# Face Detection with a 3D Model

Adrian Barbu, Nathan Lay, Gary Gramajo

✦

**Abstract**—This paper presents a part-based face detection approach where the spatial relationship between the face parts is represented by a hidden 3D model with six parameters. The computational complexity of the search in the six dimensional pose space is addressed by proposing meaningful 3D pose candidates by image-based regression from detected face keypoint locations. The 3D pose candidates are evaluated using a parameter sensitive classifier based on difference features relative to the 3D pose. A compatible subset of candidates is then obtained by non-maximal suppression. Experiments on two standard face detection datasets show that the proposed 3D model based approach obtains results comparable to or better than state of the art.

**Index Terms**—face detection

## 1 INTRODUCTION

In designing a good face detector we have mainly two choices. We might want a simple classifier (e.g. a sliding window classifier) and not care about the inner face representation or the face parts. In this case we would need to train it with features that have good discrimination power and tens of thousands of training examples to be able to cover the large variability of the faces in images due to 3D pose, illumination direction, face occlusions and other factors. For example, the face detector from [33] was trained with 100,000 faces and the detector

- *A. Barbu, and G. Gramajo are with the Department of Statistics, Florida State University, Tallahassee, Florida 32306, USA, Fax: 850-644-5271, Email: abarbu@stat.fsu.edu, ggramajo@stat.fsu.edu.*
- *N. Lay is with the National Institutes of Health, Bethesda, MD 20892, USA. Email: nathan.lay@nih.gov*

from [18] with more than 26,000 faces and their perturbations.

If we want a more interpretable model where the face parts are taken into consideration, then we are faced with the computational problem of enforcing dependencies between the face parts. A common way to handle this problem is through the DPM framework [10], and has been applied to face detection in many recent works such as [18], [28], [36]. Another way is by image-based regression, which has been used for face alignment in a number of works [3], [4], [9], [21] and has been also used for face detection recently in the Joint Cascade[5].

Besides being more interpretable, an added benefit of a part based model is that the face parts (eyes, mouth, nose, ears, chin, etc) have much smaller variability because they have simpler 3D shapes and appearances than the whole face. Thus a face detection system based on these parts could be trained with fewer training examples. For example the DPM model [36] has been trained with only 900 faces and obtains very good face detection results. The DPM based model uses 18 planar models to represent the part configurations for many possible out of plane face rotations. At test time, it tries all of them and returns the best scoring configurations above a threshold. This 2D approach leads to the question: what obstacles are there in using a single 3D face model instead of all these 2D models? We argue that the obstacles are mostly computational. The 2D models used in the DPM approaches have a tree structure so that dynamic programming

can be applied to obtain a globally optimal configuration. Thus the 2D models make some modeling compromises (many 2D tree-based models instead of a single 3D model) in order to be guaranteed that the global optimum is found. The 2D models used in the face alignment based approaches are not restricted to tree structure and use image-based regression to search for the optimal configuration in the high dimensional shape space.

In this paper we investigate an approach that uses a rigid 3D model to represent the interactions between the face parts and image based regression to search for it



Fig. 1. The face keypoints are fully connected by a simplex in our 3D model.

in images. The 3D model contains a rigid 6 parameter face 3D pose and independent deformations for the face parts from the locations predicted by the 3D pose, as illustrated in Figure 1. Many face parts can be missing, thus this model can be seen as an AND/OR graph [27] where the 3D face is the noisy AND of the face parts [34]. This model is a simplex, not a tree, and dynamic programming cannot be used for exact inference. Instead, the 6D pose space for each face will be searched using data-driven proposals made by image-based regression from the detected face part locations. These proposals are then evaluated by an energy model and the lowest energy configuration is the final result. In order for this approach to work well, the face parts need to be detected accurately.

**Contributions.** This paper brings the following contributions:

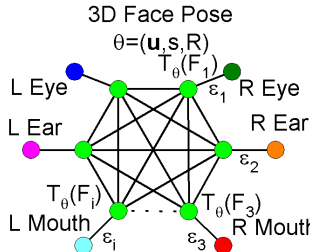– It presents a face representation that lies in a six-dimensional space, which is larger than the usual 3.1 dimensional space of position, scale, and a number of discrete poses. Such a larger dimensional space poses learning challenges to avoid overfitting and computational challenges in finding the hidden parameters quickly.

– It presents a novel computational paradigm that offers a combinatorial number of ways for detecting a face by attempting to detect many face parts and requiring only a few successful detections. This approach avoids relying on a cascade with a single computational path, which could fail if not trained appropriately.

– It introduces a novel set of local selected features that are computed based on the candidate 3D pose and have been selected during training from a larger pool by feature selection.

– It introduces a simple parameter sensitive model with a 1D parameter that is virtually as fast as its parameter insensitive counterpart and can be trained by stochastic gradient descent on millions of observations.

A detection cascade locks-in its losses at each level, since anything rejected at any level cannot be recovered. On the other hand, the noisy AND approach presented in this paper is more robust to detection failures than a cascade, as illustrated in the following example. If we assume that each face keypoint could be detected with probability at least 0.9, then the probability that a face has at most 4 face keypoints detected out of 9 is $\sum_{k=0}^{4} \binom{9}{k} \cdot 0.9^k \cdot 0.1^{9-k} < 0.001$. Thus if a face is considered detected when it has at least 5 out of 9 keypoints detected, then the probability for detecting a face based on its keypoints is at least 0.999, by the above computation. So even if the keypoint detections are not very reliable, it is unlikely that many of them will fail at the same time. At the same time, a cascade using one keypoint (e.g. face center) as an intermediate level would detect only 90% of the faces.

## 1.1 Related Work

There are different types of approaches to face or object detection and we will relate to the most relevant ones, even though this list might not be nearly close to complete.

**Multi-view models.** Some works approach the problem of detecting 3D objects with multiple viewpoints by having separate models for a range of possible views. In contrast, our work uses a single 3D model that controls the 2D positions of the features that are extracted for the face classifier, and the feature weights are also controlled by the 3D pose in the parameter sensitive face classifier.

The works [25], [18], [30] use rigid classifiers based on different types of predefined features such as Haar [25], HOG [6], ACF [8] and combinations. Very good results were obtained in [18] and [16] by training with many deformations of the positive examples.

Many works [10], [11], [36], [28] use deformable part-based models, where the relationships between parts are organized in a tree for obtaining a computationally tractable models. In this work, the parts are only related to the 3D pose, which connects all of them through a high-order simplex, as illustrated in Figure 1.

In [11] a 2D part configuration is detected using a version of the deformable part model [10] and then a 3D pose and shape is inferred from the 2D configuration. In contrast, our work directly uses the 3D pose to represent the relative positions of the parts without going through an intermediate 2D model.

