

# Human and Sheep Facial Landmarks Localisation by Triplet Interpolated Features

Heng Yang\*, Renqiao Zhang\* and Peter Robinson  
University of Cambridge

{hy306, rz264, pr10}@cam.ac.uk

## Abstract

*In this paper we present a method for localisation of facial landmarks on human and sheep. We introduce a new feature extraction scheme called triplet-interpolated feature used at each iteration of the cascaded shape regression framework. It is able to extract features from similar semantic location given an estimated shape, even when head pose variations are large and the facial landmarks are very sparsely distributed. Furthermore, we study the impact of training data imbalance on model performance and propose a training sample augmentation scheme that produces more initialisations for training samples from the minority. More specifically, the augmentation number for a training sample is made to be negatively correlated to the value of the fitted probability density function at the sample's position. We evaluate the proposed scheme on both human and sheep facial landmarks localisation. On the benchmark 300w human face dataset, we demonstrate the benefits of our proposed methods and show very competitive performance when comparing to other methods. On a newly created sheep face dataset, we get very good performance despite the fact that we only have a limited number of training samples and a set of sparse landmarks are annotated. Source code is available for academic use <sup>1</sup>.*

## 1. Introduction

Many computer vision applications require localisation of a set of landmarks for the purpose of fine-grained recognition. For example, joint localisation in human pose estimation [34], part localisation for bird [6] and dog [25] breed recognition. It is of interest to localise facial landmarks for animals and humans, given the fact that their faces hold rich information such as identity, expression, health conditions, etc. In this paper, we are interested in localising sheep and human facial landmarks for real applications.

\* indicates authors contribute equally.

<sup>1</sup><https://github.com/ChrisYang/TIFfacealignment>

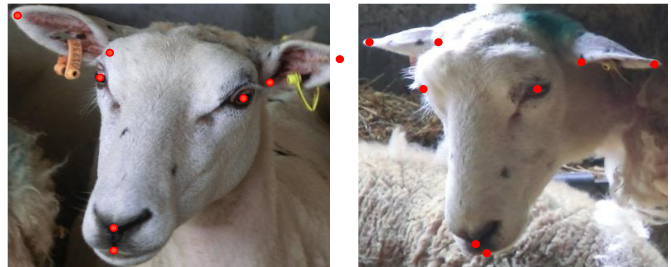


Figure 1: Normal sheep (left) vs. sheep in pain (right). The red landmarks are associated with distinguishable patterns that we intend to localise.

Sheep facial landmark localisation is new in computer vision field and has very promising potential in animal welfare. Compared to other animals, sheep have less intricate facial muscles and thus do not appear to have a wide array of facial expressions. However, researchers have linked a few specific postures with emotional experiences, for example backward ear posture, which is associated with unfamiliar and uncontrollable unpleasant situations, could express fear [7]. Identifying the pain or suffering of animal (like sheep) is an essential aspect of animal welfare and is very helpful to both researchers and farmers. As an example shown in Fig.1, the sheep on the right is suffering heavily from painfulness while the sheep on the left is in a normal condition. Experts on animal welfare research are able to pick up several distinguishable patterns of sheep-in-pain such as orbital tightening, abnormal ear position and abnormal nostril and philtrum shape. In order to identify those features automatically, localising the corresponding landmarks on sheep face is very essential, which is conceptually very similar to human facial landmarks localisation (also called face alignment). As a classical problem in computer vision, face alignment has been intensively studied in the past decades due to its wide applications for example face recognition, facial expression recognition, avatar animation, etc. Several recent methods such as [9, 8, 22, 31, 41, 39, 45] have reported close-to-human performance on the academic

databases such as LFW [21], LFPW [5] and HELEN [24].

However, we meet several obstacles when we apply the state of the art algorithms directly to real data, for both human and sheep facial landmark localisation. First, unlike the benchmark dataset for human face alignment, in which a large number of landmarks are often annotated, the number of facial landmarks in practice is usually smaller, due to the annotation cost and fewer landmarks of interest. Second, both human face and sheep face show big head pose variations in real world given the uncontrollability. It usually results in localisation failures.

In this paper we deal with the problems mentioned above. We build our localisation algorithm on top of the Cascaded Pose Regression (CPR) framework, given its good performance in facial landmarks localisation in the wild. There has been a series of works with incremental improvement one after the other including [8, 9, 16]. The most recent work RCPR ([8]) introduced interpolated shape-indexed features used in each regression. It demonstrated better robustness against large pose variations and shape deformations, compared to the closest landmark indexed feature in [9]. However, the two-point-interpolation method limits the feature extraction space, especially when the number of facial landmarks is small. Landmark sparsity is often the case when we need to annotate a new training dataset given a limited amount of time or only a small number of landmarks needed. To overcome those issues we make the following contributions.

