



# Cascaded Nonlinear Shape Model Regression

---

Pedro Martins<sup>1</sup>, Bruno Silva<sup>1</sup>, Jorge Batista<sup>2,1</sup>

<https://www.isr.uc.pt/~pedromartins>  
[pedromartins@isr.uc.pt](mailto:pedromartins@isr.uc.pt)

<sup>1</sup>Institute of Systems and Robotics (ISR)  
University of Coimbra, Portugal

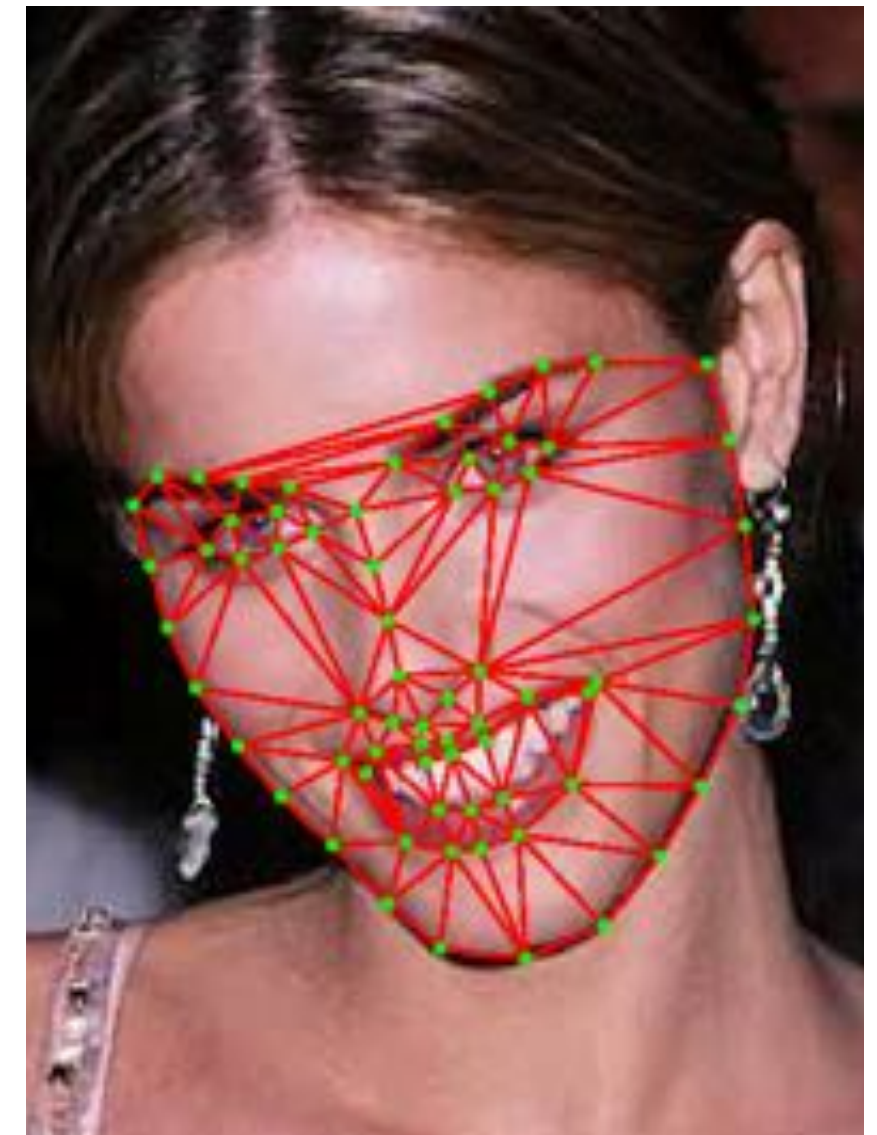
<sup>2</sup>Department of Electrical and Computer Engineering (DEEC)  
University of Coimbra, Portugal



# Introduction

---

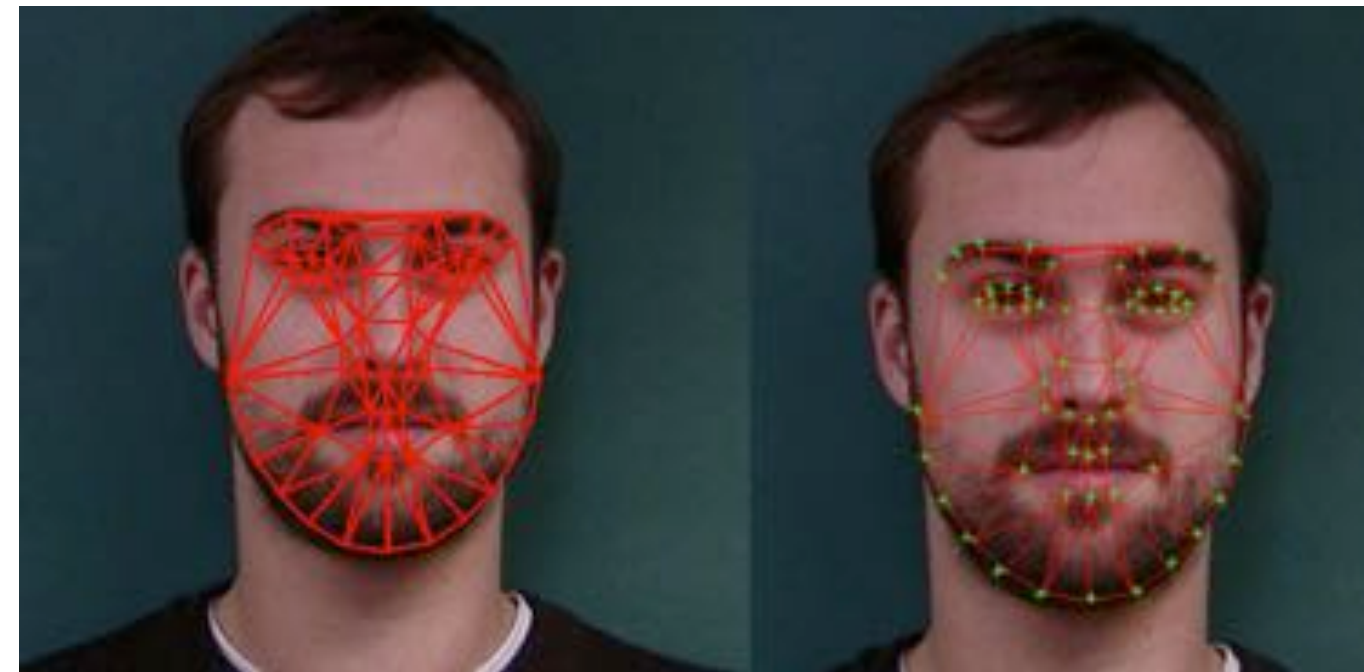
- Nonrigid face alignment (facial landmark localization) with deformable models.
- Range of Applications: Face recognition, emotion recognition, pose estimation, ...
- Cascaded Regression Framework.
  - Gradient Descent vs Cascaded Regression.
  - Review of the standard formulation.
- Nonlinear Extension of the Cascaded Regression
  - Linear Shape Model.
  - Nonlinear Regression (via Convolution Neural Network).
- Evaluation Results (LFPW, LFW, HELEN, 300W datasets).



