

# KEPLER: Keypoint and Pose Estimation of Unconstrained Faces by Learning Efficient H-CNN Regressors

Amit Kumar, Azadeh Alavi and Rama Chellappa  
Department of Electrical and Computer Engineering, CFAR and UMIACS  
University of Maryland-College Park, USA  
{akumar14,azadeh,rama}@umiacs.umd.edu

**Abstract**—Keypoint detection is one of the most important pre-processing steps in tasks such as face modeling, recognition and verification. In this paper, we present an iterative method for Keypoint Estimation and Pose prediction of unconstrained faces by Learning Efficient H-CNN Regressors (KEPLER) for addressing the face alignment problem. Recent state of the art methods have shown improvements in face keypoint detection by employing Convolution Neural Networks (CNNs). Although a simple feed forward neural network can learn the mapping between input and output spaces, it cannot learn the inherent structural dependencies. We present a novel architecture called H-CNN (Heatmap-CNN) which captures structured global and local features and thus favors accurate keypoint detection. H-CNN is jointly trained on the visibility, fiducials and 3D-pose of the face. As the iterations proceed, the error decreases making the gradients small and thus requiring efficient training of DCNNs to mitigate this. KEPLER performs global corrections in pose and fiducials for the first four iterations followed by local corrections in a subsequent stage. As a by-product, KEPLER also provides 3D pose (pitch, yaw and roll) of the face accurately. In this paper, we show that without using any 3D information, KEPLER outperforms state of the art methods for alignment on challenging datasets such as AFW [40] and AFLW [18].

## I. INTRODUCTION

Keypoint detection on unconstrained faces is one of the most studied topics in the past decade, as accurate localization of fiducials is a vital pre-processing task for variety of applications. In the last five years, keypoint localization using Convolution Neural Networks (CNN) has received great attention from computer vision researchers. This is mainly due to the availability of large scale annotated unconstrained face datasets such as AFLW [18]. Works such as [33] have hypothesized that as the network gets deeper more abstract information such as identity, pose, attributes are retained while immediate local features are lost. However, various methods [24], [34], and [36] directly use CNNs as regressors or use deep features from CNNs to design regressors for predicting keypoints.

On the other hand, an earlier method of Explicit Shape Regression (ESR) proposed by Cao et al. [8] achieved superior results by introducing the important concept of non-parametric shape regression for facial keypoint localization. Following [8], unconstrained face alignment received great deal of attention and many of its variants [20], [23], [16], [24], [19] were published later, using a variety of features producing incremental improvements over [8]. However, they are all limited by the fixed number of points on the face.

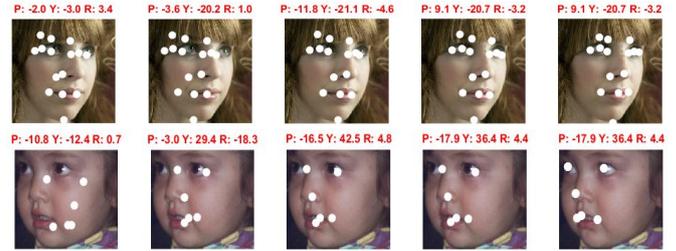


Fig. 1: Sample results generated by the proposed method. White dots represent the location of keypoints after each iteration. The first row shows an image from the AFLW dataset. The points move at subpixel level after fourth iteration. The second row is a sample image from the AFW dataset, which shows how the last stage of error correction can effectively mitigate the inconsistency of the bounding box across datasets. The numbers in red are the predicted 3D pose P:Pitch Y:Yaw R:Roll

In real life applications, there are more challenging datasets such as IJBA [17] and AFW [40], which do not always have 68 or 49 fixed points mainly due to occlusion or pose variations. As alternatives, researchers moved towards more sophisticated techniques, incorporating 3D shape models [39], [14], [13], domain learning [38], recurrent autoencoder-decoder [1] and many others. However, one question still remains unanswered: Can cascaded shape regression be applied for an arbitrary face with no prior knowledge ?

The motivation for this work is to adapt cascaded regression for predicting landmarks of arbitrary faces, while taking advantage of CNNs. We transform the cascaded regression formulation into an iterative scheme for arbitrary faces. In each iteration the regressor predicts the increment for the next stage jointly for all the points while maintaining the shape constraint. As by-products of KEPLER, we get the visibility confidence of each keypoint and 3D pose (pitch, yaw and roll) for the face image. The main contributions of this paper are:

- We design a novel GoogLeNet-based [26] architecture with a channel inception module which pools features from intermediate layers and concatenates them similar to inception module. We call the proposed architecture *Channeled Inception* in the rest of the paper. This network is used in all the stages of KEPLER.
- Inspired by [9], we present an iterative method for estimating the face landmarks using the fixed point

consolidation scheme inspired by [9]. We observe that estimating landmarks on a face is more challenging than estimating keypoints on a human body. The overview of the pipeline is shown in Figure 2.

- After each iteration, the error from ground-truth decreases, making the gradient smaller and hence different training policies are employed in every stage for the efficient training of H-CNN.
- We evaluate the performance of our keypoint estimation method on challenging datasets such as AFLW and AFW, which include faces in diverse poses and expressions. We also introduce a new protocol for evaluating the facial keypoint localization scheme on the AFLW dataset which is more challenging and usually left out while evaluating unconstrained face alignment methods.

The rest of the paper is organized as follows. Section II reviews closely related works. Section III presents the proposed method in detail. Section IV describes the experiments and comparisons, which are then followed by conclusions and suggestions for future works in section V.

## II. RELATED WORK

Following [8], we classify previous works on face alignment into two basic categories.