- We propose a new feature extraction scheme, called triplet-interpolation feature (TIF) for cascaded pose regression. It uses three anchor landmarks to calculate a shape-indexed feature. It is more robust to large head pose variation and shape deformation. More importantly, with this scheme, features can be extracted from the facial area with no restriction.
- We propose an augmentation scheme for training sample to deal with the issue of imbalanced training data distribution. This scheme sets the augmentation number of each training sample to be negatively correlated to its value in the probability density function of the training data. More intuitively, we augment the minority training samples with more random initialisations and vice versa.

We have carried out experiments on both human and sheep facial landmarks localisation and demonstrate the benefits of our proposed methods under the situation of sparse landmarks and large head pose variations. It also shows competitive overall performance comparing to other related methods.

The remainder of the paper is organized as follows. In section 2 we present related work. Then we introduce the

triplet-interpolation features and the augmentation scheme in section 3. In section 4 we evaluate our proposed methods on both human and sheep facial landmarks localisation and in section 5 we draw some useful conclusions.

## 2. Related work

### 2.1. Facial landmarks localisation

Facial landmarks localisation has made considerable progress in recent years and a large number of methods have been proposed. Two types of source information are usually used: facial appearance and shape information. Based on whether a method has an explicit detection model for an individual landmark or not, we categorise them into local-based methods and holistic-based methods. The methods in the former category usually rely on explicit discriminative local detection and usually use deformable shape models to regularise the local outputs while the methods in the latter category directly regress the shape (the representation of the facial landmark locations) in a holistic way.

Local based methods usually consist of two parts: local experts and spatial shape models. The former describes how image around each facial landmark looks like in terms of local intensity or colour patterns while the latter describes how face shape varies. There are three main types of local feature detection. (1) Classification methods include Support Vector Machine (SVM) classifier [30, 5] based on various image features such as Gabor [38], SIFT [27], Discriminative Response Map Fitting (**DRMF**) by dictionary learning [2] and multichannel correlation filter responses [18]. (2) Regression-based approaches include Support Vector Regressors (SVRs)[28] with a probabilistic MRF-based shape model, Continuous Conditional Neural Fields (**CCNF**)[4]. (3) Voting-based approaches are also introduced in recent years, including regression forests based voting methods [12, 14, 43] and exemplar based voting methods [35, 33]. One typical shape model is the Constrained Local Model (CLM) [13]. There are some other shape models such as RANSAC in [5], graph-matching in [46], Gaussian Newton Deformable Part Model (**GNDPM**) [37] and mixture of trees [48].

Holistic methods have gained higher popularity in recent years. Most of them work in a cascaded way similar to the classical Active Appearance Model (AAM) [11]. We list very recent holistic methods as well as their properties in Table 1. These methods work in a similar cascaded framework but differ from each other mainly in three aspects. First, how to set up the initialisations; Second, how to calculate the shape-indexed features; Third, what type of regressor is applied at each iteration. Feature extraction and regression are usually interdependent. Core aspects are discussed in [29]. As can be seen, several methods have investigated using simple pixel difference (diff.) features that

Table 1: Holistic methods and their properties.

| Methods   | RCPR [8]     | ESR [3]      | LBF [31]         | TREES [22]   | SDM [39] | TCDCN [45]      |
|-----------|--------------|--------------|------------------|--------------|----------|-----------------|
| features  | pixel diff.  | pixel diff.  | forest on pixels | pixel        | SIFT     | ConvNet feature |
| regressor | random ferns | random ferns | linear           | random trees | linear   | ConvNet         |

is calculated from the current shape. Random ferns and random trees are widely used for regression. Using raw pixel difference feature makes the algorithm very efficient. In our testing, the method ESR, RCPR, LBF and TREES with c++ implementation process a standard face image in milliseconds on an i7 desktop with a single core. This is a great advantage in systems that are designed to process a large number of faces, for example to analyse a group of sheep at the same time. SDM has been widely applied given its good performance of the publicly available model. It runs at around 30 frames per second. TCDCN has applied deep learning approach for face alignment by multi-task learning, but training such a model usually requires a big dataset with multiple additional annotations such as head pose, w/o glasses, etc.

There are several other approaches for holistic face alignment such as occlusion detection based methods by [19, 41], combined local and holistic method in [1], SDM variants including the global SDM [40] and shape searching in [47]. Due to their different setting and limited space, we will not compare them in our experiments. Please refer to [42] for a comprehensive study of recent face alignment methods.

## 2.2. Data imbalance

The data imbalance problem is of particular importance in real world scenarios as the available data usually follows a long tail distribution. Data imbalance has been widely studied in classification problems, i.e., a few classes are abundant while others only have a limited number of samples [20]. State of the art solutions include sampling methods (e.g. under-sampling [26] and SMOTE over-sampling [10]), cost-sensitive learning [17, 20]. On the contrary, very little attention has been paid on data imbalance in regression problem (like our facial landmark localisation). This is mainly due to the fact that the data imbalance is difficult to be noticed given the continuity and the usually high dimensionality of the output space. Thus in this paper we investigate how to adapt the approach of tackling class imbalance to regression problem.