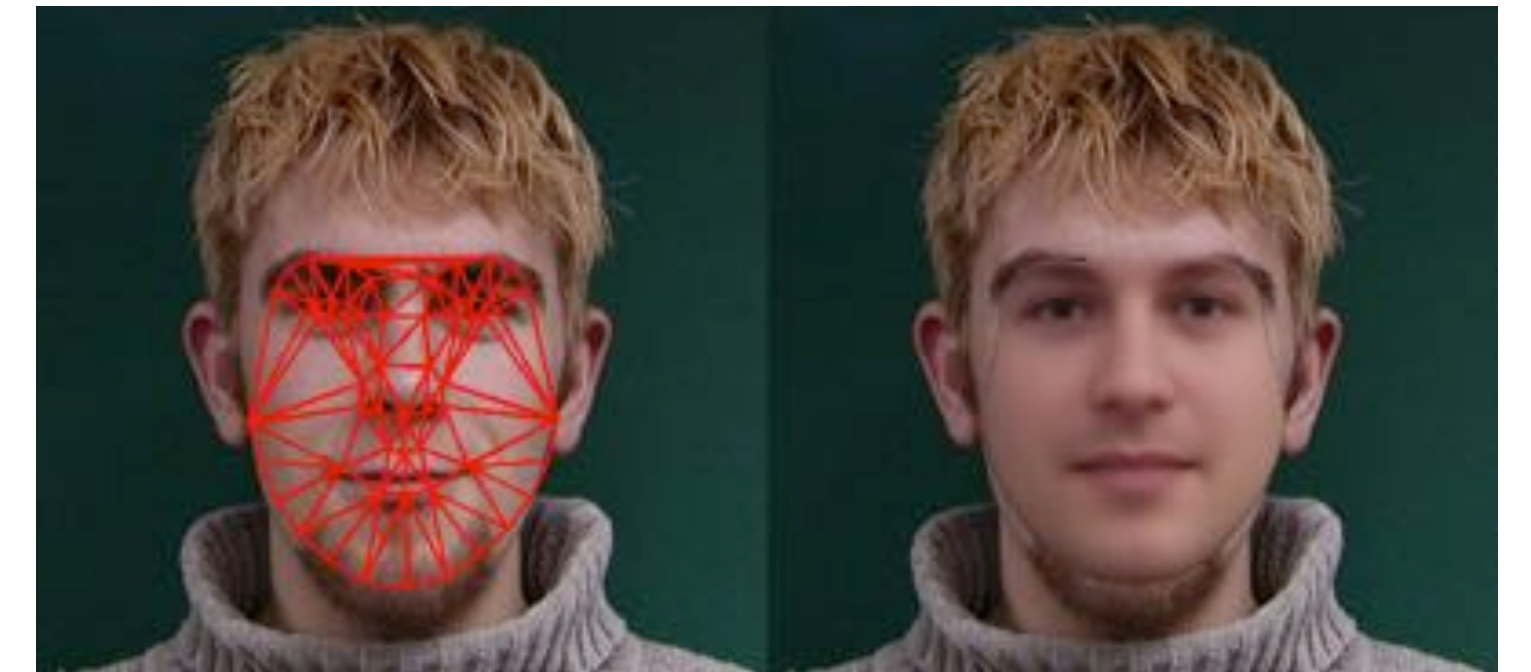
**Face alignment example**

# Related Work

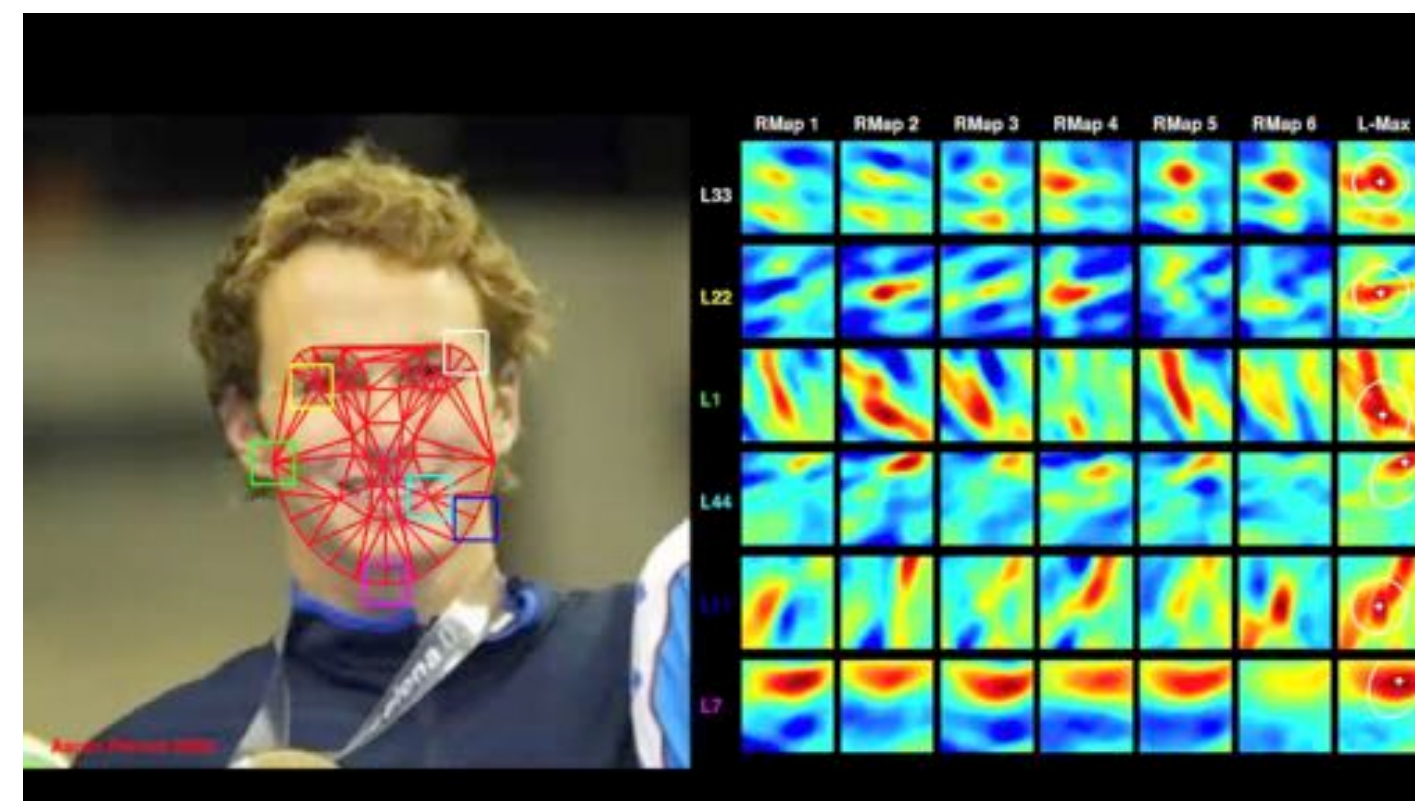
- Active Shape Model (ASM)
- Deformable Part Model (DPM)
- Active Appearance Model (AAM)
  - Project-Out Inverse Compositional (PO-IC)
  - Simultaneous Inverse Compositional (SIC)
- Constrained Local Model (CLM)
  - Convex Quadratic Fitting (CQF)
  - Subspace Constrained Mean-Shifts (SCMS)
  - Bayesian CLM (BCLM)
- Cascaded Regression (CR)
  - Supervised Descent Method (SDM)
  - Project-Out Cascade Regression (PO-CR)
  - Simultaneous Cascaded Regression (SCR)