**3D view based models.** Some works [20], [24] divide the view sphere into a number of sectors and collect templates for each view. Given a new candidate object, detection is obtained by template matching. Our work is not template based, but is based on a parameter sensitive classifier that uses features extracted relative to a 3D model.

**3D models.** Our work resembles [17], [12] in that features are extracted based on a 3D model and the object pose hypothesis. However, these approaches use complex inference algorithms (one based on EM and the other based on dynamic programming) while our approach uses regression to propose data-driven candidates from multiple channels. Furthermore,

our approach does not need any synthetic 3D models, since it constructs the 3D model from training images by energy minimization. Moreover, none of these works was used for face detection.

**Cascade approaches.** Some of the most successful face detection approaches [5], [16] are based on a detection cascade. Since a cascade locks in the losses at each stage, it has an upper limit on the detection rate. For example, the Cascade-CNN cannot reach an 86% detection rate on FDDB. Our approach is not based on a cascade but has candidates proposed through multiple channels. The combinatorial verification (e.g. requiring four face parts detected out of 9) gives more flexibility and allows many part detectors to fail while still detecting a face successfully. Consequently, our face detector can reach higher detection rate, as high as 89% or more on FDDB.

**Face alignment.** Pose candidates have been previously proposed by image based regression in the shape regression machine [35] and for face alignment [3], [4], [9], [21], however, they are not based on a 3D model. The Cascade-CNN [16] uses a convolutional neural network (CNN) to improve alignment of the detected face bounding boxes, an approach also not based on a 3D model.

Our work generalizes the Local Binary Features [21] to local (non-binary) versions that are extracted relative to our 3D model. Moreover, we use a parameter sensitive classifier for scoring the candidates instead of a standard classifier.

The Joint Cascade [5] obtains very good face detection results by alternating classification with face alignment by regression, both steps using the LBF features [21]. Our work differs in many ways. First, we use a 3D model instead of a 2D model. Second, we detect many keypoints in a bottom-up step and use them to propose many 3D pose candidates, instead of starting with a "mean face", this way we can obtain higher detection rates. Third, we use a parameter sensitive classifier instead of

a standard classifier, which adapts the feature weights to the 3D pose. It is possible that by using one or more face alignment steps we could further improve the detection rate in the low false positive region.

**Parameter sensitive classifiers.** Parameter sensitive classifiers were introduced in [31] for Boosting and linear SVM, and in [32] for SVM with multiplicative kernels. While the multiplicative kernel formulation is generic and can be used with multi-dimensional parameters, it might be too computationally expensive for a face verification classifier. This is why we introduced a simple formulation for the one-dimension parameter representation that can be solved by direct energy minimization.

**Face detection with pose estimation.** Another notable work is the face detection and pose estimation with energy based models [19] which uses a Convolutional Neural Network to directly map the input image patch into a pose manifold for faces and outside the manifold for non-faces. It would be interesting to see how this work compares to the current state of the art methods on the FDDB and AFW datasets.

## 2   PARAMETER SENSITIVE MODEL

Parameter sensitive models are models that take a feature vector $\mathbf{x} \in \mathbb{R}^M$ and a parameter $\theta \in \mathbb{R}^p$ to obtain a prediction $s(\mathbf{x}; \theta)$.

**Parameter sensitive linear model.** For example the dependence on $\mathbf{x}$ could be linear

$$s(\mathbf{x}; \theta) = \mathbf{w}(\theta)^T \mathbf{x} = \sum_{i=1}^{M} w_i(\theta) x_i \qquad (1)$$

with the linear weights $\mathbf{w}(\theta)$ depending on the parameter $\theta$. Such a model is more computationally efficient than a SVM classifier with multiplicative kernels [32] since it avoids kernel multiplication with many support vectors.

In our application the parameter $\theta$ is one dimensional, so we discretize the range of $\theta$ into $B_\theta$ equal intervals $[b_i, b_i + 1]$ to represent $w_i(\theta)$ as piecewise constant $w(\theta) = w_{ik}$ for $\theta \in [b_k, b_{k+1})$. These parameters are collected in the matrix $W = (\mathbf{w}_1, ..., \mathbf{w}_M), \mathbf{w}_i \in \mathbb{R}^{B_\theta}$.

**Non-linear model.** We can introduce nonlinearity without kernels as a sum of bivariate parameter sensitive response functions:

$$s(\mathbf{x}; \theta) = \sum_{i=1}^{M} f_i(x_i, \theta) \qquad (2)$$

where the function $f_i(x_i, \theta)$ depends on one variable $x_i$ of vector $\mathbf{x}$ and the parameter $\theta$. An example of a bivariate response function $f_i(x_i, \theta)$ that make up the nonlinear model (2) is shown in Figure 2.
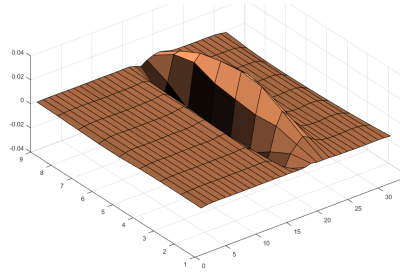


Fig. 2. One of the bivariate response functions $f_i(x_i, \theta)$ for the nonlinear model (2).

By discretizing the range of $\theta$ into $B_\theta$ bins as above and the range of $x_i$ into $B_\mathbf{x}$ bins, we use piecewise constant bivariate response functions $f_i(x_i, \theta) = w_{ijk}$ if $x_i$ is in bin $j$ and $\theta$ is in bin $k$. Again, these parameters are collected in the matrix $W = (\mathbf{w}_1, ..., \mathbf{w}_M)$ with $\mathbf{w}_i \in R^{B_\theta B_\mathbf{x}}$.

In both cases we write the parameter sensitive model as $s_W(\mathbf{x}; \theta)$ to emphasize that its parameters form the matrix $W$.

### 2.1   Training the Parameter Sensitive Model

The training examples $(B_i, \mathbf{x}_i, \theta_i), i = \overline{1, N}$ contain rectangles (boxes) $B_i$, feature vectors $\mathbf{x}_i$ and the parameters $\theta_i$. The ground truth is given as a set of face boxes $f_j$. Let $o_{ij}$ be the overlap of box $B_i$ with face box $f_j$ if they belong to the same image, otherwise is zero. Let $o_i = \max_j o_{ij}$ be the maximum overlap of box $B_i$ with a face from the same image.