## 3. Method

In this section, we first briefly review the general cascaded pose regression (CPR) approach, on which our localisation algorithm has been built. Then we introduce the triplet-interpolated features. Following that, inverse propor-

tional augmentation is discussed in details as an approach to deal with imbalanced training data.

### 3.1. General CPR and RCPR

The shape of a human or sheep face is represented as a vector of landmark locations, i.e.,  $S = (y_1, \dots, y_k, \dots, y_K) \in \mathbf{R}^{2K}$ , where  $K$  is the number of landmarks.  $y_k \in \mathbf{R}^2$  is the 2D coordinates of the  $k$ -th landmark. CPR is formed by a cascade of  $T$  regressors,  $R^{1 \dots T}$ . Shape estimation starts from an initial shape  $S^0$  and progressively refines the pose. Each regressor refines the pose by producing an update,  $\Delta S$ , which is added up to the current shape estimate, that is,

$$S^t = S^{t-1} + \Delta S. \quad (1)$$

The update  $\Delta S$  is returned by the regressor that takes the previous pose estimation and the image feature  $I$  as inputs:

$$\Delta S = R^t(S^{t-1}, I) \quad (2)$$

The CPR is summarized in Algorithm 1 [16]. This CPR framework differs from the classic boosted approaches mainly in the feature re-sampling process. More specifically, instead of using the fixed features, the input feature for regressor  $R^t$  is calculated relative to the current pose estimation, thus in turn introduces geometric invariance into the cascade process and shows good performance in practice. This is often referred as pose-indexed features as in [16]. The idea of sampling features from current pose estimation is later used in [9, 22]. To strengthen the geometric invariance, instead of extracting features from the closest landmarks, RCPR [8] utilizes a different feature-indexing method ( $h^t(I, S^{t-1})$ ), namely the interpolated shape-indexed features. The features are extracted with reference to two shape points. [8] has proven that RCPR is more robust to large pose variations than the general CPR.

---

#### Algorithm 1 Cascaded Pose (shape) Regression

---

**Require:** Image  $I$ , initial pose  $S^0$

**Ensure:** Estimated pose  $S^T$

- 1: **for**  $t=1$  to  $T$  **do**
  - 2:      $f^t = h^t(I, S^{t-1})$      ▷ Shape-indexed features
  - 3:      $\Delta S = R^t(f^t)$      ▷ Apply regressor  $R^t$
  - 4:      $S^t = S^{t-1} + \Delta S$      ▷ update shape
  - 5: **end for**
-

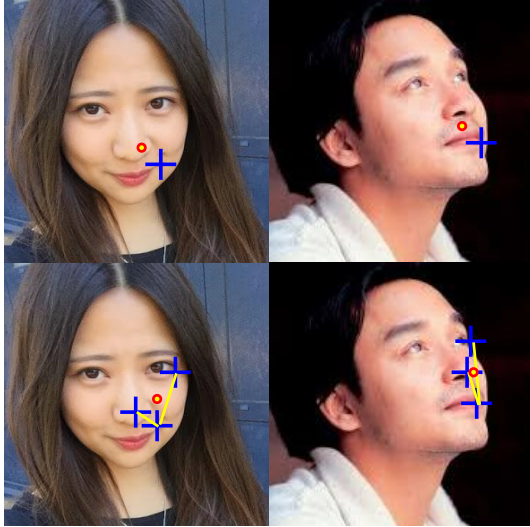


Figure 2: Pixels indexed by the same local coordinates should have the same semantic meaning. The triplet-interpolated feature shows its feature invariance to large pose variation in the right bottom figure.

### 3.2. Triplet-Interpolated Feature (TIF)

The above CPR scheme and its variants are very popular given its high computational efficiency and localisation accuracy. In each iteration, random ferns or random forests takes raw pixel values as input features, which in turn become essential to fast convergence in the cascaded learning. Prevalent pixel-indexing features intend to be invariant with respect to pose variation. That is to say, the indexed pixels referencing to same shape points are expected to have same semantic meaning across different samples. Such efforts have been made in [9], which applied shape-indexed features, and in [8], which achieved stronger geometric invariance with the interpolated shape-indexed features.

However, the interpolated shape-indexed features in RCPR has a fundamental drawback. It can only draw features that are lying on the line segment between two landmarks. As example shows in Fig 3a, features can be extracted from a rich area of the face when the landmarks are dense. However it becomes problematic when the facial landmarks are sparse. Features can only be extracted from very restricted locations (see Fig. 3b). This limits the randomness of feature extraction.