**Part-Based Deformable models:** These methods perform alignment by maximizing the confidence of part locations in a given input image. One of the major works in this category was done by Zhu and Ramanan [40], where they used a part-based model for face detection, pose estimation and landmark localization assuming the face shape to be a tree structure. [5] by Asthana et al., learned a dictionary of probability response maps followed by linear regression in a Constrained Local Model (CLM) framework. Hsu et al. [11] extended the mixture of tree model [40] to achieve better accuracy and efficiency. However, their method again assumes face shape to be a tree structure, enforcing strong constraints specific to shape variations.

**Regression-based approaches:** Since face alignment is naturally a regression problem, a multitude of regression-based approaches has been proposed in recent years. Methods reported in [21], [8], [36] are based on learning a regression model that directly maps image appearances to target outputs. However, these methods along with methods from [4], [27], [3], [6], [28] and [25] were mostly evaluated either in a lab setting or on face images where all the facial keypoints are visible. Wu et al. [30] proposed an occlusion-robust cascaded regressor to handle occlusion. Xiong et al. [31] pointed out that standard cascaded regression approaches such as Supervised Descent Method (SDM) [32] tend to average conflicting gradient directions resulting in reduced performance. Hence, [31] suggested domain dependent descent maps. Inspired by this, Cascade Compositional Learning (CCL) [38] and Ensemble of Model Regression Trees (EMRT) [37] developed head pose based and domain selective regressors respectively. [38] partitioned the optimization domain into multiple directions based on

head pose and learned to combine the results of multiple domain regressors through composition estimator function. Similarly [37] trained an ensemble of random forests to directly predict the locations of keypoints whereafter face alignment is achieved by aggregating the consensus of different models.

Recently, methods using 3D models for face alignment have been proposed. PIFA [13] by Jourabloo et al. suggested a 3D approach that employed cascaded regression to predict the coefficients of 3D to 2D projection matrix and the base shape coefficients. Another recent work from Jourabloo et al. [14] formulated the face alignment problem as a dense 3D model fitting problem, where the camera projection matrix and 3D shape parameters were estimated by a cascade of CNN-based regressors. However, [38] suggests that optimizing the base shape coefficients and projection is indirect and sub-optimal since smaller parameter errors are not necessarily equivalent to smaller alignment errors. 3DDFA [39] by Zhu et al. fitted a dense 3D face model to the image via CNN, where the depth data is modeled in a Z-Buffer.

Our work principally falls in the category of regression-based approaches and addresses the issue of adapting the cascade shape regression to unconstrained settings. KEPLER performs joint training on three fundamental tasks, namely, 3D pose, visibility of each keypoint and the location of keypoints, using only 2D color image. It also demonstrates that efficient joint training on the three tasks achieves superior performance. One of the closely related work is [35] where the authors used multi-tasking for many attributes, but did not leverage the intermediate features.

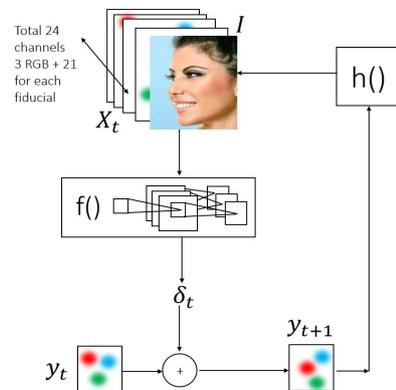


Fig. 2: Overview of the architecture of KEPLER. The function  $f()$  predicts visibility, pose and the corrections for the next stage. The representation function  $h()$  forms the input representation for the next iteration.

## III. KEPLER

KEPLER is an iterative method which at its core consists of three modules. Figure 2 illustrates the basic building blocks of KEPLER. The first module is a rendering module  $h$  which models the structure in an N-dimensional input space, with N being the maximum number of keypoints on a face. The current location of the keypoints are represented by the

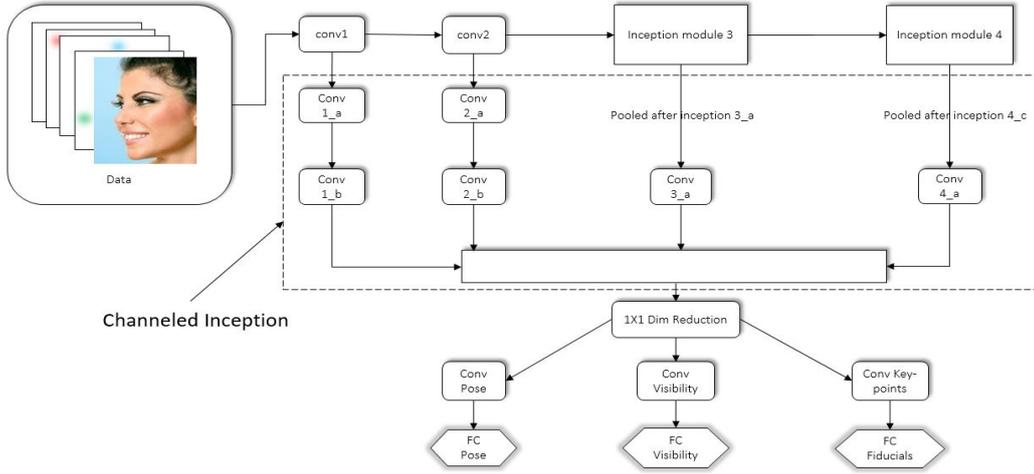


Fig. 3: The KEPLER network architecture. The dotted line shows the channeled inception network. The intermediate features are convolved and the responses are concatenated in a similar fashion as the inception module. Tasks such as pose are abstract and contained in deeper layers, however, the localization property is in the shallower layers.

vector  $\mathbf{y}_t = \{y_t^1 \dots y_t^N\}$ . The output of the rendering module is concatenated to the raw RGB input image  $\mathbf{I}$ , along the third dimension which is then fed to the function  $f$ .