**Training Cost Function.** Inspired by [26], the cost function to be minimized for training:

$$E(W) = \sum_{j} \sum_{i,o_{ij}>0.5} \frac{2o_{ij}-1}{s_j} L(s_W(\mathbf{x}_i, \theta_i))$$
$$+ \mu \sum_{i,o_i<0.3} L(-s_W(\mathbf{x}_i, \theta_i)) + \sum_{i=1}^{M} \rho(\mathbf{w}_i) \quad (3)$$

contains sums of loss terms for each true positive and can be written as a sum of per-image loss terms. The difference from [26] is that the positive boxes (that overlap at least 0.5 with a face) are weighted based on their overlap and then normalized by $s_j = \sum_{i,o_{ij}>0.5}(2o_{ij}-1)$ so that the total weight is 1 for each true face.

The prior $\rho(\mathbf{w})$ encourages smooth changes of the parameters between parameter bins

$$\rho(\mathbf{w}) = s\|\mathbf{w}\|^2 + c\sum_{k=2}^{B_\theta-1}(w_{k+1}+w_{k-1}-2w_k)^2 \quad (4)$$

for the linear model and of parameter bins and variable bins for the nonlinear model:

$$\rho(\mathbf{w}) = s\|\mathbf{w}\|^2 + c\sum_{j=1}^{B_\mathbf{x}}\sum_{k=2}^{B_\theta-1}(w_{j,k+1}+w_{j,k-1}-2w_{jk})^2$$
$$+ d\sum_{k=1}^{B_\theta}\sum_{j=2}^{B_\mathbf{x}-1}(w_{j+1,k}+w_{j-1,k}-2w_{jk})^2 \quad (5)$$
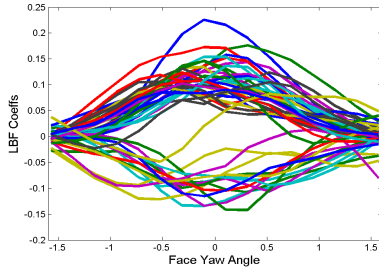


Fig. 3. Top 50 coefficients by total variation, for a parameter sensitive linear model.

An example of top learned coefficients for the LBF features is given in Figure 3, where each curve represents the coefficient of one keypoint across 16 yaw angle bins.

**Optimization.** Observe that the cost function $E(W)$ is convex when the loss function $L(x)$ is convex (such as the logistic or the hinge loss).

However, we will use the Lorenz loss function [1]

$$L(x) = \begin{cases} \ln(1+(x-1)^2) & \text{if } x < 1 \\ 0 & \text{else} \end{cases} \quad (6)$$

which is more robust to outliers than the logistic and hinge losses and was observed to work very well in this application.

Optimization is done by Feature Selection with Annealing [1], a fast iterative procedure that alternates parameter updates with progressive removal of variables according to a schedule. The parameter update step consists of one epoch of stochastic gradient descent with momentum 0.9, where the a batch contains all training examples from one image, thus the parameters are updated one image at a time.

The variable removal schedule is linear, after each parameter update epoch the same number of variables are removed, such that the desired number of variables are left after 20 epochs. In our application we start with 9000 variables and are left with 3000 after 20 epochs.

## 3 FACE DETECTION USING A 3D MODEL

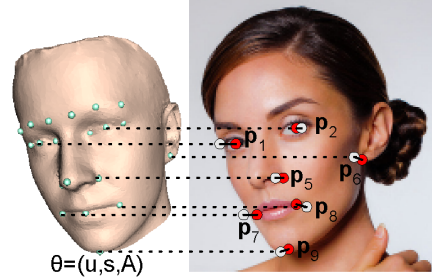Given an image, the goal is to find the faces with keypoints and 3D poses.



Fig. 4. A face is represented as a pair of a 3D pose $\theta = (\mathbf{u}, s, A)$ and a 2D point configuration $P = (\boldsymbol{p}_1, ..., \boldsymbol{p}_L)$.

**Face Representation.** The face has $L$ keypoints that are 2D points in the image, $P = (\boldsymbol{p}_1, ..., \boldsymbol{p}_L), \boldsymbol{p}_i \in \mathbb{R}^2$. The face 3D pose is represented as a projected rigid transformation $T_\theta : \mathbb{R}^3 \to \mathbb{R}^2$ with parameters $\theta = (\mathbf{u}, s, A)$ consisting of 2D translation $\mathbf{u} \in \mathbb{R}^2$, scale $s$ and 3D rotation matrix $A$, and defined as $T_\theta(\mathbf{x}) = \mathbf{u} + s\pi(A\mathbf{x})$, where $\pi : \mathbb{R}^3 \to \mathbb{R}^2, \pi(x,y,z) = (x,y)$ is the projection on the $(x,y)$ plane. Thus
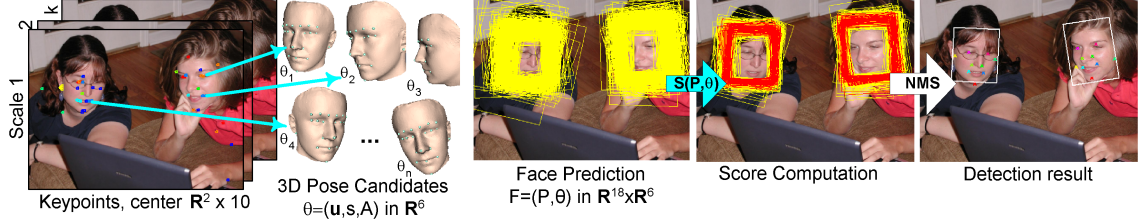
Fig. 5. Face detection using a 3D model. The face keypoints are detected independently and used to propose 3D pose candidates $\theta = (\mathbf{u}, s, A) \in \mathbb{R}^6$. Face candidates are predicted from the rigid model and evaluated using the score $S(P, \theta)$. The detected faces are obtained by non-max suppression.

each face is represented as a pair $F = (P, \theta)$ of the 2D keypoints $P = (\boldsymbol{p}_1, ..., \boldsymbol{p}_L)$ and the 3D pose $\theta = (\mathbf{u}, s, A)$.

**Face 3D model.** The 3D face model consists of $L$ 3D keypoints in a rigid configuration that can be written as a $3 \times L$ matrix $R = (\boldsymbol{r}_1, ..., \boldsymbol{r}_L), \boldsymbol{r}_i \in \mathbb{R}^3$. The 2D configuration $P$ of the face keypoints in the image is related to the pose $\theta$ through the relation

$$\boldsymbol{p}_i = T_\theta(\boldsymbol{r}_i) + \epsilon_i, i = \overline{1, L}$$

where $T_\theta$ is the 3D face pose defined above and $\epsilon_i \in \mathbb{R}^2$ are independent deformations for each keypoint. We will write this in matrix form as:

$$P = T_\theta(R) + \boldsymbol{\epsilon} \qquad (7)$$

This relation is illustrated in Figure 4 where the gray dots are the predicted point locations $T_\theta(\boldsymbol{r}_i)$ and $\boldsymbol{p}_i$ are the actual point locations.

