

# Tree-gated Deep Mixture-of-Experts For Pose-robust Face Alignment

Estèphe Arnaud<sup>1</sup>, Arnaud Dapogny<sup>2</sup>, and Kévin Bailly<sup>1, 2</sup>

<sup>1</sup>Sorbonne Université, CNRS, Institut des Systèmes Intelligents et de Robotique, ISIR, F-75005 Paris, France

<sup>2</sup>Datakalab, Paris, France

## Abstract

Face alignment consists of aligning a shape model on a face image. It is an active domain in computer vision as it is a preprocessing for a number of face analysis and synthesis applications. Current state-of-the-art methods already perform well on "easy" datasets, with moderate head pose variations, but may not be robust for "in-the-wild" data with poses up to 90°. In order to increase robustness to an ensemble of factors of variations (e.g. head pose or occlusions), a given layer (e.g. a regressor or an upstream CNN layer) can be replaced by a Mixture of Experts (MoE) layer that uses an ensemble of experts instead of a single one. The weights of this mixture can be learned as gating functions to jointly learn the experts and the corresponding weights. In this paper, we propose to use tree-structured gates which allows a hierarchical weighting of the experts (Tree-MoE). We investigate the use of Tree-MoE layers in different contexts in the frame of face alignment with cascaded regression, firstly for emphasizing relevant, more specialized feature extractors depending of a high-level semantic information such as head pose (Pose-Tree-MoE), and secondly as an overall more robust regression layer. We perform extensive experiments on several challenging face alignment datasets, demonstrating that our approach outperforms the state-of-the-art methods.

## 1 Introduction

Face alignment refers to the process of localizing a number of landmarks on a face image (lips and eyes corners, pupils, nose tip). This is an important research field of computer vision, as it is an essential preprocess for applications such as face recognition, tracking, expression analysis as well as face synthesis.

Recently, regression-based methods appeared among the most successful ones, achieving high accuracies on images with reasonable variations in head pose, illumination as well as occasional facial expressions or partial occlusions. These methods address the face alignment problem by directly learning a mapping between shape-indexed face texture and

the landmark positions. This regression is often performed in a cascaded manner: starting from an initial guess, a first stage predicts a displacement for every landmark coordinate. This prediction then gets progressively refined through successive regressors that are trained to improve the predictions of the previous stages. In the frame of such a coarse-to-fine strategy, the first cascade stages usually capture large deformations, while the last stages focus on more subtle variations. The work of Xiong *et al.* [24] proposes to use successive linear regressions based on SIFT descriptors extracted around each landmark at the current position of the model. Similarly, Ren *et al.* [15] propose to learn shape-indexed pixel intensity difference features with random forests in order to speed up the feature extraction step.

Deep learning techniques have also been investigated to tackle the face alignment problem. For example, in [19] each cascade stage is modeled using deep convolutional networks (CNNs) in order to jointly learn the representation and regression steps in an end-to-end fashion instead of training the regressor upon handcrafted features. Mnemonic Descent Method [21] improves the feature extraction process by sharing the CNN layers among all cascade stages, and the landmark trajectories through the successive cascade stages are modeled using recurrent neural networks. This results in memory footprint reduction as well as more efficient representation learning and a more optimized cascaded alignment process.

Using deep architectures within the cascaded regression framework allows to achieve high-end alignment precision. However, despite the success of these methods, these are very sensitive to extreme conditions such as large head pose, facial expression, illumination changes, or occlusions induced by objects in front of the face (e.g. glasses, hands, hairs). The appearance of the face can then drastically change and corrupt the input features fed to the displacement regressor in the first stages of the cascades, causing errors that will be hard to overcome later on. In order to address these limitations, ensemble methods can be combined with the use of deep learning techniques: using a committee of expert layers instead of just a single, strong layer improves the diversity

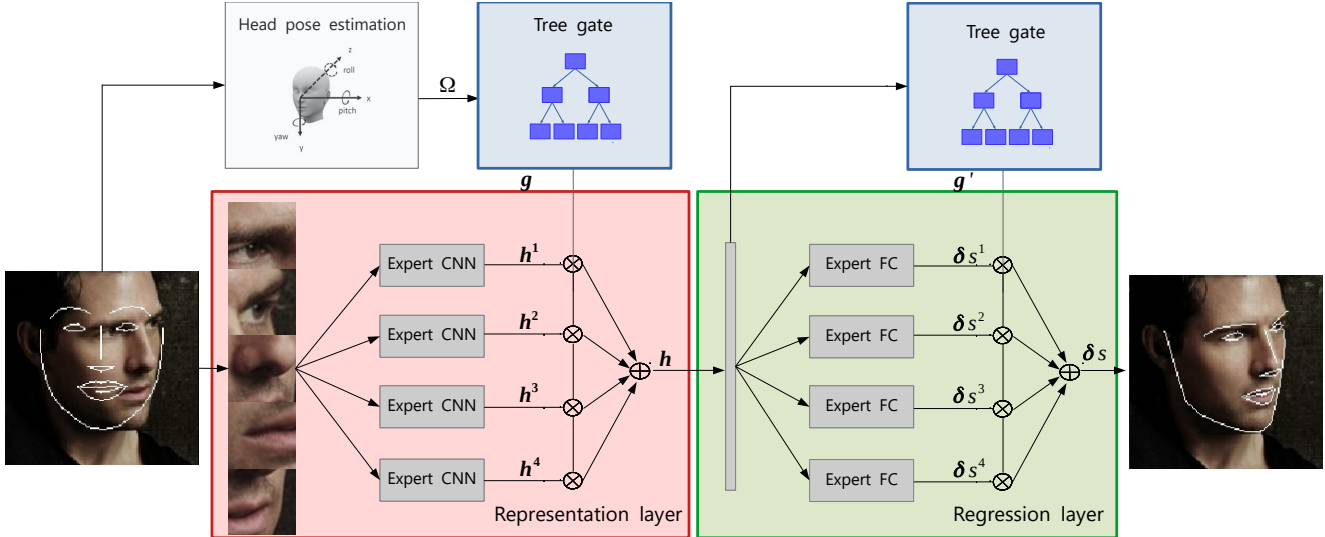


Figure 1: Architecture of the proposed Pose-Tree-MoE model for face alignment. Only one cascade stage is represented for clarity reasons. Patches are extracted around each landmark and fed to an ensemble of expert CNNs. Head pose is estimated and used to weight each expert CNN through a tree-gate. The extracted representations are then used as an input for an ensemble of expert FC layers. Contribution of each expert is, once again, weighted using a tree-gate to predict landmark displacement. These stages are then repeated through all cascade stages, allowing a pose-adaptive, robust face landmark alignment.

of possible responses, which leads to an increased overall robustness. Furthermore, an adaptive combination of the outputs of these expert layers can be learned jointly by the use of gates. Using a tree gate structure allows to learn a hierarchical clustering of the expert layers, which leads to a more efficient selection and therefore specialization of the experts. These gates can be based on either an extracted representation (e.g. when applied for regression), or a high-level semantic information (such as head pose when applied for learning representation layers).