**Active Shape Model (ASM)**



**Active Appearance Model (AAM)**

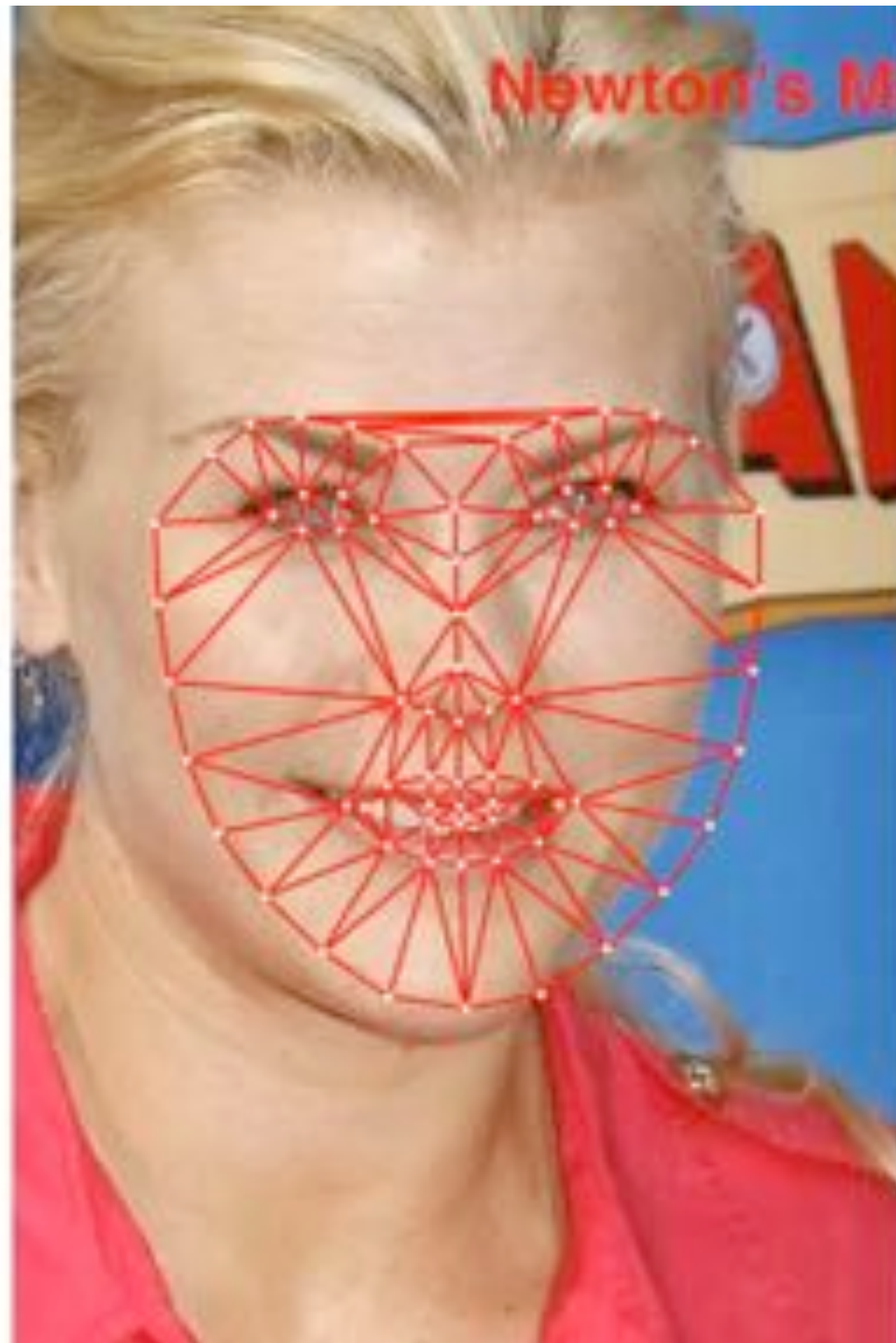


**Bayesian Constrained Local Model (BCLM)**



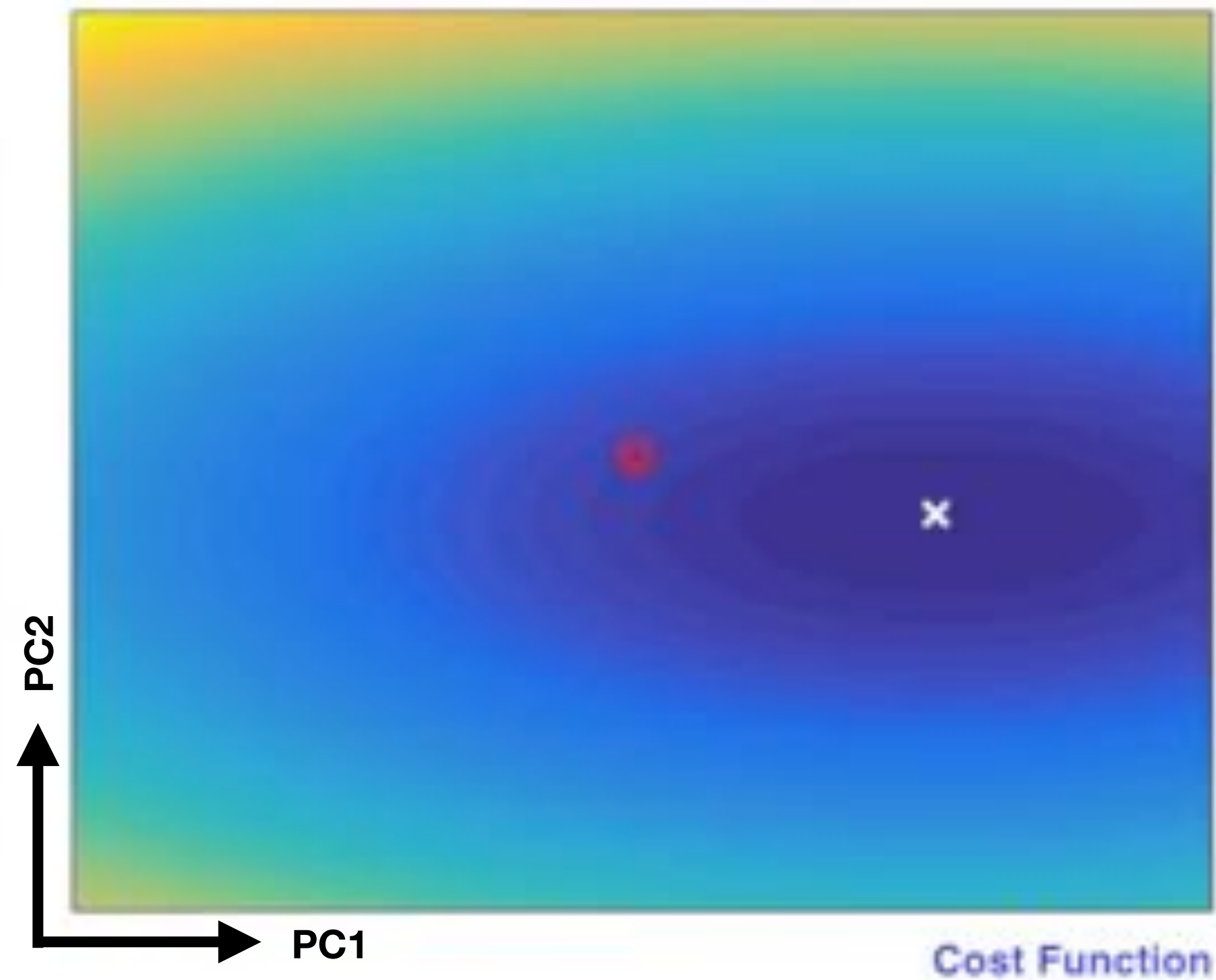
**Simultaneous Cascaded Regression (SCR)**

# Gradient Descent vs. Cascaded Regression



## Gradient Descent

- Requires 'good' initialization.
- In general, requires to compute the Jacobian at each iteration.
- Require to compute the Hessian and its inverse (2nd order methods).
- Learning Fast.
- Testing Slow.