The second module is the function  $f$  which calculates the correction to be made at the next stage. The function  $f$  is modeled by a convolution neural network whose architecture is described in section III-A.

The third module is the correction stage which adds the increments, predicted by  $f$ , to the current locations. The output goes again into the rendering module  $h$  which prepares the rendered data for the next iteration. The rendering function is not learned in this work, but represented by a 2D Gaussian with fixed variance and centered at current keypoint locations in each of the  $N$  channels. Finally, the Gaussian rendered images are stacked together with image  $\mathbf{I}$ . Therefore the overall method can be summarized by the following set of equations.

$$\delta_t = f_t(\mathbf{X}_t, \Theta_t) \quad (1)$$

$$\mathbf{y}_{t+1} = \mathbf{y}_t + \delta_t \quad (2)$$

$$\mathbf{X}_{t+1} = h(\mathbf{y}_{t+1}) \quad (3)$$

where  $f$  is a function with learned parameters  $\Theta_t$ , predicting the increments  $\delta_t$ . The prediction function  $f$  is indexed by  $t$  as it is trained separately for every iteration. In the first iteration, the function  $h$  renders Gaussians at  $y_0$ , which is the mean shape. In this work we set  $t = 5$  iterations. We perform the last iteration only to take into effect the improper bounding box across different datasets (see Figure 1). The loss functions for each task is mentioned below.

#### Keypoint localization

Keypoint localization is the task of predicting the keypoints in a face. In this paper, we consider predicting the locations of  $N = 21$  keypoints on the face. With each point is associated the visibility of that point. The loss function for this task is

given by

$$L_1(\mathbf{y}, \mathbf{g}) = \sum_{i=1}^N v^i (y_t^i - g^i)^2, \quad (4)$$

where  $y_t^i$  and  $g^i$  are the predicted and the ground truth locations of the  $i^{th}$  keypoint respectively at time  $t$ .  $v^i$  is the ground truth visibility associated with each keypoint. We discuss this loss function and its variant in section III-C.

#### Pose Prediction

Pose prediction refers to the task of estimating the 3D pose of the face. We use the Euclidean loss function for pose prediction.

$$L_2(\mathbf{p}_p, \mathbf{g}_p) = (p_{yaw} - g_{yaw})^2 + (p_{pitch} - g_{pitch})^2 + (p_{roll} - g_{roll})^2 \quad (5)$$

where  $p$  stands for predicted and  $g$  for the ground-truth.

#### Visibility

This task is associated with estimating the visibility of each keypoint. The number of keypoints visible on the face varies with pose. Hence, we use the Euclidean loss to estimate the visibility confidence of each point.

$$L_3(\mathbf{v}_p, \mathbf{v}_g) = \sum_{i=1}^N (v_{p,i} - v_{g,i})^2, \quad (6)$$

Therefore the net loss in the network is the weighted linear combination of the above loss functions.

$$L(p, g) = \lambda L_1(\mathbf{y}, \mathbf{g}) + \mu L_2(\mathbf{p}_p, \mathbf{g}_p) + \nu L_3(\mathbf{v}_p, \mathbf{v}_g) \quad (7)$$

where  $\lambda$ ,  $\mu$  and  $\nu$  are the weight parameters suitably chosen depending on the iteration.

#### A. Network Architecture

For the modeling function  $f$  we design a unique ConvNet architecture based on GoogLeNet [26] by pruning the inception network after inception\_4c. As PReLU has shown better performance in many vision tasks such as object

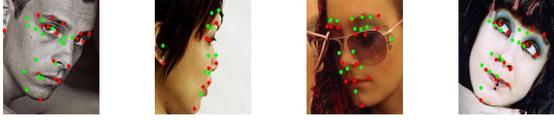


Fig. 4: Qualitative results of KEPLER after second stage. The green dots represent the predicted points after second stage. Red dots represent the ground truth. It can be seen that the visible points have taken the shape of input face image.

recognition [10], in this pruned network we first replace the ReLU non-linearity with PReLU. We pool the intermediate features from the pruned GoogLeNet. Then convolutions are performed from the output of each branch, and the output maps are concatenated similar to the inception module. We call this module the *Channeled Inception* module. Since the output maps after  $conv_1$  are larger in size, we first perform  $4X4$  convolution and then again a  $4X4$  convolution, both with the stride of 3 to finally match the dimension of the output to  $7X7$ . Similarly after  $conv_2$  we first perform  $4X4$  convolution and then  $3X3$  convolution to match the output to  $7X7$ . The former uses a stride of 4 and the latter uses 2. The most naïve way of combining features is by concatenation. However, the concatenated output blob can be very high dimensional and hence we perform  $1X1$  convolution for dimensionality reduction. This lets the network decide the weights to effectively combine the pooled features into lower dimension. It has been shown in [33] that adjacent layers are correlated and hence, we only pool features from alternate layers.

Next the network is trained on three tasks namely, pose, visibilities and the bounded error using ground truth. The joint training is helpful since it models the inherent relationship between visible number of points, pose and the amount of correction needed for a keypoint in particular pose. Choosing an architecture like GoogLeNet is based on the fact that due to fewer number of parameters the training of GoogLeNet is faster and adding to it batch normalization, even speeds up the training process. In order to further speed up the process we only use convolution layers till the last layer where we use a fully connected layer to get the final output. The architecture of the whole network is shown in Figure 3.