The representation layer, regression layer and the corresponding tree gates can be trained jointly in an end-to-end manner using neural trees. The proposed architecture is summarized on Figure 1. The contributions of this paper are thus three-folds:

- We show that integrating ensemble methods within a deep architecture is beneficial to the overall robustness of face alignment methods. In particular, we use a committee of expert layers instead of a single, strong layer for representation and regression layer respectively. This allows each expert network to be geared towards a specific alignment case.
- We propose an adaptive weighting strategy based on gates that learn to combine the contribution of each expert network. Weights can be estimated from a high level semantic information such as head pose, or di-

rectly from an embedding extracted by the representation layer. In addition, tree-structured gates allow to learn a set of hierarchically clustered expert layers. It can be learned as neural trees [11] to allow joint end-to-end training of expert networks and gates.

- We propose a real-time face alignment system that outperforms state-of-the-art approaches on several databases and is particularly robust to large poses and occlusions.

## 2 Related work

The main current challenge of face alignment problem is the robustness to strong variations, such as large pose or occlusion. A first approach may consist in better conditioning the training in order to increase generalization capability. Landmark localization can be improved by simultaneously learning other related tasks such as attributes detection [27][7]. Indeed, learning to detect attributes such as the presence of glasses on the face can improve the model robustness to occlusions. However, these techniques require additional data, either unlabelled or annotated with auxiliary attributes.

Other models aim to address robustness to only one source of variation. The architecture and the training procedure are then designed to specifically address this variation.

For example, some models are specialized in occlusion handling and explicitly predict the occluded part of the face [2][6][25][26]. However, learning such models generally requires data with occlusion labels, which is a major requirement. Head pose variations can also drastically change the appearance of the face, and some landmarks can be self-occluded. Taking this information into account can improve the face alignment process. Zhu *et al.* [31] proposed to integrate head pose information to condition and adapt Convolutional Neural Networks. Kumar *et al.* [13] proposed a Bayesian formulation in which head pose estimation allows to condition heatmaps extracted to estimate landmarks localization, with the constraint that the estimated face shape must follow a dendritic structure for effective information sharing. These approaches have demonstrated that the head pose information significantly improves the localization performance. However, in this case, head pose is treated as a *post-hoc* multiplicative variable. Conversely, we argue that taking this information into account upstream in the network leads to better representation learning, bringing more robustness, as head pose can dramatically affect face appearance.

Other approaches explicitly seek to build models that are specific to each variation type, and combine them together. Wu *et al.* [22] proposes a global framework, trained in a cascaded manner, which simultaneously performs facial landmarks localization, occlusion detection and head pose estimation with separate modules. Relationships between these allows the modules to benefit from each other. However, each module requires additional annotations in the trainset. By contrast, the proposed method only relies on facial landmarks and head pose information that can be inferred from the landmark positions.

In addition, the architecture of our model combines advantages of ensemble methods with those of deep learning techniques. Such models have already been explored in the deep learning literature: a first approach for combining deep learning and ensemble methods is to craft a differentiable ensemble architecture in order to allow end-to-end parameter learning. Kontshieder *et al.* [11] designed a differentiable deep neural forest, by unifying the divide-and-conquer principle of decision trees (allowing to cluster data hierarchically) with the representation learning from deep convolution networks. Each predictor is a binary tree, whose split nodes contain routing functions, defining the probability of reaching one of the sub-trees. The probabilistic routing functions are differentiable, allowing these neural trees to be integrated into a fully-differentiable system. The forest corresponds to a set of trees, whose final output is the simple average of the output of each tree. Since then, other models have sought to generalize neural forests [20][8], by integrating upstream convolutional or multi-layer perceptron layers within routing functions to learn more complex input partitionings, leading to higher performance. Dapogny *et al.*

[4] used neural forests for face alignment, with promising results. However, their model uses handcrafted SIFT descriptors. In addition, to adapt neural forest for regression purposes, the authors use neural trees whose leaves are fixed and correspond to a sampling of the remaining displacements from the training data (with small variations). Their model then seeks to optimally combine fixed leaves. This training procedure can theoretically lead to rigid responses, reducing the expressiveness of the model.

A second approach may be to parallelize a set of small networks instead of a single large network. The idea of using a set of regressors within an end-to-end system was firstly introduced by Jacobs *et al.* [9] and more recently taken up by Eigen *et al.* [5]. It shows promising results and is well adapted to our problem. Eigen *et al.* design a Mixture-of-Experts (MoE) layer, consisting in jointly learning a set of expert subnetworks with gates, allowing to learn to combine a number of experts depending on the input. In the same vein, Shazeer *et al.* [18] introduce sparsity in MoE in order to save computation and to increase representation capacity.

### 3 Framework overview

In this section, we introduce our Pose-Tree-MoE model for face alignment: first, in Section 3.1 we describe the head pose estimation module. We then detail in Section 3.2 the representation layer, which select relevant experts based on this head pose estimate, and extract features from patches extracted around a current feature point localization. Then, Section 3.3 shows how we predict landmark coordinate displacements from this pose specific representation. In Section 3.4 we detail the architecture of the gates used to weight the contribution of each expert network for both the representation and regression layer respectively. The whole architecture is integrated into a cascaded regression framework 3.5. This section also provides implementation and architectural details to ensure reproducibility of the results.

#### 3.1 Head pose estimation

Following [16] we use a truncated pre-trained ResNet-50 network to extract head pose from the raw face image  $I$ . We note  $\phi$  the embeddings (2048 units) of the last fully-connected layer. A naive approach would consist in using a single deep network  $\Omega(\phi(I))$  that directly predicts the 3 head pose angles, as it was done in [16]. However, as pointed out in [14], sharing all the representations layers may or may not be optimal, depending on the tasks at stake. In our case, we obtained better performance by regressing yaw  $\gamma$ , pitch  $\beta$  and roll  $\alpha$  values with separate networks:

$$\Omega = \Omega_\gamma(\phi) \parallel \Omega_\beta(\phi) \parallel \Omega_\alpha(\phi) \quad (1)$$

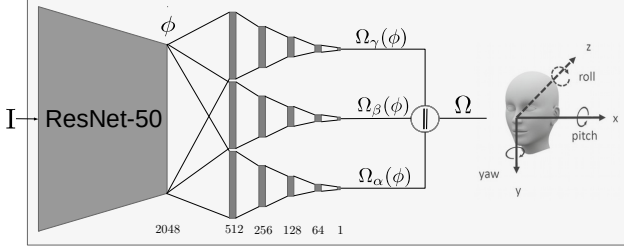


Figure 2: Overview of head pose estimator. Head pose is estimated by three distinct regressors stacked on a common representation extracted by a ResNet-50.

with  $\Omega_\gamma(\phi)$ ,  $\Omega_\beta(\phi)$ ,  $\Omega_\alpha(\phi) \in [-\pi, \pi]^3$  and  $\parallel$  the concatenation operator. Lastly, we also fine-tune the ResNet-50 backbone to improve the head pose regression accuracy.