## Cascaded Regression

- Captures the variance of the initialization.
- Precomputed Regression matrix.
- Learning Slow.
- Testing Fast.

# Cascade Regression Framework



$k$  - cascade level

$$\mathbf{s}^k = \mathbf{s}^{k-1} + \mathbf{R}^{k-1} \mathcal{F}(\mathbf{I}, \mathbf{s}^{k-1})$$

Updated shape vector

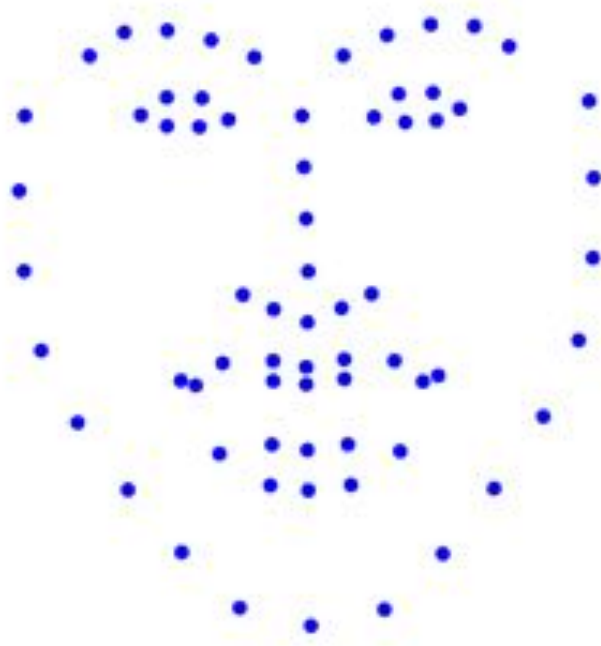
Previous shape vector

Regression Matrix

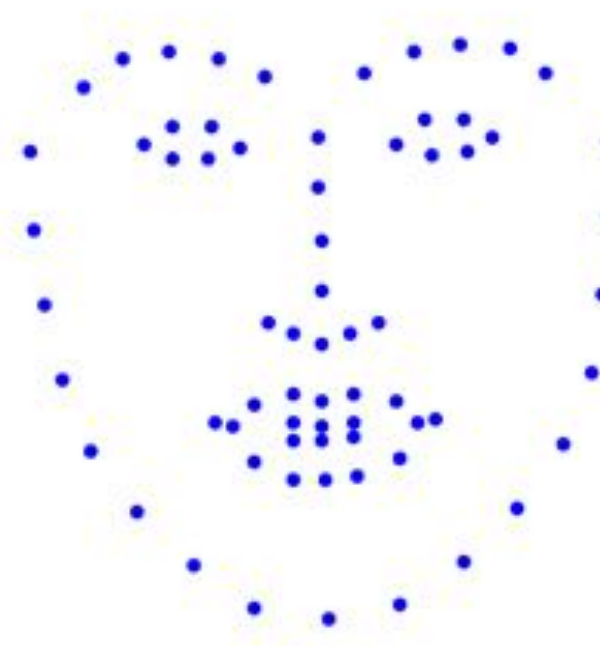
Feature Extraction

$$\mathbf{s} = \begin{pmatrix} x_0 \\ \vdots \\ x_v \\ y_0 \\ \vdots \\ y_v \end{pmatrix}$$