To combine the benefits of geometric invariance and avoid its limitations, we propose a new indexing approach, namely Triplet-interpolated feature(TIF), as shown in 3c. The indexing process works in the following way: Out of every group of three randomly selected landmarks, one is randomly chosen and assigned as the primary point. Then two vectors, from the primary to the rest two, can span the whole plane by linear combination. By setting the param-

eters of the linear combination, a position can be selected within the spanned area, as shown in Fig 3c. The location of the point  $p$  indexed by TIF is represented as:

$$p(S, i, j, k, \alpha, \beta) = y_i + (\alpha \cdot \vec{v}_{ij} + \beta \cdot \vec{v}_{ik}) \quad (3)$$

where  $S$  is the current shape and  $i, j, k$  are landmark indexes.  $\vec{v}_{ij} = y_j - y_i$  is the vector from the position of  $i$  to the position of  $j$ .  $\alpha$  and  $\beta$  are the random ratios that control the position of the indexed point. Compared to the original closest landmark indexed feature in [9], the TIF has two main advantages: 1) it is computationally cheaper since it does not have the shape transformation step; 2) it is more robust to large head pose variation given the triplet interpolation property, as shown in Fig. 2. Compared to the two-point-interpolated-feature in RCPR [8], it is able to extract features from a much wider range, especially when the landmarks are sparse. We will show the benefits of using TIF in the experiment section. Apart from the feature extraction process, we follow the cascaded pose regression process used by ESR [9] and RCPR [8]. Note that in this paper, we only use the feature extraction part of RCPR as the occlusion estimation part requires landmark-wise occlusion annotation. In this way we also make the benefits of feature extraction clearer.

### 3.3. Negatively Correlated Augmentation (NCA)

Before introducing our data augmentation scheme, we first analyse the data distribution of the benchmark database for human facial landmark localisation, i.e., 300w, which is a benchmark database for human facial landmark localisation. It consists of face images from AFW [48], HELEN [36], LFPW [5] and the newly annotated iBug [32]. We partition it to 3148 training images and 689 test images. Training images are from AFW (337 images), HELEN training set (2000 images) and LFPW training set (811 images), and test images are from HELEN test set (330 images), LFPW test set (224 images) and iBug (135 images).

Because it is impractical to analyse the data distribution directly on the output space given its high dimensionality, we ignore individual face difference and small facial deformation. Then facial landmarks distribution is mainly affected by head pose variations, which lie in low dimensional manifold. Therefore, we analyse the distribution of head poses. Since head pose is not provided by the database, estimated head pose information for each face is derived from the annotated facial landmarks. To this end, we fit a mean 3D model (68 facial points) of a head to the annotated points in the image. Then we feed the set of corresponding 3D and 2D points to the POSIT [15] algorithm which produces the head pose information.

As shown in Fig. 4, the majority of training samples distribute near frontal angles. More than 97% of the samples lie within roll angle range between  $-20^\circ$  and  $20^\circ$ . For pitch

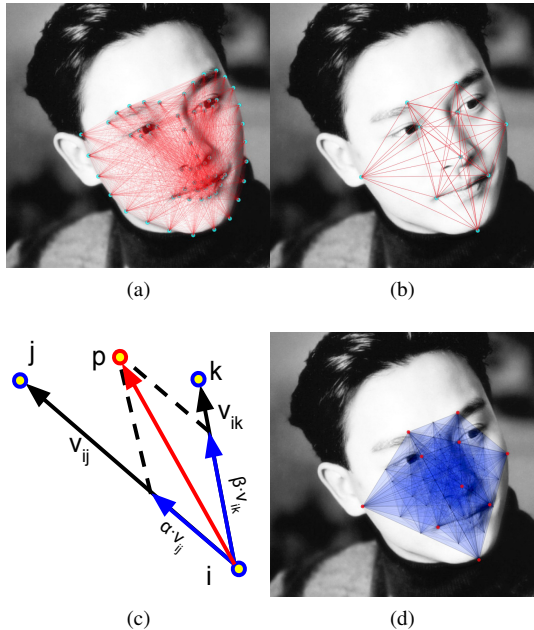


Figure 3: The red lines in (a)(b) show the available area for feature extraction when we use the linear-interpolated shape-index features. (c) and (d) illustrate the concept of our Triplet-interpolated features and its available feature region. (b)(d) together show that how the new indexing method extends the available area for feature extraction when the shape annotation is sparse.

and yaw angle, such percentages are 83% and 76% respectively. For each training sample, we calculate the most significant rotation angle, i.e., the angle with the biggest absolute value. Then we fit a Gaussian curve on all the training samples as shown in Fig. 4d.

We ran several models on the test images including the Explicit Shape Regression (ESR) [9], the Robust Cascaded Pose Regression (RCPR) [8], the Supervised Descent Method (SDM) [39], and the TCDCN [45]. Then we recorded their failures, i.e. a sample with mean localisation error bigger than 0.1 inter-ocular-distance (IOD). The overall distribution is shown in Fig. 4e. Despite these methods being modelled in very different ways, their failures are quite similar. Only a few failures are within angle range between  $-20^\circ$  and  $20^\circ$ , where the majority of training samples distribute. To this end, we can conclude that the imbalanced distribution of training data has heavy impact on testing performance, regardless of the algorithm design.

