(2): 177–190. ISSN 1573-1405. doi:10.1007/s11263-013-0667-3. URL <https://doi.org/10.1007/s11263-013-0667-3>.
- Cristinacce, D.; and Cootes, T. 2008. Automatic feature localisation with constrained local models. *Pattern Recognition* 41(10): 3054 – 3067. ISSN 0031-3203. doi:<https://doi.org/10.1016/j.patcog.2008.01.024>. URL <http://www.sciencedirect.com/science/article/pii/S0031320308000630>.
- Cristinacce, D.; and Cootes, T. F. 2006. Feature Detection and Tracking with Constrained Local Models. In *BMVC*.
- Deng, J.; Roussos, A.; Chryso, G.; Ververas, E.; Kotsia, I.; Shen, J.; and Zafeiriou, S. 2019. The Menpo Benchmark for Multi-pose 2D and 3D Facial Landmark Localisation and Tracking. *International Journal of Computer Vision* 127(6): 599–624. ISSN 1573-1405. doi:10.1007/s11263-018-1134-y. URL <https://doi.org/10.1007/s11263-018-1134-y>.
- Deng, J.; Trigeorgis, G.; Zhou, Y.; and Zafeiriou, S. 2017. Joint Multi-View Face Alignment in the Wild. *IEEE Transactions on Image Processing* 28: 3636–3648.
- Dong, X.; Yan, Y.; Ouyang, W.; and Yang, Y. 2018. Style Aggregated Network for Facial Landmark Detection. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 379–388. doi:10.1109/CVPR.2018.00047.
- Dong, X.; Yu, S.-I.; Weng, X.; Wei, S.-E.; Yang, Y.; and Sheikh, Y. 2018. Supervision-by-Registration: An Unsupervised Approach to Improve the Precision of Facial Landmark Detectors. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Felzenszwalb, P. F.; Girshick, R. B.; McAllester, D.; and Ramanan, D. 2010. Object Detection with Discriminatively Trained Part-Based Models. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 32(9): 1627–1645. doi:10.1109/TPAMI.2009.167.
- Feng, Z.-H.; Kittler, J.; Awais, M.; Huber, P.; and Wu, X.-J. 2017. Wing Loss for Robust Facial Landmark Localisation with Convolutional Neural Networks. *arXiv e-prints* arXiv:1711.06753.
- Ikeuchi, K.; Hebert, M.; and Delingette, H. 1995. A Spherical Representation for Recognition of Free-Form Surfaces. *IEEE Transactions on Pattern Analysis & Machine Intelligence* 23(07): 681–690. ISSN 1939-3539. doi:10.1109/34.391410.
- Jourabloo, A.; Ye, M.; Liu, X.; and Ren, L. 2017. Pose-Invariant Face Alignment With a Single CNN. In *The IEEE International Conference on Computer Vision (ICCV)*.
- Kahraman, F.; Gokmen, M.; Darkner, S.; and Larsen, R. 2007. An Active Illumination and Appearance (AIA) Model for Face Alignment. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, 1–7. doi:10.1109/CVPR.2007.383399.
- Kumar, A.; and Chellappa, R. 2018. Disentangling 3D Pose in a Dendritic CNN for Unconstrained 2D Face Alignment. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Köstinger, M.; Wohlhart, P.; Roth, P. M.; and Bischof, H. 2011. Annotated Facial Landmarks in the Wild: A large-scale, real-world database for facial landmark localization. In *2011 IEEE International Conference on Computer Vision Workshops (ICCV Workshops)*, 2144–2151. doi:10.1109/ICCVW.2011.6130513.
- Liu, Y.; Jourabloo, A.; Ren, W.; and Liu, X. 2017. Dense Face Alignment. In *The IEEE International Conference on Computer Vision (ICCV) Workshops*.
- Liu, Z.; Zhu, X.; Hu, G.; Guo, H.; Tang, M.; Lei, Z.; Robertson, N. M.; and Wang, J. 2019. Semantic Alignment: Finding Semantically Consistent Ground-truth for Facial Landmark Detection. *ArXiv* abs/1903.10661.
- Long, J.; Shelhamer, E.; and Darrell, T. 2015. Fully Convolutional Networks for Semantic Segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Lv, J.; Shao, X.; Xing, J.; Cheng, C.; and Zhou, X. 2017. A Deep Regression Architecture with Two-Stage Re-initialization for High Performance Facial Landmark Detection. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 3691–3700. doi:10.1109/CVPR.2017.393.
- Lv, J.; Shao, X.; Xing, J.; Cheng, C.; and Zhou, X. 2017. A Deep Regression Architecture With Two-Stage Re-Initialization for High Performance Facial Landmark Detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Ma, N.; Zhang, X.; Zheng, H.-T.; and Sun, J. 2018. ShuffleNet V2: Practical Guidelines for Efficient CNN Architecture Design. In *The European Conference on Computer Vision (ECCV)*.
- Milborrow, S.; and Nicolls, F. 2008. Locating Facial Features with an Extended Active Shape Model. In Forsyth, D.; Torr, P.; and Zisserman, A., eds., *Computer Vision – ECCV 2008*, 504–513. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN 978-3-540-88693-8.
- Newell, A.; Yang, K.; and Deng, J. 2016. Stacked Hourglass Networks for Human Pose Estimation. In Leibe, B.; Matas, J.; Sebe, N.; and Welling, M., eds., *Computer Vision – ECCV 2016*, 483–499. Cham: Springer International Publishing. ISBN 978-3-319-46484-8.

- Qian, S.; Sun, K.; Wu, W.; Qian, C.; and Jia, J. 2019. Aggregation via separation: Boosting facial landmark detector with semi-supervised style translation. *Proceedings of the IEEE International Conference on Computer Vision 2019*-Octob: 10152–10162. ISSN 15505499. doi:10.1109/ICCV.2019.01025.