### 3.1 Energy Model

For any face $F = (P, \theta)$ let $B_F$ be the bounding box of the points $P$. The best configuration of faces is obtained by energy minimization:

$$(F_1, ..., F_n) = \underset{n, F_1, ..., F_n}{\operatorname{argmin}} (E_{data}(F_1, ..., F_n)$$
$$+ E_{prior}(n, F_1, ..., F_n))$$

The data term

$$E_{data}(F_1, ..., F_n) = \sum_{j=1}^{n} (\tau - S(P_j, \theta_j))$$

is based on the scores $S(P_j, \theta_j)$ for the faces $F_j = (P_j, \theta_j)$ in the image and the parameter $\tau$ that controls the minimum score for a valid detection. The score function $S(P, \theta)$ is defined in more detail in Section 3.2.4. The prior

$$E_{prior}(n, F_1, ..., F_n) = \sum_{j=1}^{n} f(\|P_j - T_{\theta_j}(R)\|) \qquad (8)$$
$$+ E_{ovr}(n, P_1, ..., P_n)$$

has a coherence term between the poses $\theta_j$ and the points $P_j$ of the face $F_j = (P_j, \theta_j)$ and a term $E_{ovr}(n, P_1, ..., P_n)$ that enforces the constraints that the bounding boxes $B_{F_j}, j = \overline{1, n}$ have small overlap with each other.

### 3.2 Inference Algorithm

The inference algorithm is illustrated in Fig. 5 and consists of the following steps:

1) Bottom-up detection of the face keypoints.
2) Generation of 3D pose candidates $\theta_1, ..., \theta_n$ from the keypoints.
3) Generation of the face candidates by pruning the 3D pose candidate $\theta_j$ using the coherence term.
4) Computation of the face scores $S(P_j, \theta_j), j = \overline{1, n}$ and removal of low scoring candidates.
5) Non-maximal suppression to output a set of high score candidates that satisfy the overlap constraints, greedily minimizing the energy $E(n, \theta_1, ..., \theta_n)$.

These steps are described in the next subsections.

#### 3.2.1 Detecting face keypoints

The $L = 9$ face keypoints used in this work are the eye centers, nose sides, mouth corners, bottom ears, and chin. We also detected the center of the face bounding box to handle small faces

where the keypoints are not clearly visible. The keypoints are detected on a Gaussian pyramid with 4 scales per octave and minimum image size of $24 \times 24$.
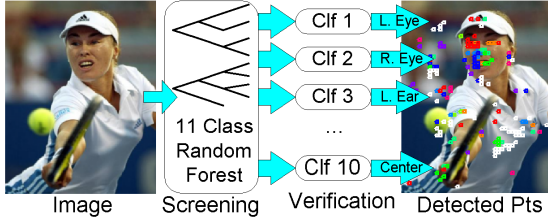


Fig. 6. Keypoint Detection involves class prediction with a Random Forest and verification with binary classifiers (one class vs. background).

The face keypoints are detected in two stages as shown in Figure 6. First, a 11-class (background, 9 keypoints, and face center) Random Forest with 100 trees of depth 11 prunes the search space. The trees use Aggregate Channel Features (ACF)[8] with 10 channels and block size of 2 in a window of size $24 \times 24$.

A FSA classifier [1] with 3000 features selected from a pool of 61000 Haar and ACF features in a $24 \times 24$ window is trained for each of the 9 keypoints and the face center, to minimize the Lorenz loss (6) and 10 iterations of hard negative mining [10]. These detectors have a detection rate on the training set of $90 - 95\%$, a false positive rate of $0.1 - 1\%$, and are used to verify the non-background classes proposed by the RF classifier. Other classifiers such as based on CNN could also be used.

### 3.2.2 Generating 3D Pose Candidates

Since the keypoints are detected for faces in a range of scales, the pose candidates are also obtained for faces in the same range. The 3D pose candidates are generated by image-based regression from the detected keypoint locations. The 3D pose $\theta = (\mathbf{u}, s, A)$ has six parameters $(\mathbf{u}, s, \boldsymbol{\varphi}) = (u^x, u^y, s, \varphi^x, \varphi^y, \varphi^z)$, where $\boldsymbol{\varphi} = (\varphi^z, \varphi^x, \varphi^y)$ is the roll-pitch-yaw decomposition of the rotation matrix $A$.

**Image based regression.** The 3D pose is predicted from a point $(x_0, y_0, s_0)$ by image based regression, predicting the relative vector

$$(\frac{u^x}{s_0} - x_0, \frac{u^y}{s_0} - y_0, \frac{s}{s_0}, \varphi^x, \varphi^y, \varphi^z) = \mathbf{y}(\mathbf{x}). \quad (9)$$

The image based regression is trained using a feature vector $\mathbf{x} = (x_1, ..., x_m)$ selected from the same feature pool of $M = 61,000$ features used in Section 3.2.1 for face keypoint verification. The range of each variable is divided into 32 equal bins and let $b_i(x)$ be the bin index function for variable $i$. The regression vector $\mathbf{y}(\mathbf{x})$ is obtained as a sum of piecewise constant functions that depend on one variable each, i.e.

$$\mathbf{y}(\mathbf{x}) = \sum_{i=1}^{m} \boldsymbol{z}_{i, b_i(x_i)}$$

where the $\boldsymbol{z}_{ib}$ is the 6D coefficient vector for variable $i$ and bin $b$. All the coefficient vectors are collected into a 3 dimensional matrix $Z$ of size $6 \times m \times 32$ that will be estimated from training examples.

**Ground truth 3D pose.** The ground truth 3D poses are obtained by least squares energy minimization as described in section 4.2. The POSIT algorithm [7] could also be used for this purpose. Then the ground truth vectors for training the 3D pose regressors are obtained as in eq. (9) for each annotated face from the fitted 3D pose $(\mathbf{u}, s, \boldsymbol{\varphi}) = (u^x, u^y, s, \varphi^x, \varphi^y, \varphi^z)$ and the keypoint location $(x_i, y_i, s_i)$.

**Training details.** The training examples are in the form of $(\mathbf{x}_i, \mathbf{y}_i) \in \mathbb{R}^M \times \mathbb{R}^6$ where $\mathbf{x}_i$ are feature vectors extracted at locations within 1 pixel from true keypoint locations and $\mathbf{y}_i$ are the relative vectors (9) based on the ground truth poses obtained as described above. Training is done by minimizing the energy

$$L(Z) = \sum_{j=1}^{n} \|\mathbf{y}_j - \sum_{i=1}^{m} \boldsymbol{z}_{i, b_i(x_{ji})}\|^2 \quad (10)$$

using the FSA algorithm [1] where $m = 2000$ features are selected from the $M = 61,000$ feature pool. A specific 6D pose regressor is trained for each keypoint for better accuracy. Observe that by using the loss function (10) the same 2000 features are used for predicting all six pose parameters, instead of using 2000

(possibly different) features for predicting each parameter.