### 3.2 Representation layer

Let's note  $\mathbf{s} \in \mathbb{R}^{2P}$  an initial shape guess (with  $P$  landmarks). As it is classical in the frame of cascaded regression for face alignment, we extract shape-indexed patches, *i.e.* patches centered at each landmark current localization. If we denote  $\mathbf{h}$  the output of the feature extraction layer, we thus have:

$$\mathbf{h} = f(I, s^1) \parallel \dots \parallel f(I, s^P) \quad (2)$$

As it is ubiquitous among recent face alignment methods, we use CNNs to model function  $f$ . Furthermore, in order to limit the number of parameters, we share the convolution kernels between the different patches, as it was done in [21]. Instead of using one large CNN to extract features around each landmark current location estimation [1], we use a committee of smaller expert CNNs. The idea is that head poses or partial occlusions can dramatically alter the appearance around specific landmarks (e.g. cheeks landmarks in case of large poses) and that we can learn expert CNNs to extract more relevant shape-indexed features in such cases. Specifically, we define  $L$  expert CNN  $\{f^l(I, s^l)\}_{l=1 \dots L}$  and the output of each expert CNN  $\mathbf{h}^l$  as follows:

$$\mathbf{h}^l = f^l(I, s^l) \parallel \dots \parallel f^l(I, s^P) \quad (3)$$

with  $\mathbf{h}^l \in \mathbb{R}^{Pn}$  and  $n$  the number of features per landmark at the output of the last CNN layer. We denote  $\mathbf{H} = [\mathbf{h}^1, \dots, \mathbf{h}^L]$  the responses of all the expert CNNs. More precisely,  $\mathbf{H}$  is a  $Pn \times L$  matrix whose columns are the responses of the  $L$  expert CNNs. Now that we have extracted expert features relatively to the neighbouring of each landmark, we want to aggregate these features: to do so a naive approach would be to simply sum these features, *i.e.* sum over all the columns of  $\mathbf{H}$ . However, a better solution would be to use a high level semantic variable, such as head pose,

to select the most relevant experts based on the output of a gate function  $g : \Omega \in [-\pi, \pi]^3 \mapsto \mathbf{g}(\Omega) \in [0, 1]^L$  such that  $\sum_{l=1}^L g_l(\Omega) = 1$ . In such a case, the output  $\mathbf{h}$  of the representation layer can be written as the sum of the contributions of the  $L$  experts, weighted by the gate value relatively to that expert:

$$\mathbf{h} = \mathbf{H} \cdot \mathbf{g}(\Omega) \quad (4)$$

### 3.3 Regression layer

Given the extracted feature vector  $\mathbf{h} \in \mathbb{R}^{Pn}$ , we now aim at regressing a displacement  $\delta \mathbf{s} \in \mathbb{R}^P$  between  $\mathbf{s}$  and the ground truth landmark localization  $\mathbf{s}^*$ . Such displacement is usually estimated using a single, large deep fully-connected network, *i.e.*  $\delta \mathbf{s} = r(\mathbf{h})$ .

Once again, instead of designing one such large network, we can use a committee of  $L'$  several smaller expert networks  $\{\delta \mathbf{s}^l = r^l(\mathbf{h})\}_{l=1 \dots L'}$ . Let's note  $\delta \mathbf{S}$  a  $\mathbb{R}^{2P \times L'}$  matrix containing all the predictions of these expert regressors:

$$\delta \mathbf{S} = [\delta \mathbf{s}^1, \dots, \delta \mathbf{s}^{L'}] \quad (5)$$

The columns of  $\delta \mathbf{S}$  contain displacements predicted by each expert regressor. More specifically, each displacements predicted by each expert indexed by  $l'$  corresponds to the output of a fully connected layer with ReLU activation:

$$\delta \mathbf{s}^l = w_1^l \cdot \max(0, \mathbf{w}_0^l \cdot \mathbf{h} + b_0^l) + b_1^l \quad (6)$$

with  $\Theta_r^l = \{\mathbf{w}_0^l, b_0^l, w_1^l, b_1^l\}$  the set of parameters of the  $l'$ -th expert. Let's suppose we now have access to the output of a gating network  $g' : \mathbf{h} \in \mathbb{R}^L \mapsto \mathbf{g}' \in [0, 1]^{L'}$  based on the extracted features  $\mathbf{h}$ . In such a case, The output of the regression layer can be written as:

$$\delta \mathbf{s} = \delta \mathbf{S} \cdot \mathbf{g}'(\mathbf{h}) \quad (7)$$

### 3.4 Gating network

In what follows, we note  $\gamma : \mathbf{z} \in \mathbb{Z} \mapsto \mathbf{x} \in \mathbb{X}$  any mapping function and  $g : \mathbf{x} \in \mathbb{X} \mapsto \mathbf{g} \in [0, 1]^L$  a gate function with  $L$  the number of expert networks associated to this gate, such that  $\sum_{l=1}^L g_l(x) = 1$ . In order to learn an adaptive combination of expert networks, two types of gates can be designed.

#### 3.4.1 Softmax gate

The most straightforward way to design a gating function is to use a softmax activation function:

$$\mathbf{g}(\mathbf{x}) = \text{softmax}(\mathbf{x}) = \left( \frac{e^{\mathbf{w}_1 \cdot \mathbf{x} + b_1}}{\sum_{l=1}^L e^{\mathbf{w}_l \cdot \mathbf{x} + b_l}}, \dots, \frac{e^{\mathbf{w}_L \cdot \mathbf{x} + b_L}}{\sum_{l=1}^L e^{\mathbf{w}_l \cdot \mathbf{x} + b_l}} \right) \quad (8)$$