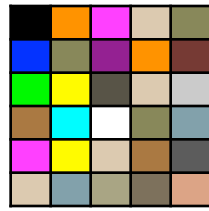
$v$  - landmarks



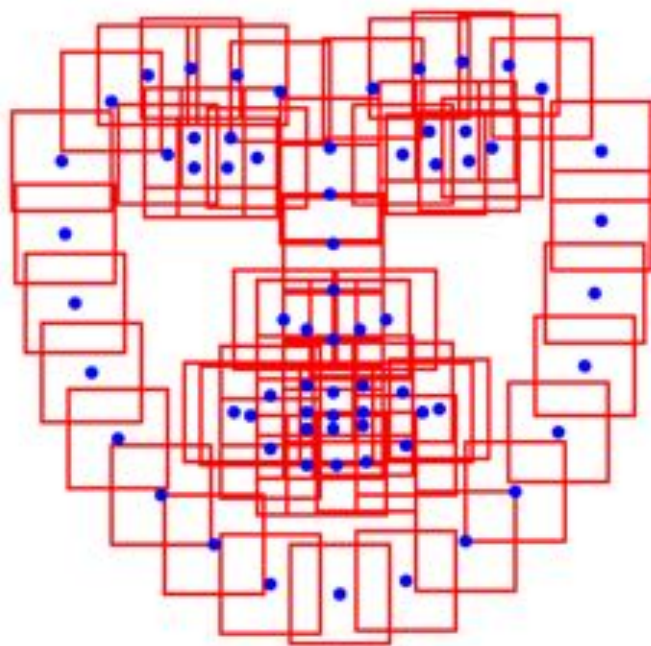
$2v \times 1$



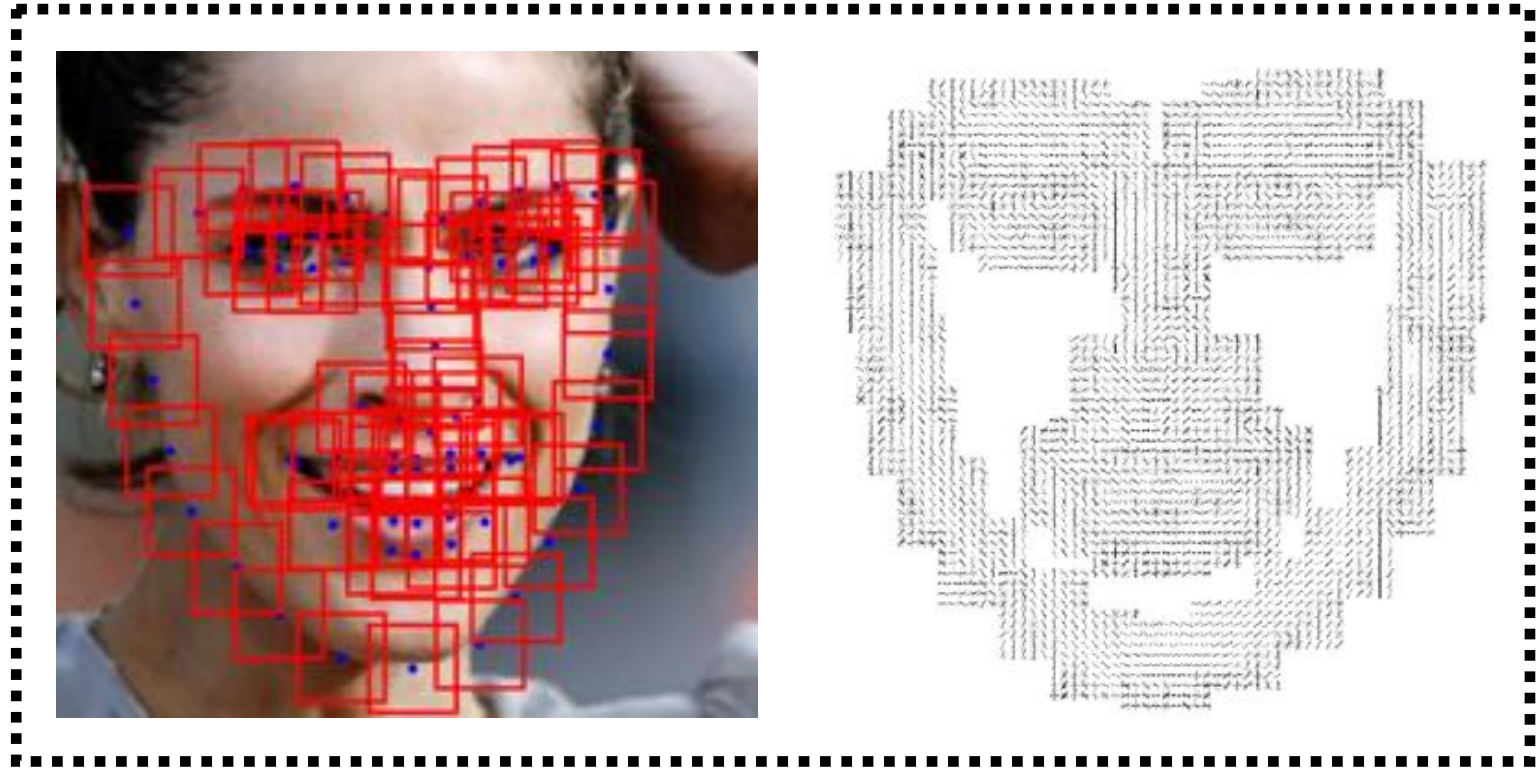
$2v \times 1$



$2v \times d$



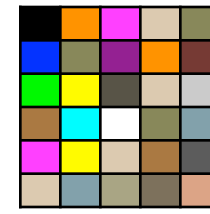
$d \times 1$



RGB

HoG

# Learning Regression Matrix



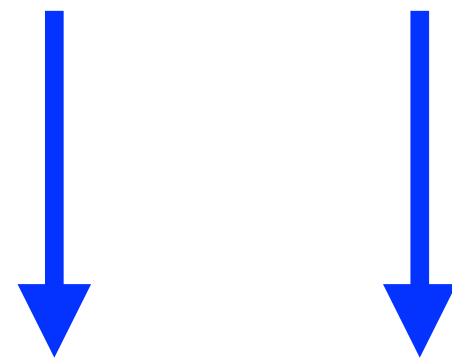
Estimate  $\mathbf{R}^k$  under Multiple Initializations

$$\arg \min_{\mathbf{R}^k} \sum_{i=1}^N \sum_{j=1}^M \|\Delta \mathbf{s}_j^k - \mathbf{R}^k \mathcal{F}(\mathbf{I}_i, \mathbf{s}_j^k)\|^2$$

$k$  - cascade level  
 $i$  - training image  
 $j$  - virtual sample

Estimate noise

$$\Sigma^k = \text{COV}(\mathbf{s}_* - \mathbf{s}_j^k)$$



Deviation from Ground Truth

Regression Labels

$$\Delta \mathbf{s}_j^k = \mathbf{s}_* - \mathbf{s}_j^k$$

Data Matrix (all features)

$$\mathbf{F} = \left[ \begin{array}{c|c|c|c|c} \color{red}\blacksquare & \color{yellow}\blacksquare & \color{green}\blacksquare & \color{cyan}\blacksquare & \color{magenta}\blacksquare \\ \color{blue}\blacksquare & \color{orange}\blacksquare & \color{purple}\blacksquare & \color{red}\blacksquare & \color{brown}\blacksquare \\ \color{cyan}\blacksquare & \color{green}\blacksquare & \color{magenta}\blacksquare & \color{blue}\blacksquare & \color{orange}\blacksquare \\ \color{magenta}\blacksquare & \color{blue}\blacksquare & \color{orange}\blacksquare & \color{red}\blacksquare & \color{brown}\blacksquare \\ \color{brown}\blacksquare & \color{orange}\blacksquare & \color{red}\blacksquare & \color{brown}\blacksquare & \color{orange}\blacksquare \end{array} \right]$$