In the framework of cascaded shape regression, data augmentation is usually carried out during training time. More specifically, for one face image sample, several initialisations are generated by Monte Carlo method. This procedure has been used in ESR, RCRP and SDM and the augmenta-

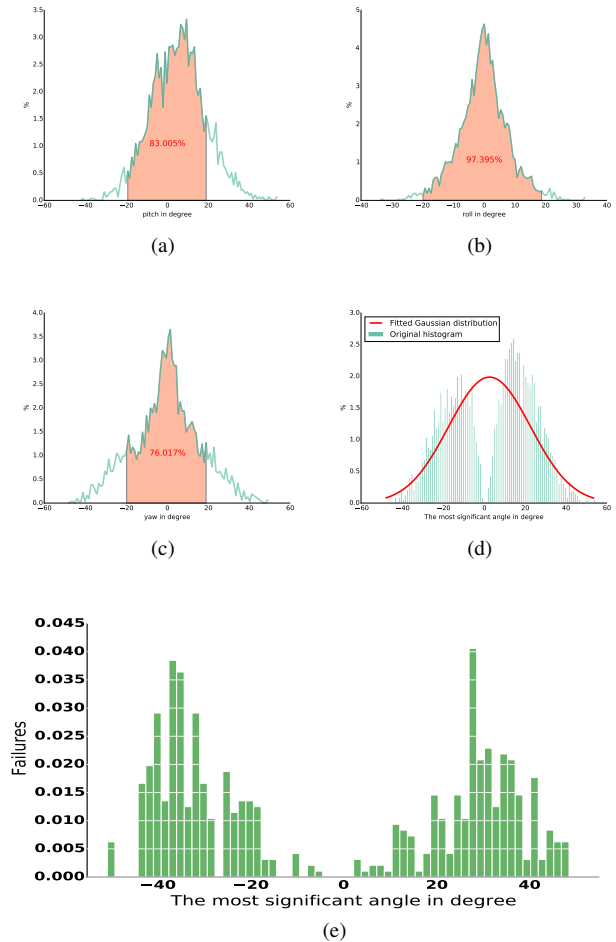


Figure 4: (a) (b) (c) show the histogram of the head pose pitch, roll and yaw angles respectively in 300w training set. (d) shows the fitted Gaussian curve and the histogram of the most significant angle of the training samples. (e) is the histogram of failures from several state of the art models trained on 300W (ESR, TREES, RCPR and SDM).

tion number is usually fixed. We propose a simple augmentation scheme, under which the amount of augmentation of each training sample is negatively correlated to the value on the fitted Gaussian curve (Fig. 4d). Conceptually, this is similar to over-sampling in classification problem but each augmented sample becomes unique in our case because of the initialisation difference. More specifically, the augmentation number  $m_x$  of training sample  $x$  is calculated as:

$$m_x = a \cdot \mathcal{N}(x_{\text{pose}}) + b \quad (4)$$

where  $x_{\text{pose}}$  the head pose of  $x$ ,  $\mathcal{N}(\cdot)$  the fitted Gaussian distribution.  $a$  is a negative variable that controls the slope and  $b$  is a bias term that controls the bounds of augmentation numbers. We use two pairs of values (the maximum and the

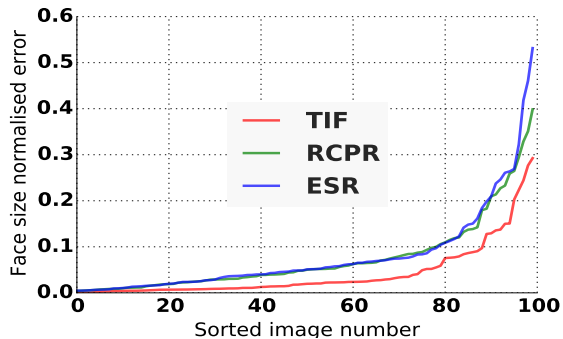


Figure 5: Performance comparison on our sheep face dataset. (Lower is better.)

minimum) to fit this linear equation with a constrain that the total number after augmentation is equal to the baseline augmentation scheme.