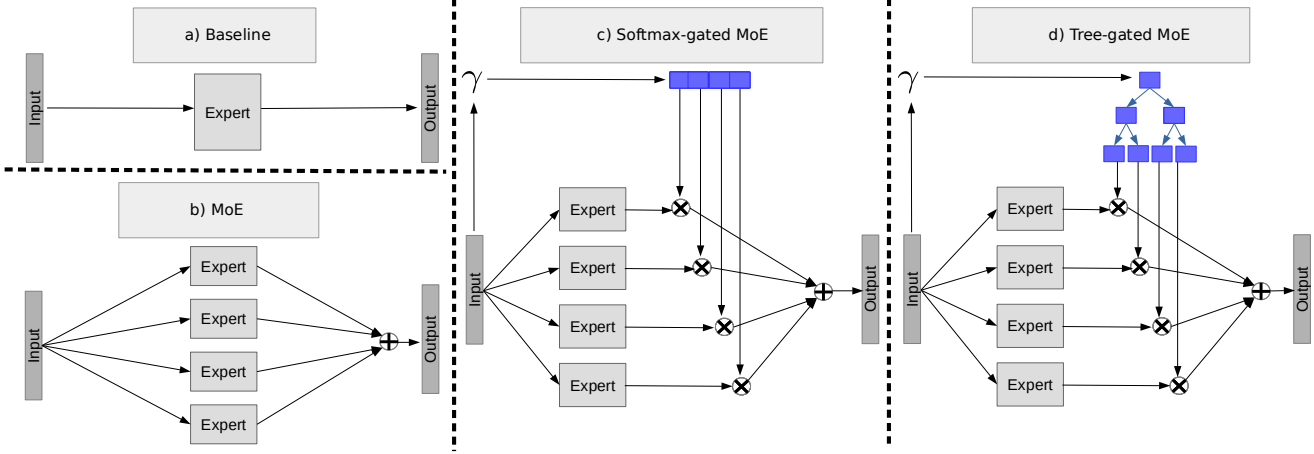


Figure 3: Architecture overview. **a) Baseline:** a baseline, single-expert network. **b) MoE:** the output is the average of  $L$  independent experts. **c) Softmax-gated MoE:** The output of  $L$  expert networks are weighted by the output of a  $L$ -dimensional softmax layer. **d) Tree-gated MoE:** The output of  $L$  expert networks are weighted by the terminal (leaf) probabilities of a  $\log_2(L)$ -deep neural tree. The function  $\gamma$  allows to gate expert networks with a transformation of the input. In our case, for the representation layer,  $\gamma$  is an head pose estimate computed on the whole image. For the regression layer,  $\gamma$  is the identity mapping.

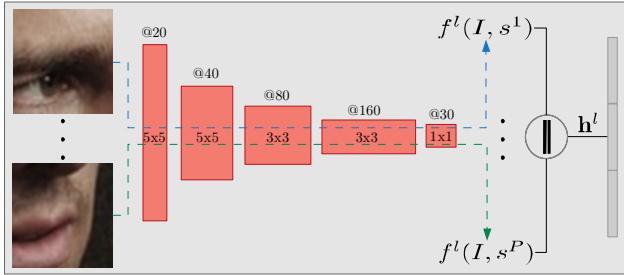


Figure 4: Overview of an expert CNN in representation layer. For each cascade stage, patches are extracted around feature points and a number of shared convolution kernels are applied to each patch. The output features are then concatenated.

with  $\Theta_g = \{\mathbf{w}_l, b_l\}_{l \in \{1, \dots, L\}}$  the parameters of the gate function. While this function is very simple, it doesn't allow to learn hierarchical partitions of  $\mathbb{X}$ .

### 3.4.2 Tree gate

In order to learn a more potent gating network we use a single neural tree. A neural tree [11] is composed of subsequent soft, probabilistic routing functions  $d_n$ , that represents the probability to reach the left child of node  $n$ . Formally,  $d_n$  is defined as a single neuron:

$$d_n(\mathbf{x}) = \sigma(\mathbf{w}_n \cdot \mathbf{x} + b_n) = \frac{e^{\mathbf{w}_n \cdot \mathbf{x} + b_n}}{1 + e^{\mathbf{w}_n \cdot \mathbf{x} + b_n}} \quad (9)$$

with  $\Theta_g = \{\mathbf{w}_n, b_n\}_{n \in \mathcal{N}}$  parameters to learn.

For an input  $\mathbf{x}$ , the probability  $\mu_l$  to reach a leaf  $l \in \{1, \dots, L\}$  is computed as a product of the successive activations  $d_n$  down the whole tree. Therefore:

$$\mu_l(\mathbf{x}) = \prod_{n \in \mathcal{N}} d_n(\mathbf{x})^{\mathbb{1}_{l \leftarrow n}} (1 - d_n(\mathbf{x}))^{\mathbb{1}_{l \rightarrow n}} \quad (10)$$

where  $l \leftarrow n$  is true if  $l$  belongs to the left subtree of node  $n$ , and  $l \rightarrow n$  is true if  $l$  belongs to the right subtree.

We define our tree-gate as the concatenation of the  $2^{\mathcal{D}}$  leaves probabilities of a single neural tree of depth  $\mathcal{D}$ :

$$\mathbf{g}(\mathbf{x}) = \mu_1(\mathbf{x}) \parallel \dots \parallel \mu_L(\mathbf{x}) \quad (11)$$

For extracting representations with a committee of experts layers, a naive solution would be to use the raw image as the gate input, *i.e.* setting  $\gamma = Identity$ . However, in this case, the raw image is too low-level and cannot be used directly, thus it is preferable to use high-level semantic information computed from  $I$ , such as head pose estimation, by setting  $\gamma = \Omega$  (see Section 3.1). By contrast, in the regression layer, the information extracted by the representation layer is already semantically abstract, thus it can directly be used as the gate input, by setting  $\gamma = Identity$ . Last but not least, the feature representation, regression and the gate parameters are all optimized jointly for each cascade stage.

## 3.5 Architectures

Similarly to other cascaded regression approaches, we learn a cascade of mappings  $\mathbf{s}^{(0)} + \sum_k \delta \mathbf{s}^{(k)}$ , where  $\mathbf{s}^{(0)} \in \mathbb{R}^{2P}$

is an initial guess (usually an average shape computed over the whole train set) and each  $\mathbf{s}^{(k)}$  is a displacement estimated by first computing representations from shape-indexed patches centered at the current landmark estimate (as provided by the displacements applied so far) for each expert CNN, and using head pose estimate to weight the expert CNNs according to their relevance *via* the gating function. From these representations, we compute the landmark displacement, once again by using a gated mixture of experts. Then, the landmark localization is updated and the subsequent representation/regression stages can be applied sequentially.

**Regarding the regression layer:** we define several architectures whose differences only lie in the architecture of the regression layer. In such a case, we use a single CNN ( $L = 1, g^l = 1 \forall l$ ) for the representation layer of each cascade stage  $k$ , which is composed of 5 strided convolution layers ( $20@5 \times 5 \rightarrow 40@5 \times 5 \rightarrow 80@3 \times 3 \rightarrow 160@3 \times 3 \rightarrow 30@1 \times 1$ ), which makes 30 features per landmark in the last layer. The output concatenated feature vector  $\mathbf{h}$  is then of size  $30 \times P$  (2040 with  $P = 68$  landmarks). The different architectures proposed for modelling the regression layer are then as follows.

- **Baseline:** a single large FC ( $L' = 1, g^{l'} = 1 \forall l'$ ) with one hidden layer  $2040 \times 8192 \rightarrow 8192 \times 136$
- **Mixture-of-Experts (MoE):** an unweighted combination of small expert FC:  $\mathbf{g}^{l'} = (\frac{1}{L'}, \dots, \frac{1}{L'})$ . We use  $L' = 64$  experts, each containing one hidden layer  $2040 \times 128 \rightarrow 128 \times 136$ .
- **Softmax-gated Mixture-of-Experts (Softmax-MoE):** a committee of  $L' = 64$  expert FC layers, each one containing a single hidden layer  $2040 \times 128 \rightarrow 128 \times 136$  gated by a single 64-dimensional softmax layer.
- **Tree-gated Mixture-of-Experts (Tree-MoE):** a committee of  $L' = 64$  experts, each containing one hidden layer  $2040 \times 128 \rightarrow 128 \times 136$ , gated by neural tree of depth 6 (with  $2^6 = 64$  leaves).