### B. Iteration 1 and 2: Constrained Training

In this section, we explain the first stage training for keypoint estimation. The first stage is the most crucial one for face alignment. Since the network is trained from scratch, precautions have to be taken on what the network should learn. Directly learning the locations of keypoints from a network is difficult because when the network gets deeper it loses the localization capability. This is due to the fact that the outputs of the final convolution layers have a larger receptive field on the input image. We devise a strategy in which the corrections for the first two stages are bounded. Let us suppose the key-points are represented by their 2D coordinates  $\mathbf{y} : \{y^i \in \mathbb{R}^2, i \in [1, \dots, N]\}$  where  $N$  is the

number of keypoints and  $y^i$  denotes the  $i^{th}$  keypoint. The bounded corrections were calculated using (8) given below.

$$\delta_t^i(g^i, y_t^i) = \min(L, \|\mathbf{u}\|) \cdot \hat{\mathbf{u}} \quad (8)$$

where  $L$  denotes the bound of correction.  $\mathbf{u} = \mathbf{g} - \mathbf{y}_t$  and  $\hat{\mathbf{u}} = \frac{\mathbf{u}}{\|\mathbf{u}\|}$  represent the error vector and error unit vector respectively. In our experiments we set the bound  $L$  to a maximum of 20 pixels. This simplifies the learning problem for the network in the first stage. According to this formulation, error correction for points for which the ground truth is far away, gets bounded by  $L$ . The interesting property of this formulation is that in the first and second stage the network only learns the direction in which the points have to shift. This can be thought of as learning the direction of the unit error vector, to which the magnitude will be added later. In addition to just having keypoint location we also have access to facial 3D pose and the visibility of each point. One-shot prediction of the location of keypoints is difficult since the input space of the ConvNet is typically nonlinear. Also, learning small corrections should be easier, when the network is being trained for the first time. Hence, to impart prior knowledge to the network we jointly learn the pose and visibility of each point. The loss functions used for the three tasks are described in the previous section.

The function  $f$  for second iteration is trained in a similar fashion with the weights initialized from the first iteration.

### C. Iteration 3: Variant of Euclidean loss

We show the outputs of the network after the second stage of training in Figure 4. Physical inspection of the outputs shows that for many of the faces, the network has already learned the magnitude and direction of the correction vector. However, there are misalignments in some images or in some keypoints in the images. But repeating the training methodology exactly as second iteration revealed that our architecture suffered from vanishing gradients. While back propagating the gradients, the loss is averaged over a batch and if there are few misalignments in a batch, there is very little gradient to be propagated. To maintain consistency we stick with the same architecture. Even though GoogLeNet [26] claims to not have vanishing gradient problem, KEPLER faced it because of the dataset being small.

This motivated us to design a loss function that satisfies both of these conditions: on the one hand, the loss function should minimize the error between prediction and the ground truth; on the other hand, it should have sufficient gradients to be propagated to make the learning process reach global minima. Towards this end, we use the following loss function.

$$L_1(\mathbf{y}, \mathbf{g}) = \frac{1}{n} \left( \sum_{i=1}^N v_i (y_i - g_i)^2 + \gamma \sum_{i=1}^N v_i |y_i - g_i| \right) \quad (9)$$

$$\frac{\delta L_1(\mathbf{y}, \mathbf{g})}{\delta \mathbf{y}} = \frac{1}{n} \left( 2 \sum_{i=1}^N v_i (y_i - g_i) + \gamma \sum_{i=1}^N v_i \frac{|y_i - g_i|}{y_i - g_i} \right) \quad (10)$$

where  $\gamma$  is a parameter which controls the strength of the gradient and  $n$  is the number of samples in a batch. We would like to emphasize that the additional term is not a

regularizer as it is added to the objective function and does not directly regularize the weights. However, this is able to provide substantial gradients for the training of ConvNet. The representation function  $h$  in this stage does not render any Gaussian in the channel for which the predicted visibility is below the threshold  $\tau$ . In this work we set this threshold  $\tau$  to 0.03 and  $\gamma$  to 0.2 obtained by cross validation. We do not constrain the amount of error corrections for the third stage training.

#### D. Iteration 4: Hard sample mining

Recently, Kabkab et al. [15] suggested that by efficiently sampling the data one can make an optimal use of training set while training ConvNets leading to improved performance. [15] developed an online data sampling method based on a convex optimization formulation and showed how their formulation can make the classifier robust in class imbalanced problem. In our case, although after the third iteration, most of the images are aligned, they lack precision in local alignment. Inspired by [15], we reuse the hard samples of the dataset to build a more robust keypoint localization system.

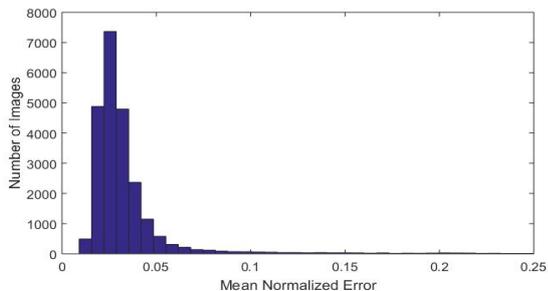


Fig. 5: Error Histogram of training samples after stage 3

Using the keypoints predicted after the third iteration, we plot the histogram (Fig.5) of normalized mean error (NME), after calculating it for all the training samples. We denote the NME on x-axis around which the maximum number of samples are centered, as  $C$ . In an ideal case, the value of  $C$  should be low, implying that the average alignment error is less. Therefore, the objective of this stage is to lower the value of  $C$  by hard sample mining. We select a threshold  $\Delta$  (0.03 in our experiments), towards the right of  $C$ , after which at least 30 – 40% of the samples lie, as the threshold for hard samples. Using  $\Delta$ , we partition the dataset into two groups of hard and easy samples. We first select equal number of samples from both groups to form a batch which is then presented to ConvNet for training. This effectively results in reusing the hard samples. Then, to counter the group imbalance we finetune the network with entire dataset again with a lower learning rate. We use the loss function as in (9) with  $\gamma = 0.1$  for this stage.