## 4. Evaluation

### 4.1. Sheep face experiment

We collected 600 sheep face photos from an animal research centre. We manually labelled the bounding boxes and 8 landmarks on faces as shown in Fig.1. We trained a structured SVM sheep face detector based on HOG features using dlib [23]. Using a few hundred sheep face images is sufficient to train a sheep face detector which can be used in real videos. In our sheep facial landmarks localisation, as usual, we assumed the face bounding boxes are available. We randomly split the 600 sheep faces into a training set (500) and a testing set (100). Then we trained our TIF model, ESR and RCPR using the same training set. We set the augmentation number to 20 for all these methods. We repeated this random process for 5 times, and recorded all the results. Since our test set is not big, we directly report the sorted sample-wise mean error (normalised by sheep face size) of the 100 images. For each index, the value is the average over 5 runs. As can be seen, on a small dataset with sparse landmarks, our method outperforms the baseline methods by a large margin. Around 90% of the sheep images are localised with mean error less than 10% of the face size. Some example images with comparison to other methods are shown in Fig. 6. The sheep face image in our collected dataset exhibits a wide range of diversity: sheep breed, facial colour, lighting condition, background, occlusion, head pose, ear posture, etc.

### 4.2. Human face experiment

In order to further evaluate the proposed schemes, we carry out experiments on human face alignment benchmark database, i.e., 300w. Recall that we split the publicly available database into training set (3148 images) and testing set

(689 images). We have implemented and trained the baseline models (ESR and RCPR) on the same training images. Note that when implementing the RCPR algorithm we only used their method of feature indexing (interpolation by two landmarks) but not their occlusion modelling since there is no occlusion annotation for training. Thus for ESR, RCPR and our TIF method, the only difference is their feature extraction step. During testing time, we also initialised them with the same random shapes for a fair comparison. We carried out two groups of experiments. In the first group, the model was trained on 68 facial landmarks, and in the second group, we only used very sparse landmarks, to simulate the case of the sheep facial landmarks localisation. The mark-up of sparse landmarks is shown 3d, which distribute almost uni-formally among the original 68 landmarks on the face. Note that we use the face bounding box detected by dlib face detector [23], followed by manual check for each face image. This is more realistic in practice than using the tight bounding boxes calculated from the annotated facial landmarks. In order to make a fair comparison, we trained our model as well as most competitive models (highlighted in Section 2) including the RCPR, SDM, TREES, ESR, CCNF, LBF, with the same setting. More specifically, we use the same training set and the same bounding box definition. For TCDCN, GNDPM, DRMF we use their initial trained models as their performance is less competitive.

As shown in Fig. 7a, 1) Our method (NCA + TIF) gets the best performance despite the improvement over the baseline RCPR method is not huge; 2) Only using TIF does not show superior performance over RCPR on dense landmarks setting, which is as expected. The benefit of using TIF is more clear on sparse landmarks, as shown in Fig. 7b. Our proposed TIF improves the baseline RCPR method as well as the similar ESR method by a large margin. Note that, there are some tricks that are able to make the cascaded pose regression methods more robust such as the smart-restart in [8] and the mirrorability based restart in [44], which are naturally compatible to our TIF method as well. In this evaluation we are more interested in the benefits brought by the TIF.

We evaluate the NCA scheme in three methods, our proposed TIF, the RCPR and ESR, since they use the same way of data augmentation. We set the smallest augmentation number to 11 and the biggest to 40 for the training samples in our NCA method, which makes the total number equal to  $20N$ , where  $N$  is the number of training samples, 20 is the augmentation number used by the baseline methods. In this evaluation, we are more concerned with test samples with big head pose variations. Therefore, we record the successful localisation rates (SLR), i.e. the percentage of test samples are with mean localisation error smaller than  $0.1IOD$ . As shown in Fig.8, the proposed NCA scheme is able to improve the SLR effectively. Among the 689 test samples, it is



Figure 6: Landmarks localisation on example sheep face images. From top to bottom show the result of our TIF method, RCPR and ESR respectively. The final column shows a failure example of our method.

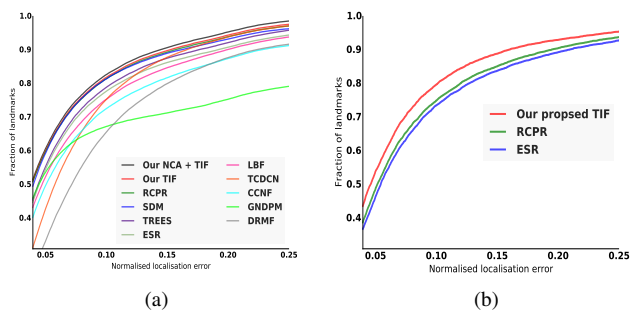


Figure 7: Performance comparison on dense (left) and sparse (right) facial landmarks. (Higher is better.)

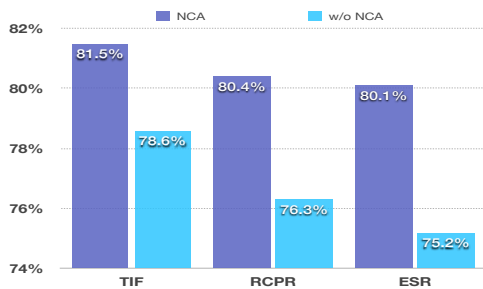


Figure 8: Successful-localisation-rate comparison for methods with and w/o our proposed Negatively Correlated Augmentation (NCA).

able to successfully localise more than around 30 samples. This is very significant given the fact that the failures from methods without NCA are already very difficult.