**Regarding the representation layer:** furthermore, instead of using a single deep CNN for learning representations, we use a committee of  $L_h = 8$  expert CNNs, each being composed of 5 strided convolution layers with fewer feature maps but as many features per landmark in the last layer ( $7@5 \times 5 \rightarrow 14@5 \times 5 \rightarrow 28@3 \times 3 \rightarrow 56@3 \times 3 \rightarrow 30@1 \times 1$ ). We use a tree-gate based on the previously computed head pose estimate as the gating network for learning representations and refer to this model as **Pose-Tree-MoE**. We can also use softmax-gate and refer to this model as **Pose-Softmax-MoE**.

**Implementation details:** with this configuration, each model has roughly the same total number of parameters (18 million parameters total for each cascade stage and with representation/regression layers), allowing a fair comparison between the models. Training is done by optimizing a  $\mathcal{L}_2$ -loss with  $150 \times 150$  grayscale images and  $32 \times 32$  patches centered around the 68 landmarks. The images are normalized so that they take values in  $[-1, 1]$ . We train 4 cascade stages and apply data augmentation as it is traditionally done in the literature: for each image we augment the initial mean shape by a random translation factor  $t \sim \mathcal{N}(0, 10)$  and a random scaling factor  $s \sim \mathcal{N}(1, 0.1)$ , and half the time a horizontal flip of the image is performed. The parameters (for the representations and regression layers, as well as the gating function) are optimized jointly in an end-to-end manner by applying ADAM optimizer [28] with a learning rate of 0.001.

## 4 Experiments

In this section, we validate our model both qualitatively and quantitatively. First, in Section 4.1 we present the datasets that we use to train or test the proposed approaches. We validate the architectural choices in Section 4.2 on frontal head poses. We then compare our model with state-of-the-art approaches in Section 4.3 for both 2D and 3D face alignment. In Section 4.4, we qualitatively assess the relevance of the proposed approach both for the representation layer and in the regression layer. Finally, in Section 4.4.4 we evaluate the runtime of our model.

### 4.1 Datasets

We evaluate the effectiveness of the proposed approach both for 2D and 3D face alignment. In the first case, the ground truth landmarks correspond to projections on the visible part of the face. In the latter case, the ground truth corresponds to the real 2D coordinates (without the depth component) of the landmarks, which are often occluded due to large pose variations.

#### 4.1.1 2D Face alignment

The 300W dataset was introduced by the I-BUG team [17] and is considered the benchmark dataset for training and testing face alignment models, with moderate variations in head pose, facial expressions and illuminations. It also embraces a few occluded images. The 300W dataset consists of four datasets: LFPW (811 images for train / 224 images for test), HELEN (2000 images for train / 330 images for test), AFW (337 images for train) and IBUG (135 images for test). As it is classically done in the literature for 2D face alignment, we train our models on a concatenation on

AFW, LFPW and HELEN trainsets, which makes a total of 3148 images for train. For comparison with state-of-the-art methods, we refer to LFPW and HELEN test sets as the common subset and I-BUG as the challenging subset of 300W, as it is commonly done in the literature.

The COFW dataset [2], is an "in-the-wild" dataset containing only occluded data. It is a benchmark dataset to test the robustness of models w.r.t. partial occlusions. COFW contains 500 images for train and 507 images for test. The models are trained with 68 landmarks annotated for each image of the datasets. However, COFW only contains images with 29 annotated landmarks. Thus, we use the method proposed in [6] to perform a linear mapping between the predictions made on the 68 landmarks to the 29 landmarks, as it is a common practice on this dataset.

For 2D face alignment, the evaluation metric used is the normalized mean error (NME), corresponding to the average point-to-point distance between the ground truth and the predicted shape, normalized by the inter-pupil distance, as it is classically done in the literature:

$$NME = \frac{1}{N} \sum_{i=1}^N \frac{\|\hat{\mathbf{s}}_i - \mathbf{s}_i\|_2}{\|\mathbf{g}_{i,l} - \mathbf{g}_{i,r}\|_2} \quad (12)$$

where  $\hat{\mathbf{s}}_i$  the prediction,  $\mathbf{s}_i$  the ground truth, and  $\mathbf{g}_{i,l}$ ,  $\mathbf{g}_{i,r}$  the left and right pupil centers respectively.

#### 4.1.2 3D Face alignment

The 300W-LP database is a large-pose dataset synthesized from 300W, and contains face images with large variations in pose on the yaw axis, ranging from  $-90^\circ$  to  $+90^\circ$ . The database contains a total of 61225 images obtained by generating additional views of the images from AFW, LFPW, HELEN and I-BUG, using the algorithm from [30]. As it is done in state-of-the-art approaches [31], we train on the augmented images corresponding to 300W trainset as well as their flipped counterparts, making a total of 101144 images for train.

The AFLW2000-3D dataset consists of fitted 3D faces and large-pose images for the first 2000 images of the AFLW database [12]. As it was done in [31], we evaluate the capacities of our method to deal with non-frontal poses by training on 300W-LP and testing on AFLW2000-3D. This database consists of 1306 examples in the  $[0, 30]$  absolute degree yaw range, 462 examples in the  $[30, 60]$  range and 232 examples in the  $[60, 90]$  range. As in [31], we report accuracy for each pose range separately, as well as the mean across those three pose ranges.

3D face alignment consists in localizing the  $(x, y)$  coordinates of the "true" landmarks, as opposed to 2D alignment in which the landmarks are projected on the visible part of the face (e.g. cheeks in case of rotations around the yaw axis). In such a case, the evaluation metric used is also the

normalized mean error, but the normalization is the size of ground truth bounding box, as it is introduced in [31]:

$$NME = \frac{1}{N} \sum_{i=1}^N \frac{\|\hat{\mathbf{s}}_i - \mathbf{s}_i\|_2}{\sqrt{h_i \times w_i}} \quad (13)$$

where  $h_i$ ,  $w_i$  the height and width of the face bounding box, respectively.

## 4.2 Architectures comparison

First, we compare the different architectures detailed in Section 3.5. For a fair comparison, all models have roughly the same number of parameters. Pose-Tree-MoE and Pose-Softmax-MoE use a pretrained head pose estimation model. Once head pose is predicted, each of these models uses it to gate the representations extracted by the committee. The other models proposed in comparison (Baseline, MoE, Softmax-MoE, Tree-MoE) do not use head pose estimation. The results are shown on Table 1.

Table 1: Comparison of different architectures in term of NME (%).