← N images x M virtual samples →

Least Squares Solution

$$\mathbf{R}^k = \Delta \mathbf{S} \left( \mathbf{F}^T \mathbf{F} \right)^{-1} \mathbf{F}^T$$

Linear Regression



Data Collection (F matrix)

$$\mathbf{s}_j^k \sim \mathcal{N}(\mu^k, \Sigma^k)$$



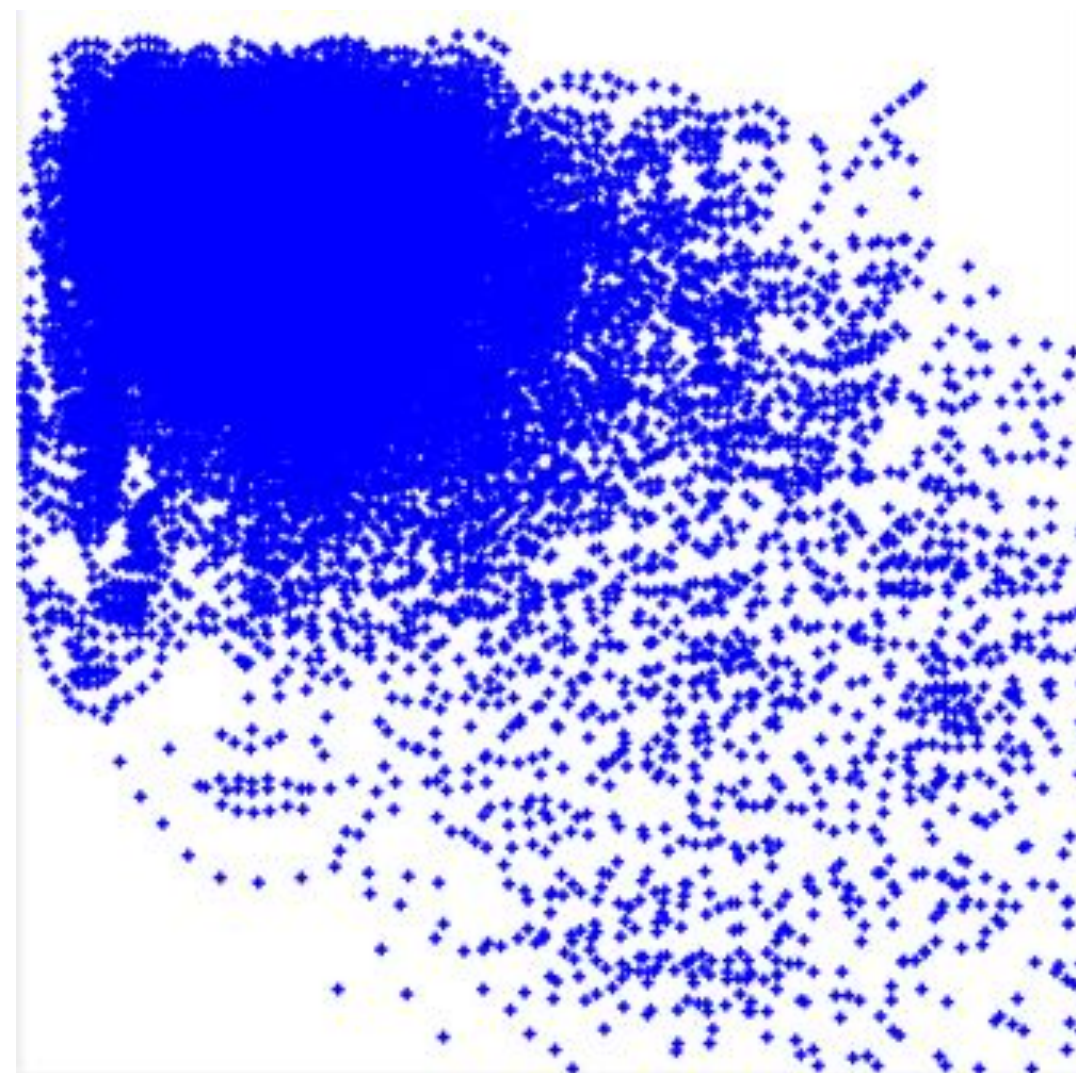
virtual sample

# Linear Shape Model

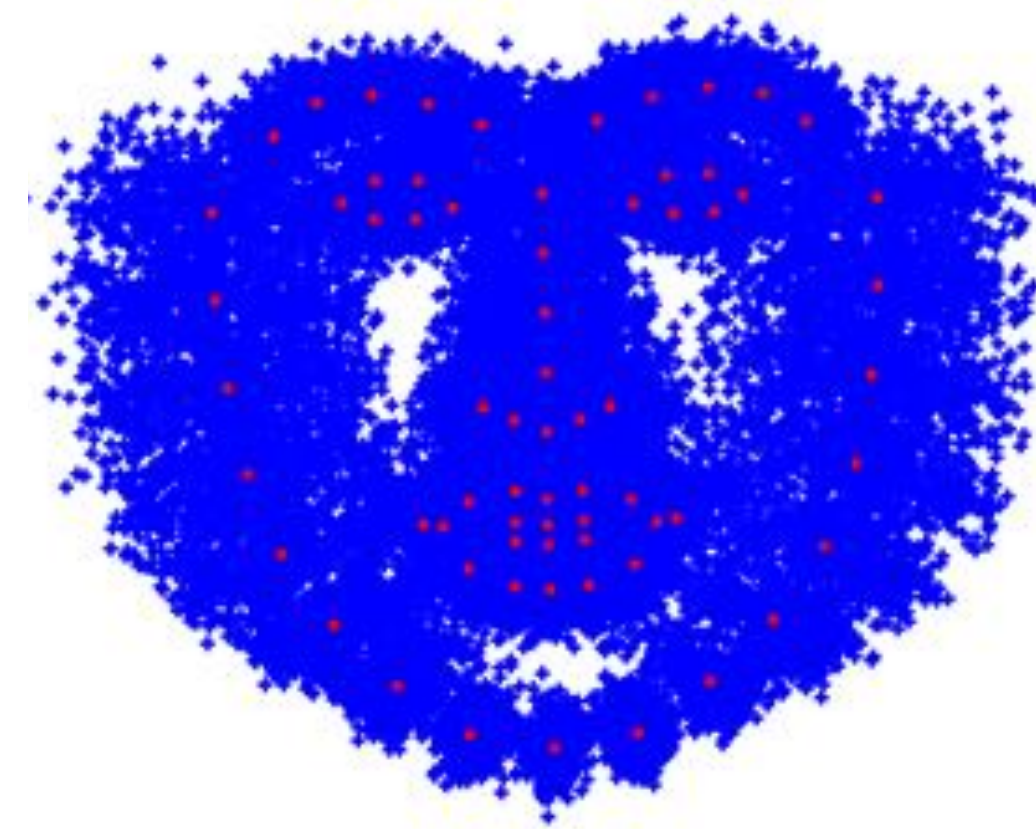
'In the Wild' Image Database



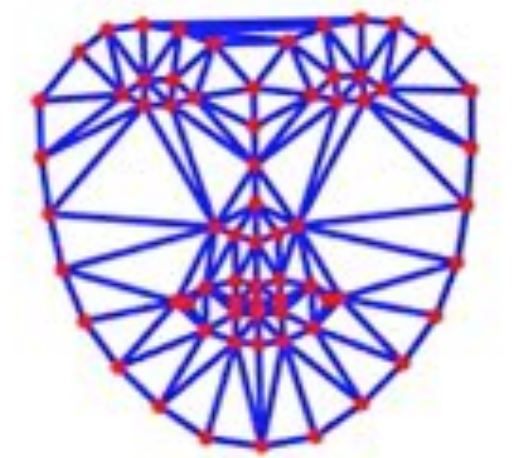
RAW Shape Data



Procrustes Alignment



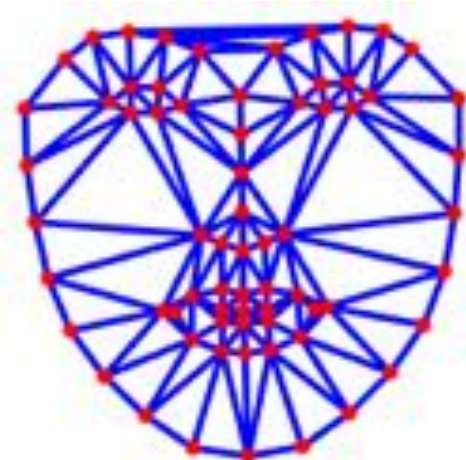