The percentage of variance explained $R^2$ for the pose regression dimensions ranges between 0.3 and 0.8, where the lowest scores were for predicting the pitch and roll angles. We choose FSA because it produced better 3D pose candidates than Random Forest regression, as will be seen in Sections 5.1 and 5.2.

### 3.2.3 Generating Face Candidates

The 3D pose candidate generation step from the previous section obtains a number of 3D pose candidates $\theta_1, ... \theta_n$. The face candidates need to contain besides the 3D poses the predicted locations in 2D of the face keypoints.

From a 3D pose $\theta$ the keypoint locations $P = (\boldsymbol{p}_1, ..., \boldsymbol{p}_L)$ are first predicted directly from the rigid 3D model as $P = T_\theta(R)$ , even though this prediction might not be very accurate. The IED (inter-eye distance) of this 3D candidate can also be computed from the 3D position of the eyes before projection.

**Keypoint support.** The coherence term from eq. (8) is based on the number of keypoints that support the candidate. The support for candidate $(P, \theta)$ at some scale of the image pyramid is the number of rescaled $\boldsymbol{p}_i$ that have a keypoint detection of type $i$ in the image at that scale at distance at most $0.1 \cdot$IED (also rescaled to that scale). The overall support for the candidate is the maximum (over all scales) of the candidate support at one scale. The face keypoints that support the candidate become the new positions of the corresponding $\boldsymbol{p}_i$.

The face candidates with a support less than a value $N^{supp}$ are eliminated. The value of $N^{supp}$ is usually 4 in our experiments, unless otherwise specified. The number of generated face candidates for a $480 \times 320$ image ranges ranges from about 4000 for $N^{supp} = 1$ to about 1200 for $N^{supp} = 4$.

### 3.2.4 Scoring the Face Candidates

A face candidate $F = (P, \theta)$ contains the 3D pose $\theta = (\mathbf{u}, s, A)$ and predicted locations

$P = (\boldsymbol{p}_1, ..., \boldsymbol{p}_L)$ of the L keypoints. The score is obtained by a parameter-sensitive model depending on the yaw angle $\varphi^y$, as described in Section 2.

**Local difference features.** Inspired by the Local Binary Features (LBF) [21] we generate a number of difference features relative to the 3D pose $\theta$ of a face candidate $F = (P, \theta)$ instead of the 2D shape. For that, an approximate tangent plane at each keypoint of the 3D face model is assigned a system of coordinates. This coordinate system is projected to a 2D coordinate system based on the 3D pose, which is used to define a skewed point grid centered at the predicted keypoint location $\boldsymbol{p}_i$ as illustrated in Figure 7. On a $11 \times 11$ grid we obtain about
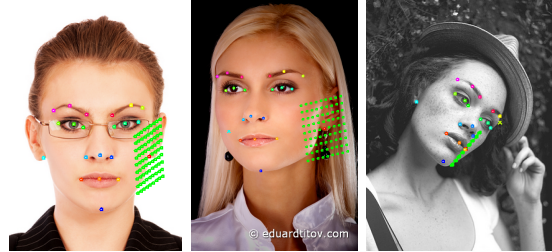


Fig. 7. Sampling grid patterns based on the 3D face pose used to generate local features.

72,000 such difference features around each keypoint.

**Modified LBF features.** Modified LBF features can be obtained as the leaf indexes for the Random Forest regression trees trained with these local difference features. The LBF were trained with 100 Random Trees of depth 6 for each of the 9 keypoints, for a total of $100 \times 32 \times 9 = 28,800$ features.

**Local selected features (LSF).** Another approach we will use is to first prune the 72,000 local difference features at each keypoint to a more manageable size (e.g. 1,000) which will be the input of the parameter sensitive model described in Section 2. The pruning of the local difference features is done by FSA regression to predict local translations of the grid center on the training examples, similar with the LBF training but using FSA instead of RF for regression.

**Special features.** For each face candidate we also created 2 special features for each keypoint including the face center. First is the distance (as percentage of the inter-eye distance of the candidate) between the predicted keypoint location of the candidate and the closest detected keypoint to that location. If there is not such detected keypoint or the distance is larger that $0.32$ IED, the value is $0.32$. The second feature is the detection score of the closest detected keypoint, with a minimum of $-1.1$. For the 9 keypoints plus face center we obtain this way 20 special features. We also added the candidate support as the 21-st special feature.

In both cases, let $\mathbf{x}(P, \theta)$ be the vector of LBF or LSF features+special features extracted from the image for the candidate $F = (P, \theta)$.

**Score function.** The score $S(P, \theta)$ of the 3D face candidate with face keypoints $P$ and pose $\theta$ is based on the feature vector $\mathbf{x}(P, \theta)$, either by a linear model (1) for the LBF features or by a nonlinear model (2) for the LSF features.

The models depend parametrically on the yaw angle $\varphi^y$ of the rotation $A$. The yaw angle ranges between $-\pi/2$ and $\pi/2$, being $0$ for frontal faces and $\pm\pi/2$ for profile faces. For this application, it is discretized into $B_\theta = 16$ bins for the LBF features and $B_\theta = 9$ bins for the LSF features.

### 3.2.5  Non-Maximal Suppression

The non-maximal suppression step iterates the following until convergence:

1. Select the face candidate $(P, \theta)$ with the largest score above a threshold and finds the bounding box $B$ of its points $P$.

2. Remove the candidates that have at least $50\%$ overlap with $B$.

## 4  FITTING 3D MODELS

In this section we show how the 3D model is fitted to and learned from 2D annotations.

### 4.1  Fitting a Rigid Projection Transformation.

Given a matrix $3 \times L$ matrix $F$ and a set of 2D points $P = (\boldsymbol{p}_1, ..., \boldsymbol{p}_L)$ in the form of a $2 \times L$ matrix, the goal is to find a rigid transformation $\theta = (\mathbf{u}, s, R)$ to minimize:
$$E(\mathbf{u}, s, R) = \|\mathbf{u}\mathbf{1} + s\pi(RF) - P\|^2$$
where $\pi((\mathbf{x}, \mathbf{y}, \boldsymbol{z})^T) = (\mathbf{x}, \mathbf{y})^T$ and $\mathbf{1} \in \mathbb{R}^L$ is the row vector with all entries $1$.

The algorithm uses hidden variables for the $z$ coordinates of the points $\boldsymbol{p}_i$ and iterates fitting the rigid transformation with updating the $z$-values.