Method	LFPW	HELEN	I-BUG	COFW
Baseline	3.92	4.42	8.89	5.9
MoE	3.79	4.25	8.95	5.87
Softmax-MoE	3.74	4.2	8.8	5.84
Tree-MoE	<b>3.74</b>	4.2	8.38	5.76
Pose-Softmax-MoE	3.84	4.27	7.99	6.12
Pose-Tree-MoE	3.82	<b>4.17</b>	<b>7.5</b>	<b>5.58</b>

In particular, the non gated regressor ensemble is more robust than a single regressor: performance is improved by 3.8% on 300W-Common (LFPW + HELEN). Moreover, adding gates improves performance, especially with tree-gates. Robustness to pose is improved thanks to softmax-gates (8.95  $\rightarrow$  8.8 on I-BUG) and significantly improved thanks to tree-gates (8.95  $\rightarrow$  8.38 on I-BUG). Robustness to occlusions is slightly improved thanks to softmax-gates (5.87  $\rightarrow$  5.84 on COFW) and significantly improved thanks to tree-gates (5.87  $\rightarrow$  5.76 on COFW). This shows that using tree-gated ensembles of regressors allows to substantially increase the overall robustness of the model, particularly in the case of partial occlusions.

Furthermore, using head pose to gate MoE CNN models (Pose-Softmax-MoE and tree-gated-MoE) allows to significantly increase the alignment accuracy on I-BUG, which contains several examples of non-frontal head poses. This is however not the case for Pose-Softmax-MoE model on COFW database, which contains occluded examples. By contrast, Pose-Tree-MoE model generalizes better on COFW (5.76  $\rightarrow$  5.58), and I-BUG (8.38  $\rightarrow$  7.5), without signifi-

cantly degrading performance on frontal faces (4.01  $\rightarrow$  4.03 on average on LFPW and HELEN testsets).

These results show that using ensemble of experts allows for greater robustness, for modelling both the regression and representation layers. The use of gates also allows each expert to be more specialized for a given representation, leading to greater robustness. Last but not least, the hierarchical aspect of tree-gates further improves the use of expert regressors. Conditioning the learned representation to head pose estimation and taking advantage of using ensemble methods all the while learning the gates and expert layers jointly allows these experts to better co-adapt, leading to maximum robustness and accuracy.

### 4.3 Comparison with state-of-the-art approaches

#### 4.3.1 2D Face alignment

Table 2: Comparison with state-of-the-art approaches in term of NME (%).

Method	Common	Challenging	COFW
RCPR [2]	6.18	17.3	8.50
SDM [24]	5.57	15.4	7.70
PIFA [10]	5.43	9.98	-
LBF [15]	4.87	11.98	13.7
TCDCN [27]	4.80	8.60	-
CSP-dGNF [4]	4.76	12.00	-
RAR [23]	4.12	8.35	6.03
RCN <sup>+</sup> [7]	4.20	7.78	-
DRDA [26]	-	10.79	6.46
SFLD [22]	-	-	6.40
PCD-CNN [13]	<b>3.67</b>	7.62	5.77
Pose-Tree-MoE	4.03	<b>7.5</b>	<b>5.58</b>

Table 2 shows a comparison between our approach and other recent state-of-the-art methods on both 300W (common and challenging subsets) and COFW databases. Our model outperforms these approaches on both 300W and COFW databases. The results on COFW show the robustness of our model to occlusions. The alignment error is similar to the human performance on this dataset (5.60 [2]). PCD-CNN [13] essentially uses head pose estimation as a multiplicative variable in a post-hoc processing fashion. By contrast, in Pose-Tree-MoE, head pose is used to select more relevant specialist CNNs to extract adequate features for each head pose range. As one can see, while PCD-CNN is better on 300W-Common, Pose-Tree-MoE significantly outperforms it on both 300W-Challenging and COFW. Therefore, using ensemble of tree-gated experts appears as a more robust way to adapt a face alignment network using head

pose information, that leads to an overall better robustness to large variations in the data.

#### 4.3.2 3D Face alignment

Table 3: Comparison with state-of-the-art approaches on AFLW2000-3D in term of NME (%) for several yaw ranges.

Method	[0, 30]	[30, 60]	[60, 90]	Mean
LBF [15]	8.15	9.49	12.91	10.19
ESR [3]	4.60	6.70	12.67	7.99
CFSS [29]	4.77	6.71	11.79	7.76
RCPR [2]	4.26	5.96	13.18	7.80
MDM [21]	4.85	5.92	8.47	6.41
SDM [24]	3.67	4.94	9.76	6.12
3DDFA [31]	2.84	<b>3.57</b>	4.96	<b>3.79</b>
Tree-MoE	2.84	4.01	4.93	3.92
Pose-Tree-MoE	<b>2.78</b>	3.97	<b>4.76</b>	3.84

Table 3 shows a comparison between our approach and other recent state-of-the-art methods for 3D face alignment on AFLW2000-3D. The state-of-the-art is achieved by the extended version of 3DDFA [31], which fits a 3D dense face model before estimating a sparse set of 68 landmarks. Our Pose-Tree-MoE achieves similar performance as compared to 3DDFA [31], all the while substantially outperforming it on large head poses in the [60, 90] range. Our model also significantly outperforms all other state-of-the-art approaches on this dataset [15, 3, 29, 2, 21, 24]. Furthermore, contrary to [31], our approach only aligns a sparse set of landmarks, thus only requires ground truth landmarks for training, as opposed to the parameters of a morphable model. This shows that using a tree-gated committee of expert CNNs allows to learn relevant experts for each pose range, that produce suitable representations upon which the tree-gated MoE layer can adaptively align the facial landmarks. Conditioning the representation using the head pose estimate significantly improves the results on large poses. This is confirmed by the comparison between Pose-Tree-MoE and Tree-MoE.

In what follows, we propose a number of qualitative experiments to assess that the head pose clustering of the expert CNNs behaves as expected.

### 4.4 Qualitative evaluation

In this section, we conduct some experiment to provide insight on how the gated models behave, by visualizing the contributions of tree-gates:

- Interpretability through hierarchical clustering visualization in representation layer. This allows to study how the model splits the poses space in order to extract the representation. In addition, this ensures consistency