#### E. Iteration 5: Local Error Correction

There is a lot of inconsistency among the bounding boxes provided by different datasets. AFLW [18] provides larger bounding box annotations compared to AFW [40].

Regression-based alignment methods are dependent on the mean shape initialization, which is scaled to the bounding box size. Also it is impractical to come up with a heuristic which tries to determine compatible bounding boxes. Almost all the existing methods perform data augmentation by randomly perturbing the bounding boxes by some amount. However, it is not clear by how much the bounding boxes should be perturbed to obtain reasonably good bounding boxes during testing which is consistent with the dataset the network was trained on. We train our networks on a larger bounding box provided by AFLW. AFLW bounding boxes tend to be square and for almost all the images the nose tip appears at the center of the bounding box. This is a big limitation for the deployment of the system in real world scenarios. It is worthy to note that the previous four stages are trained on full images and hence produce global corrections. Our last stage of local correction is optional, which depends

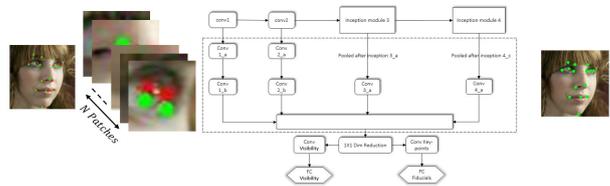


Fig. 6: Red dots in the left image represent the ground truth while green dots represent the predicted points after the fourth iteration. Local patches centered around predicted points are extracted and fed to the network. The network shown in Fig 3 (see section III-E for details) is trained on the task of local fiducial correction and visibility of fiducials inside the patch. The image on the right shows the predictions after local correction.

upon the test set and the bounding box annotations that it comes with. We train an exactly similar network as in Fig 3. but only for the tasks of predicting the visibility and corrections in the local patches. Predicting the pose with a local patch of say  $WXW$  pixels is difficult which can lead the network to learn improper weights. We choose all the  $N$  patches irrespective of the visibility factor. Learning visibility and corrections is important because we do not want the network to propagate any gradient if the point is invisible. We observe during experimentation that training the ConvNet on two tasks together achieves significantly better performance than when the network is trained only for the task of error correction. We again partition the dataset into easy and hard sample groups according to the strategy explained in the previous section. We finally finetune the network with the whole dataset with lower learning rate.

## IV. EXPERIMENTS AND COMPARISON

### A. Datasets

We select two challenging datasets with their most recent benchmarks.

**In-the-wild datasets:** To make the system robust for images in real life scenarios such as challenging shape variations and significant view changes, we select AFLW [18] for training and, AFLW and AFW [40] as the main test sets.

**AFLW** contains 24,386 in-the-wild faces (obtained from *Flickr*) with head pose ranging from  $0^\circ$  to  $120^\circ$  for yaw and upto  $90^\circ$  for pitch and roll with extremely challenging shape variations and deformations. Along with this AFLW also demonstrates external-object occlusion. There are a total of 21% invisible landmarks caused by occlusion, larger than 13% on COFW [7] where only internal object-occlusion is exhibited. In addition, one important point to note is that COFW also provides the annotations for the invisible landmarks while in the case of AFLW the invisible landmarks are absent. **AFW** is a popular benchmark for the evaluation of face alignment algorithms. AFW contains 468 in-the-wild faces (obtained from *Flickr*) with yaw degree up to  $90^\circ$ . The images are well diverse in terms of pose, expression and illumination. The number of visible points also varies depending on the image, but the location of occluded points are to be predicted as well.

AFLW provides at most 21 points for each face. It excludes coordinates for invisible landmarks, which we consider to be the best, because there is no way of correctly knowing the exact location of those points. In many cases such invisible points are mostly hallucinated and annotated thereafter.

### Testing Protocols:

**(I)AFLW-PIFA:** We follow the protocol used in PIFA [13]. We randomly select 23,386 images for training and the remaining 1,000 for testing. We divide the testing images in three groups as done in [13]:  $[0^\circ, 30^\circ]$ ,  $[30^\circ, 60^\circ]$  and  $[60^\circ, 90^\circ]$  where the number of images in each group are taken to be equal.

**(II)AFLW-Full:** We also test on the full test set of AFLW of sample size 1,000.

**(III)AFLW-All variants:** In the next experiment, to have more rigorous analysis, we perform the test on all variants of images from (I) above. To create all variants images, we first rotate the whole images from (I) at angles of  $15^\circ, 30^\circ, 45^\circ$  and  $60^\circ$ . We do the same with the horizontally flipped version of these images. We then rotate the bounding box coordinates and the key-points also at the same angles and crop the faces. This is done for all the images following the AFLW-PIFA protocol. One important effect of this rotation is that some of the images have smaller face compared to others due to rotated bounding box. This experiment tests the robustness of the algorithm on faces of different effective sizes and orientations.

**(IV)AFW:** We only use AFW for testing purposes. We follow the protocol as stated in [40]. AFW provides 468 images in total, out of which 341 faces have height greater than 150 pixels. We only evaluate on those 341 images following the protocol of [40].