## 5. Conclusion and discussion

In this paper, we have addressed the problems of localising key landmarks on sheep and human faces. We proposed a new feature extraction scheme by triplet interpolation (TIF), which is more effective under the conditions of large head pose variation and landmark sparsity. On our new sheep face dataset of only 600 images, our proposed method works considerably well on a large diversity of sheep faces. We also studied the issue of training data imbalance and proposed a sample augmentation strategy to improve the performance on test samples that have big variations.

Though we have pushed forward the state of the art method for facial landmarks localisation and decreased the failures, there are still failures that are mainly caused by head pose variation or heavy occlusions. It is an open question whether we need to address these challenges explicitly or provide more data similar to the failure cases. Regarding the sheep face analysis, we have only localised the landmarks of interest, there are still many problems to tackle in order to build an automatic computer vision system to identify the sheep in pain. We believe these are all interesting and valuable problems for both computer vision and animal welfare community.

## References

- [1] J. Alabort-i Medina and S. Zafeiriou. Unifying holistic and parts-based deformable model fitting. In *CVPR*, pages 3679–3688, 2015.
- [2] A. Asthana, S. Zafeiriou, S. Cheng, and M. Pantic. Robust discriminative response map fitting with constrained local models. In *CVPR*, pages 3444–3451, 2013.
- [3] A. Asthana, S. Zafeiriou, S. Cheng, and M. Pantic. Incremental face alignment in the wild. In *CVPR*, pages 1859–1866, 2014.
- [4] T. Baltrušaitis, P. Robinson, and L.-P. Morency. Continuous conditional neural fields for structured regression. In *ECCV*, pages 593–608. Springer, 2014.
- [5] P. Belhumeur, D. Jacobs, D. Kriegman, and N. Kumar. Localizing parts of faces using a consensus of exemplars. In *CVPR*, pages 545–552, 2011.
- [6] T. Berg, J. Liu, S. W. Lee, M. L. Alexander, D. W. Jacobs, and P. N. Belhumeur. Birdsnap: Large-scale fine-grained visual categorization of birds. In *CVPR*, pages 2019–2026. IEEE, 2014.
- [7] A. Boissy, A. Aubert, L. Désiré, L. Greiveldinger, E. Delval, and I. Veissier. Cognitive sciences to relate ear postures to emotions in sheep. *Animal Welfare*, 20(1):47, 2011.
- [8] X. P. Burgos-Artizzu, P. Perona, and P. Dollár. Robust face landmark estimation under occlusion. In *ICCV*, pages 1513–1520, 2013.
- [9] X. Cao, Y. Wei, F. Wen, and J. Sun. Face alignment by explicit shape regression. In *CVPR*, pages 177–190. Springer, 2012.
- [10] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, pages 321–357, 2002.
- [11] T. Cootes, G. Edwards, and C. Taylor. Active appearance models. *IEEE Trans. Pattern Anal. Mach. Intell.*, 23(6):681–685, 2001.
- [12] T. Cootes, M. C. Ionita, and S. P. Robust and accurate shape model fitting using random forest regression voting. In *ECCV*, pages 278–291. Springer, 2012.
- [13] D. Cristinacce and T. Cootes. Feature detection and tracking with constrained local models. In *Proc. Brit. Mach. Vis. Conf.*, volume 2, page 6, 2006.
- [14] M. Dantone, J. Gall, G. Fanelli, and L. Van Gool. Real-time facial feature detection using conditional regression forests. In *CVPR*, pages 2578–2585, 2012.
- [15] D. F. DeMenthon and L. S. Davis. Model-based object pose in 25 lines of code. In *ECCV*, pages 335–343. Springer, 1992.
- [16] P. Dollár, P. Welinder, and P. Perona. Cascaded pose regression. In *CVPR*, pages 1078–1085, 2010.
- [17] C. Elkan. The foundations of cost-sensitive learning. Citeseer.
- [18] H. K. Galoogahi, T. Sim, and S. Lucey. Multi-channel correlation filters. In *ICCV*, pages 3072–3079, 2013.
- [19] G. Ghiasi and C. C. Fowlkes. Occlusion coherence: Localizing occluded faces with a hierarchical deformable part model. In *CVPR*, pages 1899–1906, 2014.
- [20] H. He, E. Garcia, et al. Learning from imbalanced data. *Knowledge and Data Engineering, IEEE Transactions on*, 21(9):1263–1284, 2009.
- [21] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled Faces in the Wild: A Database for Studying Face Recognition in Unconstrained Environments. Technical Report 07-49, University of Massachusetts, Amherst, Oct. 2007.