Shape Model



$$\mathcal{B}(\mathbf{s}; \mathbf{b}) = \mathbf{s}_0 + \sum_{i=1}^n \phi_i b_i$$

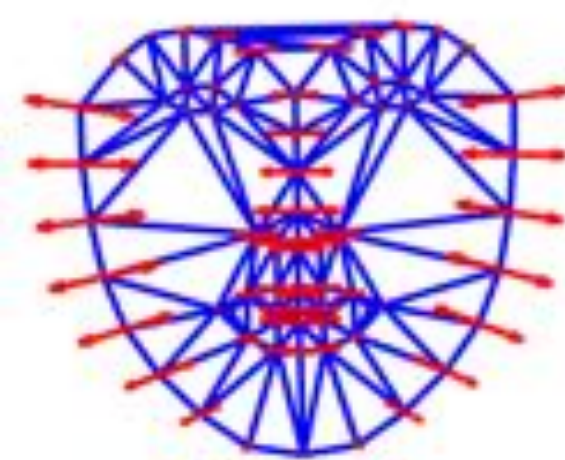
↑  
shape parameters

Mean Shape

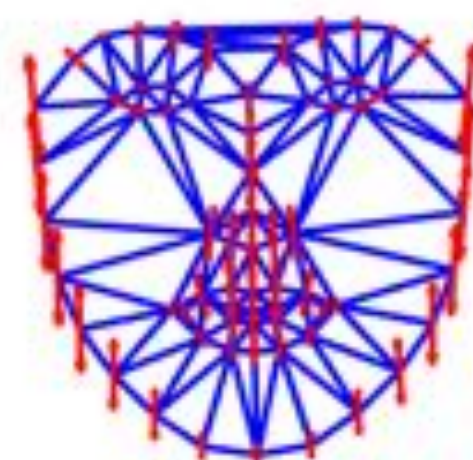


$\mathbf{s}_0$

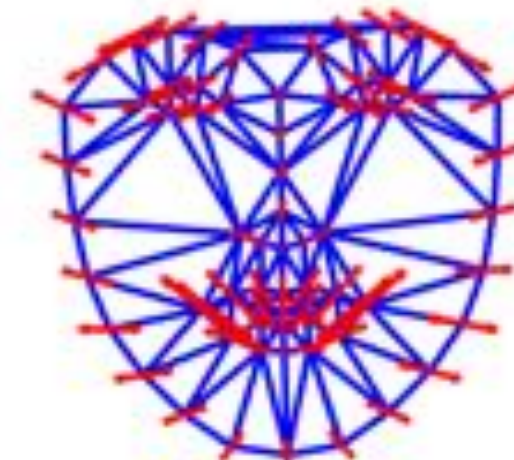
Shape Basis



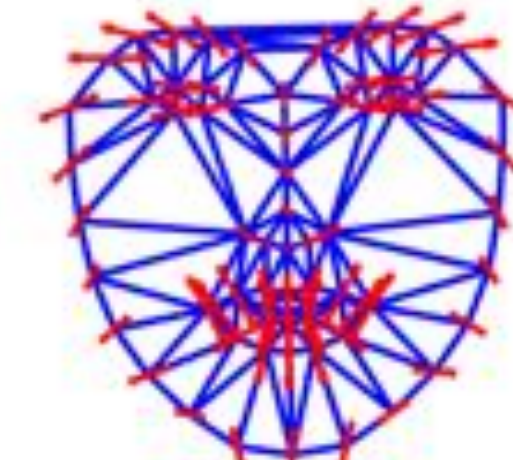
$\phi_1$



$\phi_2$



$\phi_3$



$\phi_4$

Similarity Transform

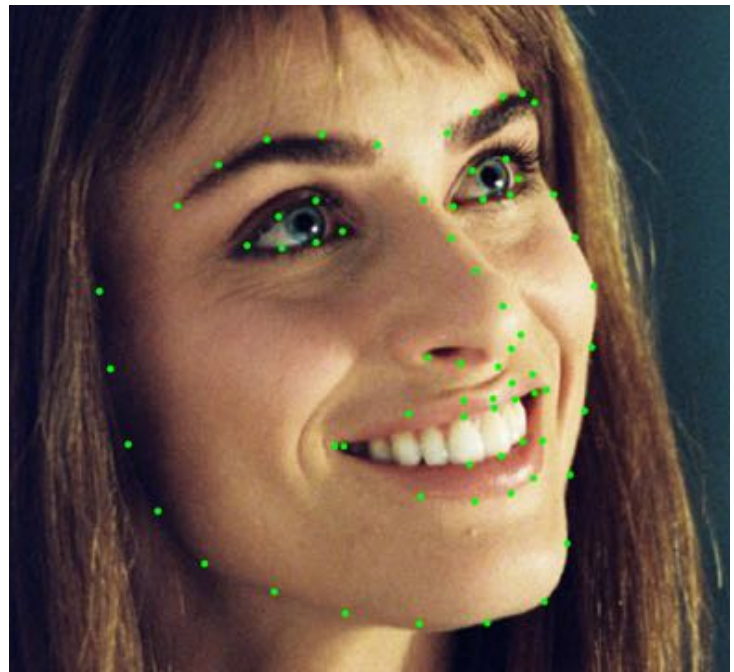
$$\mathcal{S}(\mathbf{s}; \mathbf{q}) = \mathbf{s} + \sum_{j=1}^4 \psi_j q_j$$