**Evaluation metric:** Following most previous works, we obtain the error for each test sample via averaging normalized errors for all annotated landmarks. We demonstrate our results with mean error over all samples, or via Cumulative Error Distribution (CED) curve. For

pose, we evaluate on continuous pose predictions as well as their discretized versions rounded to nearest  $15^\circ$ . We report the continuous mean absolute error for the AFLW testset and plot the Cumulative Error Distribution curve for AFW dataset. All the experiments including training and testing were performed using the Caffe [12] framework and two Nvidia TITAN-X GPUs. Our method can process upto 3-4 frames per second, which can be higher in batch mode.

	AFLW	AFW
Method	NME	NME
TSPM [40]	-	11.09
CDM [2]	12.44	9.13
RCPR [7]	7.85	-
ESR [8]	8.24	-
PIFA [13]	6.8	8.61
3DDFA [39]	5.32	-
LPFA-3D [14]	4.72	7.43
EMRT [37]	4.01	3.55
CCL [38]	5.85	2.45
Rec Enc-Dec [1]	>6	-
<b>KEPLER</b>	<b>2.98</b>	3.01

TABLE I: Comparison of KEPLER with other state of the art methods. NME stands for normalized mean error. For AFLW, numbers for other methods are taken from respective papers following the PIFA protocol. For AFW, numbers are taken from respective works published following the protocol of [40].

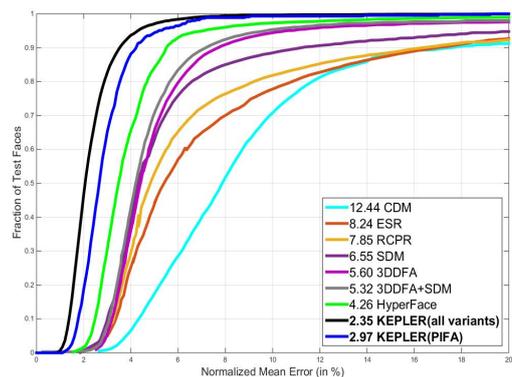


Fig. 7: Cumulative error distribution curves for landmark localization on the AFLW dataset. The numbers in the legend are the average normalized mean error normalized by the face size.

### B. Results

Table I compares the performance of KEPLER compared to other existing methods. Table II summarises the performance of KEPLER under different protocols of AFLW testset. Table III shows the mean error in degrees, in estimating the 3D pose of a face image. Figures 7 and 8 show the cumulative error distribution in predicting keypoints on the AFLW and AFW test sets. Figure 9 shows the cumulative error distribution in pose estimation on AFW.

**Comparison with CCL [38]:** It is clear from the tables that KEPLER outperforms all state of the art methods on the AFLW dataset. It also outperforms all state of the art methods except CCL [38] on the AFW dataset. Visual

inspection of our results suggests that KEPLER is a little farther from ground truth on invisible points. We note that CCL [38] manually annotates the AFLW dataset with 19 landmarks along with the invisible landmarks, leaving the earpoints. In our experiments we prefer to use the dataset as provided by AFLW [18], although we believe that CCL-kind of reannotation may boost the performance (since during AFW evaluation the locations of occluded points also need to be predicted). In KEPLER there is no loss propagated for the invisible points. We believe that training KEPLER on the revised annotation by [38] would make the prediction of occluded points more precise.

Method	AFLW-PIFA	AFLW-FULL	AFLW-Allvariants	AFW
<b>KEPLER</b>	2.98	2.90	2.35	3.01

TABLE II: Summary of performance on different protocols of AFLW and AFW by KEPLER.

Method	AFLW				AFW
	Yaw	Pitch	Roll	MAE	Accuracy ( $\leq 15^\circ$ )
Random Forest [29]	-	-	-	12.26°	83.54%
<b>KEPLER</b>	<b>6.45°</b>	<b>5.85°</b>	<b>8.75°</b>	<b>6.45°</b>	<b>96.67%</b>

TABLE III: Comparison of Mean error in 3D pose estimation by KEPLER on AFLW testset. For AFLW [29] only compares mean average error in Yaw. For AFW we we compare the percentage of images for which error is less than  $15^\circ$ .

We also verify our claim that iteration 5 is optional and only required for transferring the algorithm to other datasets with different bounding box annotations. To support our claim we calculate the normalized mean error after iteration 4 for both datasets and compare with the error obtained after iteration 5. The error after iteration 4 for AFLW testset was 0.0369 (which is already lower than all existing works) and after fifth iteration it was 0.0299, bringing the performance up by 18%. On the other hand the improvement in AFW (whose bounding box annotation is different from AFLW) was close to 60%. The error after iteration 4 on AFW dataset was 0.0757 which decreases to 0.0301 after fifth iteration. We demonstrate some qualitative results from AFLW and AFW test sets in Figure 10.

## V. CONCLUSIONS AND FUTURE WORKS

In this work we show that by efficiently capturing the structure of face through additional channels, we can obtain precise keypoint localization on unconstrained faces. We propose a novel *Channeled Inception* deep network which pools features from intermediate layers and combines them in the same manner to the Inception module. We show how cascade regressors can outperform other recently developed works. As a byproduct of KEPLER, 3D facial pose is also generated which can be used for other tasks such as pose dependent verification methods, 3D model generation and many others. In conclusion, KEPLER demonstrates that by improved initialization and multitask training, cascade regressors outperforms state of the art methods not only in predicting the keypoints but also for head pose estimation.