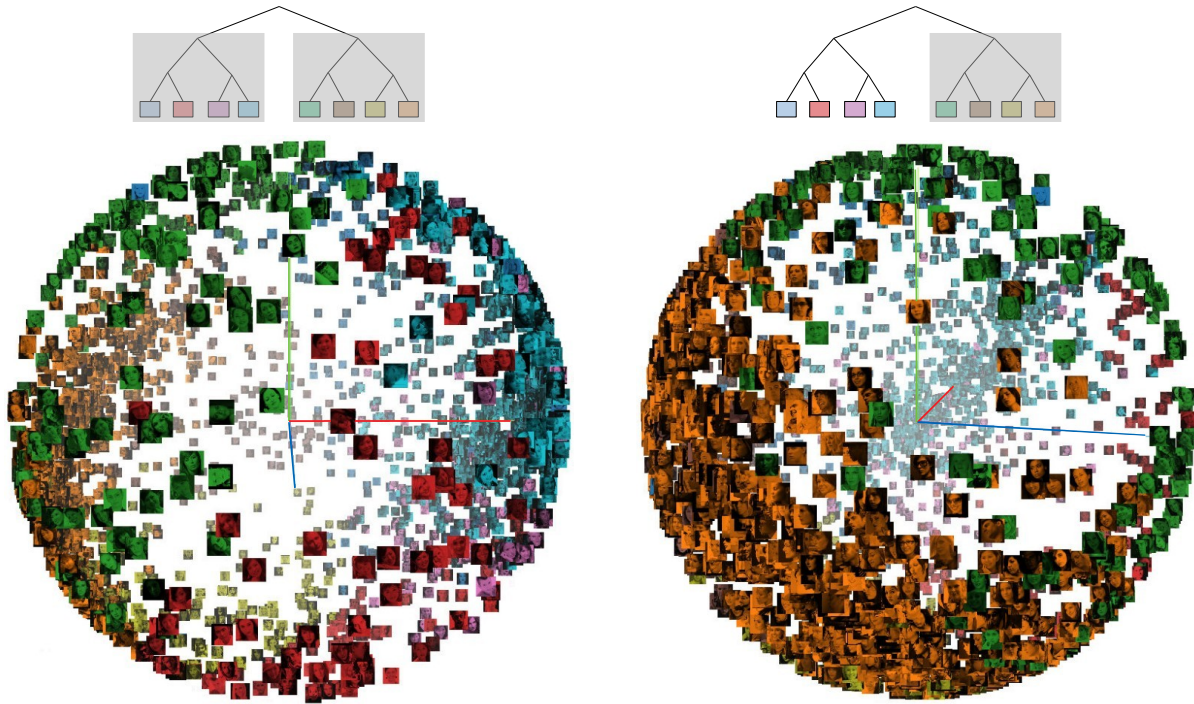


Figure 5: Visualization of representation clusters (Pose-Tree-MoE on AFLW2000-3D). **Top:** Associated color for each cluster (leaf of the 3-depth tree). Each image is then colored according to the cluster with the highest weight in the final output. **Bottom:** Dispersion in head pose space (red axis: yaw, green axis: pitch, blue axis: roll). For more visibility, the data is normalized by shifting each point by the centroid and making it unit norm. At the first level of the tree, the split is mainly due to the yaw orientation, as illustrated by the comparison of the purple/blue images in the left graph with the orange/brown images in the right graph. For the right sub-tree, the split mainly uses the yaw and roll intensities, as shown on the right graph.

between the spatial distribution of poses and the use of expert CNNs.

- Efficiency through the distribution and use of expert FCs in regression layer. This allows for fewer regressors to be used, whose predictions are accurate.

#### 4.4.1 Representation layer

As seen in the previous section, integrating head pose information to extract representations significantly improves the robustness to strong variations in pose, and results in a model that exceeds the state-of-the-art. It might be interesting to introspect the model in order to study how it behaves. To do this, we propose to visualize hierarchical clustering performed by our model on a dataset with maximum pose variability, such as AFLW2000-3D. Figure 5 represents the faces of AFLW2000-3D in the pose space, where each face is colored according to the expert CNN with the most weight in the committee. Since a unique color is given for each expert CNN, we can observe the splitting performed by the gates on the dataset. Figure 5 illustrates the repartition in

head pose space from a Pose-Tree-MoE model trained on 300W-LP. We can then observe that on the first level of the tree, the red axis representing the yaw allows to separate the data associated with the two subtrees respectively. The same can be said for the second tree level, therefore the model learned to split the head pose space according to the yaw orientation primarily. This is consistent with the fact that the model was trained with 300W-LP, whose images essentially augmented with yaw.

#### 4.4.2 Regression layer

Figure 6 shows the average of the cumulative sum of the gates probabilities of expert CNNs, sorted in descending order on I-BUG. Notice on the right part that the tree-gated model allows 20% of the regressors to explain more than 90% of the final prediction, while the softmax-gated model needs about 40% of the regressors to explain 90% of the prediction. Thus, tree-gates allows to output a correct alignment using less expert regressors: the better repartition of expert regressors towards the specific alignment cases makes it pos-

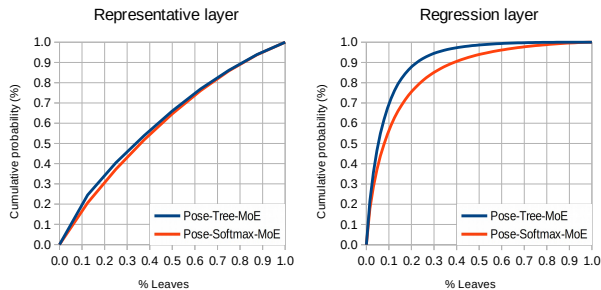


Figure 6: Cumulative top-scoring expert networks distribution for representation and regression layer at the first stage in the cascade for both Pose-Softmax-MoE and Pose-Tree-MoE.

sible to better specialize each expert regressor, and to use fewer of these to obtain a better representation. The results for the representation layer are less clear-cut, the Pose-Tree-MoE model lying marginally above the Pose-Softmax-MoE model. This is likely due to the lower number of experts (8 vs 64 for representation vs prediction), indicating that the difference between tree and softmax-gated MoE models becomes more conspicuous as the number of experts increases. All in all, tree gates promotes a more efficient repartition of the experts, and a better specialization thereof, which, in turns, leads to an overall higher accuracy and robustness.

#### 4.4.3 Visualizations

Figure 7 shows the predictions of our model (5<sup>th</sup> and 11<sup>th</sup> columns) for large pose examples on AFLW2000-3D database, as compared with the ground truth markup (6<sup>th</sup> and 12<sup>th</sup> columns). The prediction and ground truth landmark localizations are also plotted along with the estimated and ground truth head pose. For most examples, the head pose estimate is very close to its ground truth counterpart: this allows to select relevant expert CNNs in the representation layers, which give rise to high-quality landmark alignment even on examples exhibiting large pose variations.

Furthermore, rows 1 to 4 and 7 to 10 shows the displacements outputted by only one top-scoring regressor (as indicated by the associated tree-gate value), from the current shape at the corresponding cascade iteration. It should be noted that using a single regressor, our Pose-Tree-MoE model can achieve reasonable alignment accuracy, which will justify investigating the use of a restricted (top-k) numbers of experts in future work, e.g. using greedy evaluation as in [4].

#### 4.4.4 Runtime evaluation

Last but not least, our method is very fast as it operates at 17.54 ms per image on a NVIDIA GTX 1080 GPU, and thus can run at 57 fps. Furthermore, In [18], MoE are used to reduce the computational load by keeping only a small number (top-k) of experts. With hierarchical gates, an interesting direction would be to evaluate the tree-gate in a greedy layerwise fashion as in [4], and keep only the regressor corresponding to the maximum probability leaf to further reduce the computational cost.

## 5 Conclusion

In this paper, we have proposed to integrate ensemble methods within a deep architecture in order to increase the overall robustness of the model to large variations in the data. The use of a committee of experts neural networks instead of a single one allows an overall greater robustness. Furthermore, we showed that using a gate function to weight the responses of each expert network allows each of these networks to be more expert for a given context. In particular, the use of tree-gates makes it possible to jointly learn a committee of expert networks and a hierarchical clustering of the use of these experts. Additionally using neural trees to model the tree-gates allows to learn both the ensemble and associated gating network in an end-to-end manner.