---
**Algorithm 1 Fit Rigid Projection**

---
**Input:** $3 \times L$ matrix $F$ and $2 \times L$ matrix $P$.
**Output:** Scalar $s$, $3 \times 3$ rotation matrix $R$ and 2D vector $\mathbf{u}$ to minimize $\|\mathbf{u}\mathbf{1} + s\pi(RF) - P\|^2$

1: Initialize $L \times 3$ matrix $B = (P^T, 0)$.
2: **for** i=1 to $N^{iter}$ **do**
3:   Call Algorithm 2 to find $\mathbf{u}, s, R$ to minimize $\|\mathbf{1}^T\mathbf{u}^T + sF^TR - B\|^2$
4:   Extract third column $\boldsymbol{c}_3 = (C_{i3})_i$ of $C = sF^TR$
5:   Update $B = (P^T, \boldsymbol{c}_3)$
6: **end for**
7: Change $R$ to $R^T$ and discard the $z$-component of $\mathbf{u}$.

---

The algorithm to fit a rigid transformation between two sets of points of the same dimension $d$ is due to Schonemann [22] and is presented in Algorithm 2.

---
**Algorithm 2 Fit 3D Rigid Transformation**

---
**Input:** Matrices $A, B$ of size $p \times d$.
**Output:** Scalar $s$, $d \times d$ rotation matrix $R$ and $d \times 1$ vector $\mathbf{u}$ to minimize $\|\mathbf{1}^T\mathbf{u}^T + sAR - B\|^2$

1: Compute the column means $\bar{\alpha} = \mathbf{1}A/p$, $\bar{\beta} = \mathbf{1}B/p$ and column centered matrices $A^* = A - \mathbf{1}^T\bar{\alpha}$ and $B^* = B - \mathbf{1}^T\bar{\beta}$.
2: Decompose $A^{*T}B^* = UDV^T$ by SVD, where $U, V$ are rotation matrices and $D$ is a diagonal matrix.
3: Obtain $R = UV^T$, $\mathbf{u} = \bar{\beta} - s\bar{\alpha}R$ and $s = tr[R^TA^{*T}B^*]/tr(A^{*T}A^*)$

---

Given two sets of points $A, B$ of the same dimension $d$, it finds a rigid transformation $(\mathbf{u}, s, R)$ represented by a translation vector $\mathbf{u}$, scaling $s$, and rotation matrix $R$ and to minimize $\|\mathbf{1}^T\mathbf{u}^T + sAR - B\|^2$.

TABLE 1
Evaluation of face candidates on the FDDB dataset.

| Experiment Number | Method Name | Keypoints | Keypoint Verification | Pose Regression Method | $N^{supp}$ | False Positive Rate %<0.3 | %<0.5 | Detection Rate %>0.5 | %>0.7 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | S0S0 | Center only | FSA | FSA | 0 | 60.2 | 77.5 | 96.6 | 90.8 |
| 2 | R9S3 | 9+Center | - | FSA | 3 | 47.7 | 70.6 | 90.0 | 86.4 |
| 3 | RS9S2 | 9+Center | FSA | FSA | 2 | 28.8 | 52.1 | 97.1 | 95.0 |
| 4 | RS9S3 | 9+Center | FSA | FSA | 3 | 9.2 | 30.4 | 94.0 | 91.7 |
| 5 | RS9S4 | 9+Center | FSA | FSA | 4 | 2.3 | 16.6 | 90.5 | 88.1 |
| 6 | RS9R3 | 9+Center | FSA | RF | 3 | 12.1 | 35.5 | 94.0 | 90.8 |

## 4.2 Learning a 3D Model from 2D Annotations

The face 3D model matrix $R$ is obtained from a number of 2D face images where the $L$ keypoints have been annotated.

Let $P_i = (p_{i1}, ..., p_{iL})$ be the 2D coordinates of the $L$ keypoints for face $i, i = \overline{1, n}$. Write $P_i = \begin{pmatrix} \mathbf{x}_i \\ \mathbf{y}_i \end{pmatrix}$, obtaining the row vectors $\mathbf{x}_i, \mathbf{y}_i$ as the $x$ and $y$ coordinates of the $L$ keypoints.

The goal is to find the matrix $F$ of size $3 \times L$ and the projected rigid transformations $\theta_i, i = \overline{1, n}$ for the annotated faces, to minimize

$$E(\theta, F) = \sum_{i=1}^{n} \|\mathbf{u}_i \mathbf{1} + s_i \pi(R_i F) - P_i\|^2$$
$$= \sum_{i=1}^{n} \|T_i^x F - X_i\|^2 + \sum_{i=1}^{n} \|T_i^y F - Y_i\|^2 \quad (11)$$

where each $T_{\theta_i}$ is a $2 \times 3$ matrix with rows $T_i^x, T_i^y$.

The minimization starts with a random F and alternates two steps until convergence:

1) Given current $F$, fit the projected rigid transformations $\theta_i$ using Algorithm 1 for each face
2) Given current $\theta_i$, find F by minimizing eq. (11)

This approach is a simplified version of the Stratified Procrustes Analysis [2].

## 5 EXPERIMENTS

**Method nomenclature.** The proposed method consists of three main parts: a) face keypoint and center detection, b) 3D pose regression+ support pruning and c) face candidate evaluation. Since each has a number of variants, we will represent each combination with a shortcut:

$$\underbrace{RS9}_{a)} \underbrace{S4}_{b)} \underbrace{PLSF}_{c)}$$

where

a) Keypoint detection can be R (Random Forest), S (FSA) or RS (Random Forest +FSA), and 9 is the number of keypoints that will be detected (plus face center). The RF screening and the FSA verification are covered in section 3.2.1.
b) 3D Pose regression can be R (Random Forest) or S (FSA), as described in Section 3.2.2. The number 4 comes from pruning the 3D candidates by the support, as described in Section 3.2.3, thus $N^{supp} = 4$.
c) Face candidate scoring (verification) has been described in Section 3.2.4, and it can be PLSF = Parametric sensitive model with LSF+special features, PLBF = parametric sensitive model with LBF features, or LBF= standard LBF classifier.

**Training dataset.** For training we used 11,000 images from the AFLW dataset [14], containing about 14,300 faces. All models (keypoint detectors, 3D model, LBF features, parameter sensitive classifier, etc) were trained on these images and their 21 point annotation.