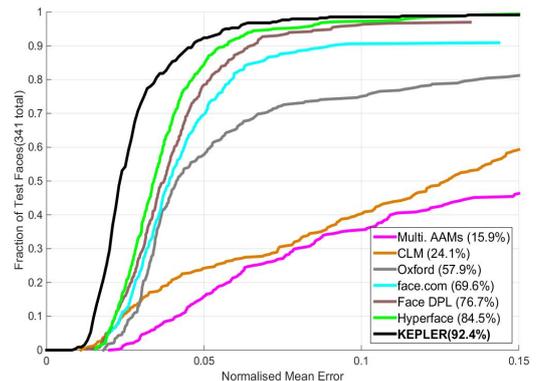


Fig. 8: Cumulative error distribution curves for landmark localization on the AFLW dataset. The numbers in the legend are the fraction of testing faces that have average error below (5%) of the face size.

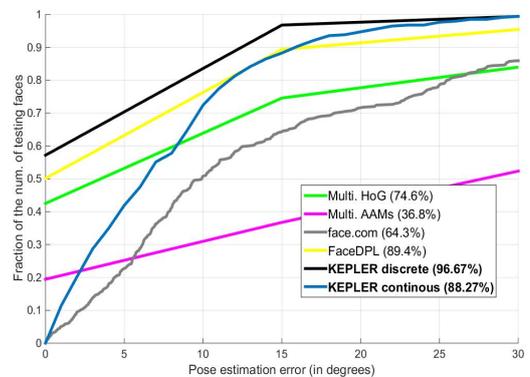


Fig. 9: Cumulative error distribution curves for pose estimation on AFW dataset. The numbers in the legend are the percentage of faces that are labeled within  $\pm 15^\circ$  error tolerance

## VI. ACKNOWLEDGMENT

This research is based upon work supported by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via IARPA R&D Contract No. 2014-14071600012. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon.

## REFERENCES

- [1] A recurrent autoencoder-decoder for sequential face alignment. <http://arxiv.org/abs/1608.05477>. Accessed: 2016-08-16.
- [2] *Pose-free Facial Landmark Fitting via Optimized Part Mixtures and Cascaded Deformable Shape Model*, 2013.
- [3] E. Antonakos, J. A. i medina, and S. Zafeiriou. Active pictorial structures. In *CVPR*, pages 5435–5444, Boston, MA, USA, June 2015.
- [4] E. Antonakos, P. Snape, G. Trigeorgis, and S. Zafeiriou. Adaptive cascaded regression. In *ICIP'16*, Phoenix, AZ, USA, September 2016.

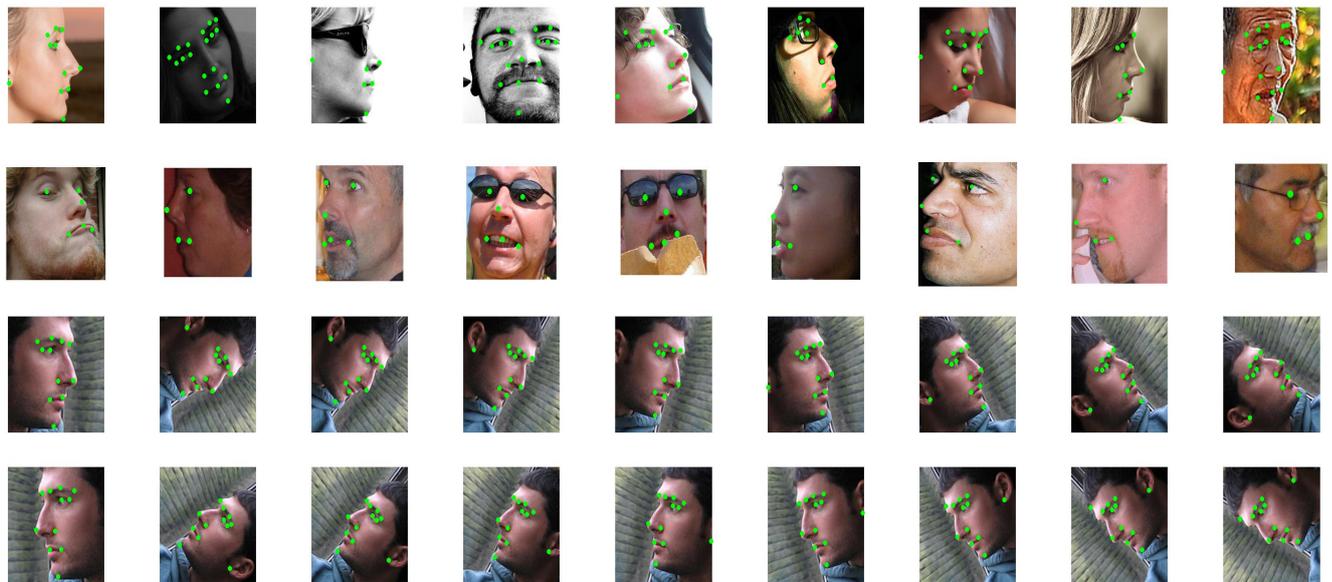


Fig. 10: Qualitative results of KEPLER after last stage. The green dots represent the final predicted points after fifth iteration. First row are the test samples from AFLW. Second row shows the samples from AFW dataset. The last two rows are the results of KEPLER after last stage from AFLW testset for all variants protocol. The green dots represent the final predicted points after fifth iteration.