As such, we showed that tree-gated MoE models can be used for modelling the regressors as well as the feature representation layers, by using high-level semantic information such as head pose as a proxy variable. These tree-gates allows a more efficient clustering and specialization of the experts, leading to a higher performance. Furthermore, thorough experimental validation, we demonstrated that, when applied for face alignment in the frame of cascaded regression, the proposed approach yields high accuracies, most notably on challenging data in term of head pose and occlusion, while keeping a reasonable computational cost.

As a future work, we will investigate the use of a limited number (top-k) of experts for both the representation and regression layers, e.g. using greedy evaluation [4], in order to further decrease the runtime. Furthermore, the tree-MoE architecture introduced in this paper is very generic and could be applied to a wide range of other computer vision problem, such as image classification, semantic segmentation, or object detection.

## Acknowledgment

This work has been supported by the French National Agency (ANR) in the frame of its Technological Research JCJC program (FacIL, project ANR-17-CE33-0002).

## References

- [1] E. Arnaud, A. Dapogny, and K. Bailly. Tree-gated Deep Regressor Ensemble For Face Alignment In The Wild. In *FG*, 2019. 4
- [2] X. P. Burgos-Artizzu, P. Perona, and P. Dollár. Robust face landmark estimation under occlusion. In *ICCV*, pages 1513–1520, 2013. 3, 7, 8
- [3] X. Cao, Y. Wei, F. Wen, and J. Sun. Face alignment by explicit shape regression. In *CVPR*, pages 2887–2894, 2012. 8
- [4] A. Dapogny, K. Bailly, and S. Dubuisson. Face alignment with cascaded semi-parametric deep greedy neural forests. *Pattern Recognition Letters*, 102:75–81, 2018. 3, 8, 10
- [5] D. Eigen, M. Ranzato, and I. Sutskever. Learning factored representations in a deep mixture of experts. *arXiv preprint arXiv:1312.4314*, 2013. 3
- [6] G. Ghiasi and C. C. Fowlkes. Occlusion coherence: Localizing occluded faces with a hierarchical deformable part model. In *CVPR*, pages 1899–1906, 2014. 3, 7
- [7] S. Honari, P. Molchanov, S. Tyree, P. Vincent, C. Pal, and J. Kautz. Improving landmark localization with semi-supervised learning. In *CVPR*, pages 1546–1555, 2018. 2, 8
- [8] Y. Ioannou, D. P. Robertson, D. Zikic, P. Kotschieder, J. Shotton, M. R. Brown, and A. Criminisi. Decision forests, convolutional networks and the models in-between. *arXiv preprint arXiv:1603.01250*, 2016. 3
- [9] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, G. E. Hinton, et al. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991. 3
- [10] A. Jourabloo, M. Ye, X. Liu, and L. Ren. Pose-invariant face alignment with a single cnn. In *ICCV*, pages 3200–3209, 2017. 8
- [11] P. Kotschieder, M. Fiterau, A. Criminisi, and S. R. Bulò. Deep neural decision forests. In *ICCV*, pages 1467–1475, 2015. 2, 3, 5
- [12] M. Köstinger, P. Wohlhart, P. M. Roth, and H. Bischof. Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization. In *ICCV Workshops*, pages 2144–2151, 2011. 7
- [13] A. Kumar and R. Chellappa. Disentangling 3d pose in a dendritic cnn for unconstrained 2d face alignment. In *CVPR*, pages 430–439, 2018. 3, 8
- [14] I. Misra, A. Shrivastava, A. Gupta, and M. Hebert. Cross-stitch networks for multi-task learning. In *CVPR*, 2016. 3
- [15] 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. 1, 8
- [16] N. Ruiz, E. Chong, and J. M. Rehg. Fine-grained head pose estimation without keypoints. In *CVPR Workshops*, pages 2155–215509, 2018. 3
- [17] C. Sagonas, E. Antonakos, G. Tzimiropoulos, S. P. Zafeiriou, and M. Pantic. 300 faces in-the-wild challenge: database and results. *Image and Vision Computing*, 47:3–18, 2016. 6
- [18] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. V. Le, G. E. Hinton, and J. Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017. 3, 10
- [19] Y. Sun, X. Wang, and X. Tang. Deep convolutional network cascade for facial point detection. In *CVPR*, pages 3476–3483, 2013. 1
- [20] R. Tanno, K. Arulkumaran, D. C. Alexander, A. Criminisi, and A. V. Nori. Adaptive neural trees. *arXiv preprint arXiv:1807.06699*, 2018. 3
- [21] G. Trigeorgis, P. Snape, M. A. Nicolaou, E. Antonakos, and S. P. Zafeiriou. Mnemonic descent method: A recurrent process applied for end-to-end face alignment. In *CVPR*, pages 4177–4187, 2016. 1, 4, 8
- [22] Y. Wu, C. Gou, and Q. Ji. Simultaneous facial landmark detection, pose and deformation estimation under facial occlusion. In *CVPR*, pages 5719–5728, 2017. 3, 8
- [23] S. Xiao, J. Feng, J. Xing, and H. Lai. Robust Facial Landmark Detection via Recurrent Attentive-Refinement Networks. In *ECCV*, volume 1, pages 57–72, 2016. 8
- [24] X. Xiong and F. De la Torre. Supervised descent method and its applications to face alignment. In *CVPR*, pages 532–539, 2013. 1, 8
- [25] X. Yu, Z. L. Lin, J. Brandt, and D. N. Metaxas. Consensus of regression for occlusion-robust facial feature localization. In *ECCV*, 2014. 3
- [26] J. J. Zhang, M. Kan, S. Shan, and X. Chen. Occlusion-free face alignment: Deep regression networks coupled with de-corrupt autoencoders. In *CVPR*, pages 3428–3437, 2016. 3, 8
- [27] Z. Zhang, P. Luo, C. C. Loy, and X. Tang. Learning deep representation for face alignment with auxiliary attributes. *PAMI*, 38(5):918–930, 2015. 2, 8
- [28] Z.-H. Zhou and J. Feng. Deep forest: Towards an alternative to deep neural networks. In *IJCAI*, 2017. 6
- [29] S. Zhu, C. C. Li, C. C. Loy, and X. Tang. Face alignment by coarse-to-fine shape searching. In *CVPR*, pages 4998–5006, 2015. 8
- [30] X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li. Face alignment across large poses: A 3d solution. In *CVPR*, pages 146–155, 2016. 7

- [31] X. Zhu, X. Liu, Z. Lei, and S. Z. Li. Face alignment in full pose range: A 3d total solution. *PAMI*, 41:78–92, 2019. [3](#), [7](#), [8](#)



Figure 7: Visualisations of the predictions outputted for each cascade stage with only the top (maximum value of tree-gate) regressor. Head pose estimation is also displayed, as well as the ground truth. Images from AFLW2000-3D.