- [22] V. Kazemi and J. Sullivan. One millisecond face alignment with an ensemble of regression trees. In *CVPR*, pages 1867–1874. IEEE, 2014.
- [23] D. E. King. Dlib-ml: A machine learning toolkit. *Journal of Machine Learning Research*, 10:1755–1758, 2009.
- [24] V. Le, J. Brandt, Z. Lin, L. Bourdev, and T. S. Huang. Enhanced pictorial structures for precise eye localization under uncontrolled conditions. In *ECCV*, pages 1621–1628, 2012.
- [25] J. Liu, A. Kanazawa, D. Jacobs, and P. Belhumeur. Dog breed classification using part localization. In *ECCV*, pages 172–185. Springer, 2012.
- [26] X.-Y. Liu, J. Wu, and Z.-H. Zhou. Exploratory undersampling for class-imbalance learning. *Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on*, 39(2):539–550, 2009.
- [27] D. G. Lowe. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vis.*, 60(2):91–110, 2004.
- [28] B. Martinez, M. Valstar, X. Binefa, and M. Pantic. Local Evidence Aggregation for Regression Based Facial Point Detection. *IEEE Trans. Pattern Anal. Mach. Intell.*, pages 1149–1163, 2012.
- [29] C. Qu, H. Gao, E. Monari, J. Beyerer, and J.-P. Thiran. Towards robust cascaded regression for face alignment in the wild. In *CVPRW*, 2015.
- [30] V. Rapp, T. Senechal, K. Bailly, and L. Prevost. Multiple kernel learning svm and statistical validation for facial landmark detection. In *Proc. IEEE Int'l Conf. on Autom. Face Gesture Recognit.*, pages 265–271, 2011.
- [31] S. Ren, X. Cao, Y. Wei, and J. Sun. Face alignment at 3000 fps via regressing local binary features. In *CVPR*, pages 1685–1692, 2014.
- [32] C. Sagonas, G. Tzimiropoulos, S. Zafeiriou, and M. Pantic. 300 faces in-the-wild challenge: The first facial landmark localization challenge. In *ICCV*, pages 397–403, 2013.
- [33] X. Shen, Z. Lin, J. Brandt, and Y. Wu. Detecting and aligning faces by image retrieval. In *CVPR*, pages 3460–3467, 2013.
- [34] J. Shotton, T. Sharp, A. Kipman, A. Fitzgibbon, M. Finocchio, A. Blake, M. Cook, and R. Moore. Real-time human pose recognition in parts from single depth images. *Communications of the ACM*, 56(1):116–124, 2013.
- [35] B. M. Smith, J. Brandt, Z. Lin, and L. Zhang. Nonparametric context modeling of local appearance for pose-and expression-robust facial landmark localization. In *CVPR*, pages 1741–1748, 2014.
- [36] X. Tan, F. Song, Z. H. Zhou, and S. Chen. Enhanced pictorial structures for precise eye localization under uncontrolled conditions. In *CVPR*, pages 1621–1628, 2009.
- [37] G. Tzimiropoulos and M. Pantic. Gauss-newton deformable part models for face alignment in-the-wild. In *CVPR*, pages 1851–1858, 2014.
- [38] D. Vukadinovic and M. Pantic. Fully automatic facial feature point detection using gabor feature based boosted classifiers. In *Proc. IEEE Int'l Conf. Systems, Man, and Cybernetics*, pages 1692–1698, 2005.
- [39] X. Xiong and F. De la Torre. Supervised descent method and its applications to face alignment. In *CVPR*, pages 532–539, 2013.
- [40] X. Xiong and F. De la Torre. Global supervised descent method. In *CVPR*, pages 2664–2673, 2015.
- [41] H. Yang, X. He, X. Jia, and I. Patras. Robust face alignment under occlusion via regional predictive power estimation. *IEEE Trans. Image Processing*, 2015.
- [42] H. Yang, X. Jia, C. C. Loy, and P. Robinson. An empirical study of recent face alignment methods. *arXiv preprint arXiv:1511.05049*, 2015.
- [43] H. Yang and I. Patras. Sieving regression forests votes for facial feature detection in the wild. In *Proc. Int'l Conf. Computer Vision*. IEEE, 2013.
- [44] H. Yang and I. Patras. Mirror, mirror on the wall, tell me, is the error small? In *CVPR*. IEEE, 2015.
- [45] Z. Zhang, P. Luo, C. C. Loy, and X. Tang. Facial landmark detection by deep multi-task learning. In *ECCV*, pages 94–108. Springer, 2014.
- [46] F. Zhou, J. Brandt, and Z. Lin. Exemplar-based graph matching for robust facial landmark localization. In *ICCV*, pages 1025–1032, 2013.
- [47] S. Zhu, C. Li, C. C. Loy, and X. Tang. Face alignment by coarse-to-fine shape searching. In *CVPR*, 2015.
- [48] X. Zhu and D. Ramanan. Face detection, pose estimation and landmark localization in the wild. In *CVPR*, pages 2879–2886, 2012.