- [5] A. Asthana, S. Zafeiriou, S. Cheng, and M. Pantic. Robust discriminative response map fitting with constrained local models. In *CVPR, CVPR '13*, pages 3444–3451, Washington, DC, USA, 2013. IEEE Computer Society.
- [6] A. Asthana, S. Zafeiriou, S. Cheng, and M. Pantic. Incremental face alignment in the wild. In *CVPR 2014*, 2014.
- [7] X. P. Burgos-Artizzu, P. Perona, and P. Dollar. Robust face landmark estimation under occlusion. *ICCV*, 0:1513–1520, 2013.
- [8] X. Cao, Y. Wei, F. Wen, and J. Sun. Face alignment by explicit shape regression. *IJCV*, 107(2):177–190, 2014.
- [9] J. Carreira, P. Agrawal, K. Fragkiadaki, and J. Malik. Human pose estimation with iterative error feedback. 2015.
- [10] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *CoRR*, abs/1502.01852, 2015.
- [11] G. S. Hsu, K. H. Chang, and S. C. Huang. Regressive tree structured model for facial landmark localization. In *ICCV*, Dec 2015.
- [12] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [13] A. Jourabloo and X. Liu. Pose-invariant 3d face alignment. In *ICCV*, Santiago, Chile, December 2015.
- [14] A. Jourabloo and X. Liu. Large-pose face alignment via cnn-based dense 3d model fitting. In *CVPR*, Las Vegas, NV, June 2016.
- [15] M. Kabkab, A. Alavi, and R. Chellappa. Dcnns on a diet: Sampling strategies for reducing the training set size. *CoRR*, abs/1606.04232, 2016.
- [16] V. Kazemi and J. Sullivan. One millisecond face alignment with an ensemble of regression trees. In *CVPR*, 2014.
- [17] B. F. Klare, B. Klein, E. Taborisky, A. Blanton, J. Cheney, K. Allen, P. Grother, A. Mah, M. Burge, and A. K. Jain. Pushing the frontiers of unconstrained face detection and recognition: larpa janus benchmark a. June 2015.
- [18] M. Koestinger, 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 *First IEEE International Workshop on Benchmarking Facial Image Analysis Technologies*, 2011.
- [19] A. Kumar, R. Ranjan, V. M. Patel, and R. Chellappa. Face alignment by local deep descriptor regression. *CoRR*, abs/1601.07950, 2016.
- [20] D. Lee, H. Park, and C. D. Yoo. Face alignment using cascade gaussian process regression trees. In *CVPR*, pages 4204–4212, June 2015.
- [21] L. Liang, R. Xiao, F. Wen, and J. S. . Face alignment via component-based discriminative search. In D. A. Forsyth, P. H. S. Torr, and A. Zisserman, editors, *ECCV*, volume 5303 of *Lecture Notes in Computer Science*, pages 72–85. Springer, 2008.
- [22] R. Ranjan, V. M. Patel, and R. Chellappa. Hyperface: A deep multi-task learning framework for face detection, landmark localization, pose estimation, and gender recognition. *CoRR*, abs/1603.01249, 2016.
- [23] 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.
- [24] Y. Sun, X. Wang, and X. Tang. Deep convolutional network cascade for facial point detection. In *CVPR, CVPR '13*, pages 3476–3483, Washington, DC, USA, 2013. IEEE Computer Society.
- [25] Y. Sun, X. Wang, and X. Tang. Deep convolutional network cascade for facial point detection. In *CVPR*, pages 3476–3483, June 2013.
- [26] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. E. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. *CoRR*, abs/1409.4842, 2014.
- [27] G. Trigeorgis, P. Snape, M. A. Nicolaou, E. Antonakos, and S. Zafeiriou. Mnemonic descent method: A recurrent process applied for end-to-end face alignment. In *CVPR*, Las Vegas, NV, USA, June 2016.
- [28] G. Tzimiropoulos and M. Pantic. Gauss-newton deformable part models for face alignment in-the-wild. In *CVPR*, pages 1851–1858, June 2014.
- [29] R. Valle, J. M. Buenaposada, A. Valdés, and L. Baumela. *Head-Pose Estimation In-the-Wild Using a Random Forest*, pages 24–33. Springer International Publishing, Cham, 2016.
- [30] Y. Wu and Q. Ji. Robust facial landmark detection under significant head poses and occlusion. In *ICCV*, pages 3658–3666, Dec 2015.
- [31] X. Xiong and F. D. la Torre. Global supervised descent method. In *CVPR*, 2015.
- [32] Xuehan-Xiong and F. De la Torre. Supervised descent method and its application to face alignment. In *CVPR*, 2013.
- [33] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. *CoRR*, abs/1311.2901, 2013.
- [34] J. Zhang, S. Shan, M. Kan, and X. Chen. Coarse-to-fine auto-encoder networks (cfan) for real-time face alignment. In D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, editors, *ECCV*, volume 8690 of *Lecture Notes in Computer Science*, pages 1–16. Springer International Publishing, 2014.
- [35] Z. Zhang, P. Luo, C. C. Loy, and X. Tang. Facial landmark detection by deep multi-task learning. In *ECCV*, pages 94–108, 2014.
- [36] S. Zhu, C. Li, C. Change Loy, and X. Tang. Face alignment by coarse-to-fine shape searching. June 2015.
- [37] S. Zhu, C. Li, C. C. Loy, and X. Tang. Towards arbitrary-view face alignment by recommendation trees. *CoRR*, abs/1511.06627, 2015.
- [38] S. Zhu, C. Li, C.-C. Loy, and X. Tang. Unconstrained face alignment via cascaded compositional learning. In *CVPR*, June 2016.
- [39] X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li. Face alignment across

large poses: A 3d solution. *CoRR*, abs/1511.07212, 2015.

- [40] X. Zhu and D. Ramanan. Face detection, pose estimation, and landmark localization in the wild. In *CVPR*, pages 2879–2886, June 2012.