- Ren, S.; Cao, X.; Wei, Y.; and Sun, J. 2016. Face Alignment via Regressing Local Binary Features. *IEEE Transactions on Image Processing* 25(3): 1233–1245. doi:10.1109/TIP.2016.2518867.
- Sagonas, C.; Tzimiropoulos, G.; Zafeiriou, S.; and Pantic, M. 2013. 300 Faces in-the-Wild Challenge: The First Facial Landmark Localization Challenge. *2013 IEEE International Conference on Computer Vision Workshops* 397–403.
- Sagonas, C.; Tzimiropoulos, G.; Zafeiriou, S.; and Pantic, M. 2013. 300 Faces in-the-Wild Challenge: The First Facial Landmark Localization Challenge. In *2013 IEEE International Conference on Computer Vision Workshops*, 397–403. doi:10.1109/ICCVW.2013.59.
- Sun, K.; Xiao, B.; Liu, D.; and Wang, J. 2019. Deep high-resolution representation learning for human pose estimation. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition 2019*-June: 5686–5696. ISSN 10636919. doi:10.1109/CVPR.2019.00584.
- Sun, Y.; Wang, X.; and Tang, X. 2013. Deep Convolutional Network Cascade for Facial Point Detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Trigeorgis, G.; Snape, P.; Nicolaou, M. A.; Antonakos, E.; and Zafeiriou, S. 2016. Mnemonic Descent Method: A Recurrent Process Applied for End-to-End Face Alignment. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 4177–4187. doi:10.1109/CVPR.2016.453.
- Wang, X.; Bo, L.; and Fuxin, L. 2019. Adaptive Wing Loss for Robust Face Alignment via Heatmap Regression URL <http://arxiv.org/abs/1904.07399>.
- Wu, W.; Qian, C.; Yang, S.; Wang, Q.; Cai, Y.; and Zhou, Q. 2018. Look at Boundary: A Boundary-Aware Face Alignment Algorithm. In *CVPR*.
- Wu, W.; and Yang, S. 2017. Leveraging Intra and Inter-Dataset Variations for Robust Face Alignment. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*.
- Xiao, S.; Feng, J.; Liu, L.; Nie, X.; Wang, W.; Yan, S.; and Kassim, A. 2017. Recurrent 3D-2D Dual Learning for Large-Pose Facial Landmark Detection. In *The IEEE International Conference on Computer Vision (ICCV)*.
- Xiong, F.; Zhang, B.; Xiao, Y.; Cao, Z.; Yu, T.; Zhou Tianyi, J.; and Yuan, J. 2019. A2J: Anchor-to-Joint Regression Network for 3D Articulated Pose Estimation from a Single Depth Image. In *Proceedings of the IEEE Conference on International Conference on Computer Vision (ICCV)*.
- Xu, X.; and Kakadiaris, I. A. 2017. Joint Head Pose Estimation and Face Alignment Framework Using Global and Local CNN Features. In *2017 12th IEEE International Conference on Automatic Face Gesture Recognition (FG 2017)*, 642–649. doi:10.1109/FG.2017.81.
- Xu, Z.; Li, B.; Yuan, Y.; and Dang, A. 2020. Beta R-CNN: Looking into Pedestrian Detection from Another Perspective. *Advances in Neural Information Processing Systems* 33.
- Yang, H.; Jia, X.; Loy, C. C.; and Robinson, P. 2015. An Empirical Study of Recent Face Alignment Methods. *CoRR* abs/1511.05049. URL <http://arxiv.org/abs/1511.05049>.
- Yang, J.; Liu, Q.; and Zhang, K. 2017. Stacked Hourglass Network for Robust Facial Landmark Localisation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*.
- Yu, X.; Huang, J.; Zhang, S.; Yan, W.; and Metaxas, D. N. 2013. Pose-Free Facial Landmark Fitting via Optimized Part Mixtures and Cascaded Deformable Shape Model. In *2013 IEEE International Conference on Computer Vision*, 1944–1951. doi:10.1109/ICCV.2013.244.
- Zafeiriou, S.; Trigeorgis, G.; Chrysos, G.; Deng, J.; and Shen, J. 2017. The Menpo Facial Landmark Localisation Challenge: A Step Towards the Solution. *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)* 2116–2125.
- Zhang, K.; Zhang, Z.; Li, Z.; and Qiao, Y. 2016. Joint Face Detection and Alignment Using Multitask Cascaded Convolutional Networks. *IEEE Signal Processing Letters* 23(10): 1499–1503. doi:10.1109/LSP.2016.2603342.
- Zhang, Z.; Luo, P.; Change Loy, C.; and Tang, X. 2014. Learning Deep Representation for Face Alignment with Auxiliary Attributes. *arXiv e-prints* arXiv:1408.3967.
- Zhu, M.; Shi, D.; Zheng, M.; and Sadiq, M. 2019. Robust Facial Landmark Detection via Occlusion-Adaptive Deep Networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhu, S.; Li, C.; Loy, C. C.; and Tang, X. 2015. Face Alignment by Coarse-to-Fine Shape Searching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 4998–5006.
- Zhu, S.; Li, C.; Loy, C.-C.; and Tang, X. 2016a. Unconstrained Face Alignment via Cascaded Compositional Learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhu, X.; Lei, Z.; Liu, X.; Shi, H.; and Li, S. Z. 2016b. Face Alignment Across Large Poses: A 3D Solution. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhu, X.; and Ramanan, D. 2012. Face detection, pose estimation, and landmark localization in the wild. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, 2879–2886. doi:10.1109/CVPR.2012.6248014.