### 5.1 Evaluation of Face Candidates

Before evaluating the whole system, we first evaluate the face candidate generator to get an idea of what to expect from the face verification step. For a face candidate $F = (P, \theta)$ with 3D pose $\theta$ one can compute a face bounding box based on the predicted keypoints $P$. The face candidates can be evaluated by computing the overlap of the face bounding boxes with the ground truth face bounding boxes.
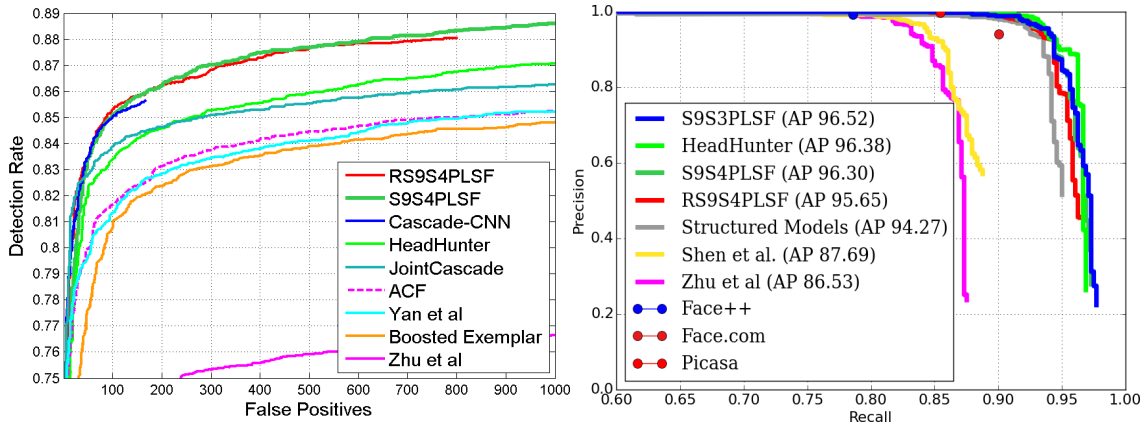
Fig. 8. Results and comparisons with other methods. Left: the FDDB dataset ( 2845 images with 5171 faces). Right: the AFW dataset (205 images with 552 faces).

In Table 1 are shown two measures of "false positive rate": the percentage of candidate boxes with overlap $< 0.3$ or $< 0.5$ with the ground truth face bounding boxes and two "detection rate" measures: the percentage of GT faces that have candidate bounding boxes with overlap $> 0.5$ or $> 0.7$.

To see the importance of the keypoint detections, experiment 1 shows the face candidates obtained only from the face center and no support pruning. Experiment 4 has the candidates predicted from all nine keypoints plus the center and $N^{supp} = 3$, clearly better than the ones from the center only in both higher detection rate (91.7% vs 90.8%) and lower false positive rate (9.2% vs 60.2%). Experiment 2 shows the candidates obtained without the keypoint verification step from Figure 6 are clearly inferior to Experiment 4. Experiment 6 shows the pose regression based on a Random Forest with 100 trees of depth 10 using the same features as the FSA pose regression. One can see that the FSA pose regression from Experiment 4 obtains higher detection rate (91.7% vs 90.8%) and lower false positive rate (9.2% vs 12.1%).

## 5.2 Face Detection Results

We present results on two standard face detection datasets: The FDDB dataset [13] with 2845 images containing 5171 faces and the AFW dataset [36] with 205 images and 552 faces.

The FDDB evaluation used the evaluation code provided on the FDDB website. The AFW evaluation used the code provided by [18]. In both cases a detection with less than 50% overlap with any annotated face is considered a false positive.

The results on the FDDB dataset are shown in Figure 8, left. Our results are: "S9S4PLSF" and "RS9S4PLSF" with LSF+special features with FSA respectively RF+FSA keypoint detection, FSA pose regression and $N^{supp} = 4$. Also shown are results from the Cascade-CNN [16], Joint Cascade [5], HeadHunter [18], Boosted Exemplar [15], ACF [30], Yan et al [28] and Zhu [36]. One can see that the proposed method obtains results comparable to the state of the art for up to 100 false positives and outperforms the other methods for the regime with at least 100 false positives. In Figures 11 and **??** are shown detection results on a few images of the FDDB dataset.

The results on the AFW dataset are shown in Figure 8, right. Also shown are results from the Head Hunter [18], Shen et al [23], Structured Models [29] and Zhu [36]. The algorithm performs well in the high recall regime and is outperformed by the Head Hunter [18] in the high precision regime.

**Evaluation of design decisions.** In Figure 9, are shown evaluations to support the decisions to use multiple keypoints, FSA pose
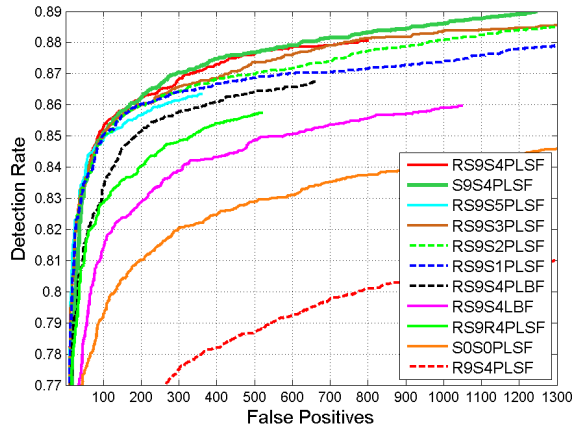
Fig. 9. Evaluation of design decisions on FDDB.



Fig. 10. Failure example from the FDDB dataset where many small, blurry and occluded faces have not been detected.

regression, support pruning, and a parameter sensitive classifier. The result "S0S0PLSF" has candidates predicted only from the detected face centers, with $N^{supp} = 0$. It shows that predicting poses from multiple keypoints and support pruning has a considerable increase in detection accuracy. The result "R9S4PLSF" detects keypoints with the RF and without the FSA verification classifiers from Figure 6, clearly showing the importance of the keypoint verification step. The result "RS9S4LBF" uses a parameter insensitive classifier trained with the logistic loss, and performs inferior to the parameter sensitive classifier "RS9S4PLBF", which in turn is inferior to the parameter sensitive classifier "RS9S4PLSF" based on LSF+special features. The result "RS9R4PLSF" uses the Random Forest described in Section 5.1 for 3D pose regression and the parameter sensitive classifier for verification. Again, it performs inferior to the algorithm "RS9S4PLSF" with FSA-based pose regression. Also shown are results with the proposed method with $N^{supp} = 1, 2, 3, 5$.

**Failure modes.** We observed that the proposed 3D model-based face detection approach has difficulties mainly with detecting small, blurry and occluded faces, as shown in Figure 10. Indeed, a 3D model based on face keypoints doesn't make much sense for small faces since the keypoints might be too small to be visible.
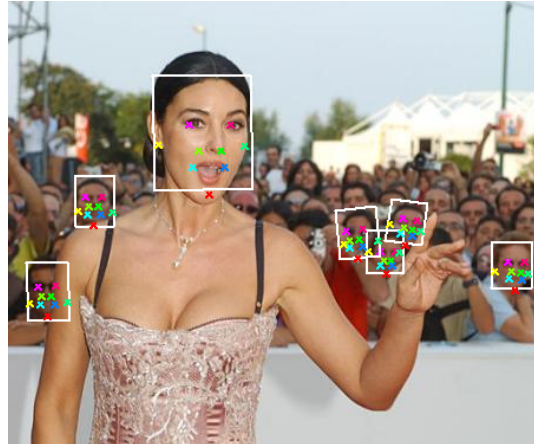
**Detection time.** The detection time for a 480x320 image is about 3 seconds when using the RF screening and 15 seconds without the RF screening. The C++ code has not been optimized for speed and most of the time is used for detecting the keypoints. We expect to obtain speedups of 10-100 times with a GPU implementation and code optimization.

## 6 CONCLUSION

In this paper we presented a method for face detection that uses a 3D model to represent the face hypotheses, extract 3D pose-aligned features and to specify the yaw parameter value for a parameter sensitive classifier. The 3D face candidates are proposed by image based regression starting from a number of face keypoints that are detected first. From experiments we observed that generating 3D face candidates from multiple face keypoints and pruning them based on the number of keypoints detected results in considerable improvements in detection accuracy compared to generating candidates only from the face center. The results obtained are comparable with the state of the art in the low false positive regime and outperforms the cascade-based state of the art methods for the regime of at least 0.04 false positives per image.

# REFERENCES

[1] A. Barbu, Y. She, L. Ding, and G. Gramajo. Feature selection with annealing for big data learning. *arXiv preprint arXiv:1310.2880*, 2013.

[2] A. Bartoli, D. Pizarro, and M. Loog. Stratified generalized procrustes analysis. *IJCV*, 101(2):227–253, 2013.

[3] X. P. Burgos-Artizzu, P. Perona, and P. Dollár. Robust face landmark estimation under occlusion. ICCV, 2013.

[4] X. Cao, Y. Wei, F. Wen, and J. Sun. Face alignment by explicit shape regression. In *CVPR*, pages 2887–2894, 2012.

[5] D. Chen, S. Ren, Y. Wei, X. Cao, and J. Sun. Joint cascade face detection and alignment. In *ECCV*, pages 109–122. 2014.

[6] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *CVPR*, volume 1, pages 886–893, 2005.

[7] D. F. Dementhon and L. S. Davis. Model-based object pose in 25 lines of code. *IJCV*, 15(1-2):123–141, 1995.

[8] P. Dollár, R. Appel, S. Belongie, and P. Perona. Fast feature pyramids for object detection. *IEEE Trans. PAMI*, 36(8):1532–1545, 2014.

[9] P. Dollár, P. Welinder, and P. Perona. Cascaded pose regression. In *CVPR*, pages 1078–1085, 2010.

[10] P. Felzenszwalb, R. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *PAMI*, 32(9):1627–1645, 2010.

[11] M. Hejrati and D. Ramanan. Analyzing 3d objects in cluttered images. In *NIPS*, pages 602–610, 2012.

[12] W. Hu. Learning 3d object templates by hierarchical quantization of geometry and appearance spaces. In *CVPR*, pages 2336–2343, 2012.

[13] V. Jain and E. Learned-Miller. Fddb: A benchmark for face detection in unconstrained settings. Technical Report UM-CS-2010-009, University of Massachusetts, Amherst, 2010.

[14] M. Koestinger, P. Wohlhart, P. Roth, and H. Bischof. Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization. In *IEEE Workshop on Benchmarking Facial Image Analysis Technologies*, 2011.

[15] H. Li, Z. Lin, J. Brandt, X. Shen, and G. Hua. Efficient boosted exemplar-based face detection. In *CVPR*, pages 1843–1850, 2014.

[16] H. Li, Z. Lin, X. Shen, J. Brandt, and G. Hua. A convolutional neural network cascade for face detection. In *CVPR*, pages 5325–5334, 2015.

[17] J. Liebelt and C. Schmid. Multi-view object class detection with a 3d geometric model. In *CVPR*, pages 1688–1695, 2010.

[18] M. Mathias, R. Benenson, M. Pedersoli, and L. Van Gool. Face detection without bells and whistles. In *ECCV*, pages 720–735. 2014.

[19] M. Osadchy, Y. L. Cun, and M. L. Miller. Synergistic face detection and pose estimation with energy-based models. *JMLR*, 8:1197–1215, 2007.

[20] N. Payet and S. Todorovic. From contours to 3d object detection and pose estimation. In *ICCV*, pages 983–990, 2011.

[21] 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.

[22] P. Schönemann and R. Carroll. Fitting one matrix to another under choice of a central dilation and a rigid motion. *Psychometrika*, 35(2):245–255, 1970.

[23] X. Shen, Z. Lin, J. Brandt, and Y. Wu. Detecting and aligning faces by image retrieval. In *CVPR*, pages 3460–3467, 2013.

[24] H. Su, M. Sun, L. Fei-Fei, and S. Savarese. Learning a dense multi-view representation for detection, viewpoint classification and synthesis of object categories. In *ICCV*, pages 213–220, 2009.

[25] P. Viola and M. J. Jones. Robust real-time face detection. *IJCV*, 57(2):137–154, 2004.

[26] L. Wan, D. Eigen, and R. Fergus. End-to-end integration of a convolution network, deformable parts model and non-maximum suppression. In *CVPR*, pages 851–859, 2015.

[27] T. Wu and S.-C. Zhu. A numerical study of the bottom-up and top-down inference processes in and-or graphs. *International journal of computer vision*, 93(2):226–252, 2011.

[28] J. Yan, Z. Lei, L. Wen, and S. Z. Li. The fastest deformable part model for object detection. In *CVPR*, pages 2497–2504, 2014.

[29] J. Yan, X. Zhang, Z. Lei, and S. Z. Li. Face detection by structural models. *Image and Vision Computing*, 32(10):790–799, 2014.

[30] B. Yang, J. Yan, Z. Lei, and S. Z. Li. Aggregate channel features for multi-view face detection. In *Biometrics (IJCB), 2014 IEEE International Joint Conference on*, pages 1–8, 2014.

[31] Q. Yuan, A. Thangali, V. Ablavsky, and S. Sclaroff. Parameter sensitive detectors. In *CVPR*, pages 1–6, 2007.

[32] Q. Yuan, A. Thangali, V. Ablavsky, and S. Sclaroff. Learning a family of detectors via multiplicative kernels. *IEEE Trans. PAMI*, 33(3):514–530, 2011.

[33] C. Zhang and Z. Zhang. Winner-take-all multiple category boosting for multi-view face detection. 2009.

[34] Y. Zhao and S.-C. Zhu. Image parsing with stochastic scene grammar. In *Advances in Neural Information Processing Systems*, pages 73–81, 2011.

[35] S. Zhou and D. Comaniciu. Shape regression machine. In *IPMI*, pages 13–25, 2007.

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

Fig. 11. Detected faces on the FDDB dataset.