↑  
pose parameters

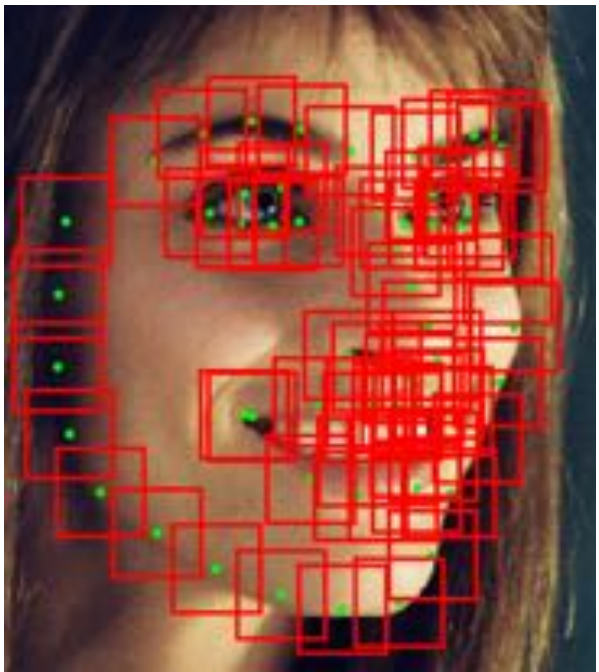
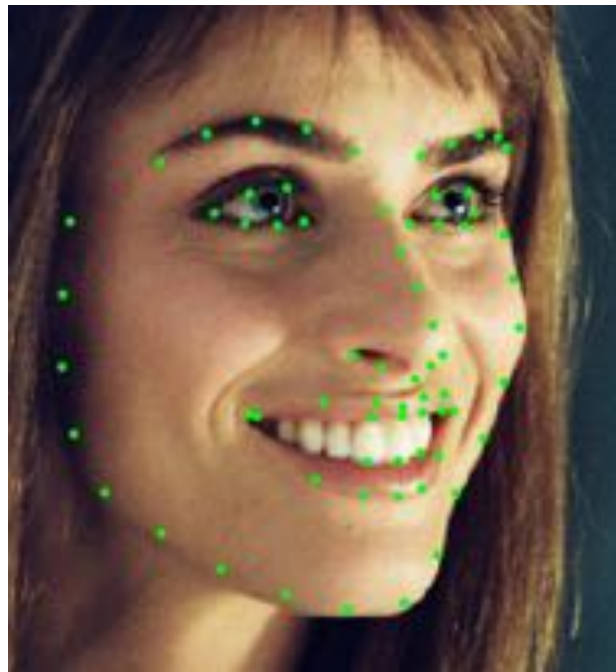
Full Shape Model

$$\mathbf{s} = \mathcal{S}(\mathcal{B}(\mathbf{b}); \mathbf{q}) \quad 7$$

# Local Appearance Regions



Similarity Warp  
( $s, \theta, t_x, t_y$ )



Local Appearance Regions

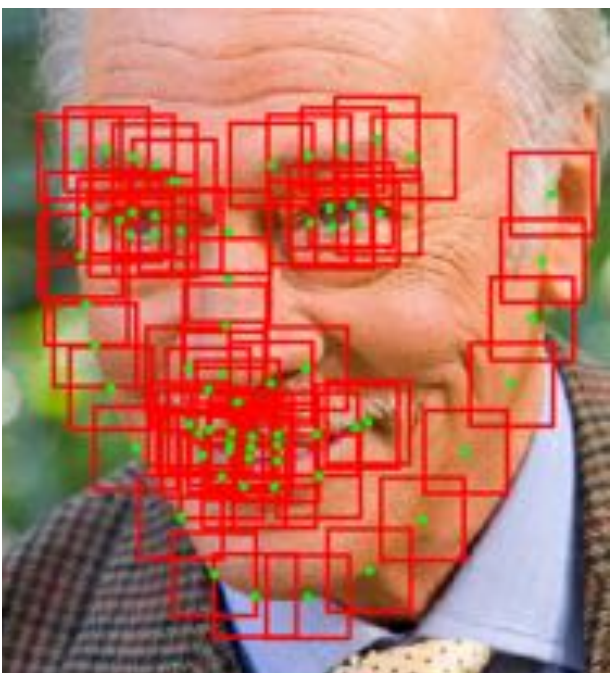
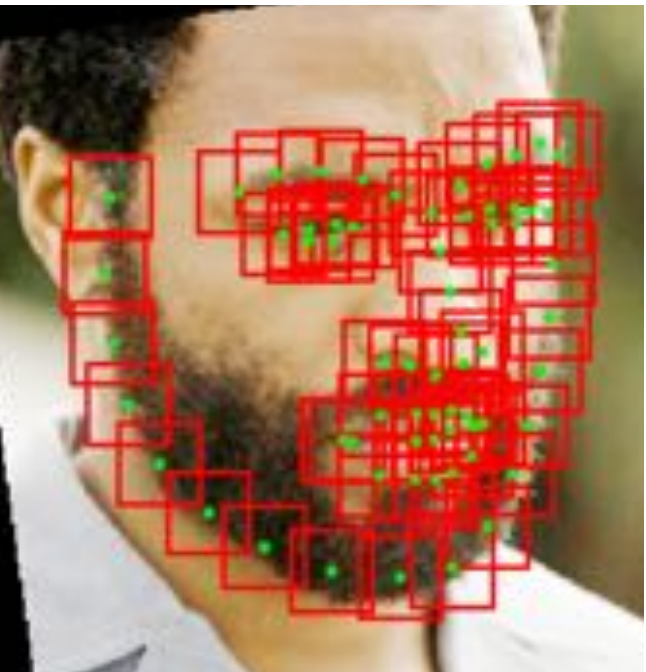


Image + Landmarks

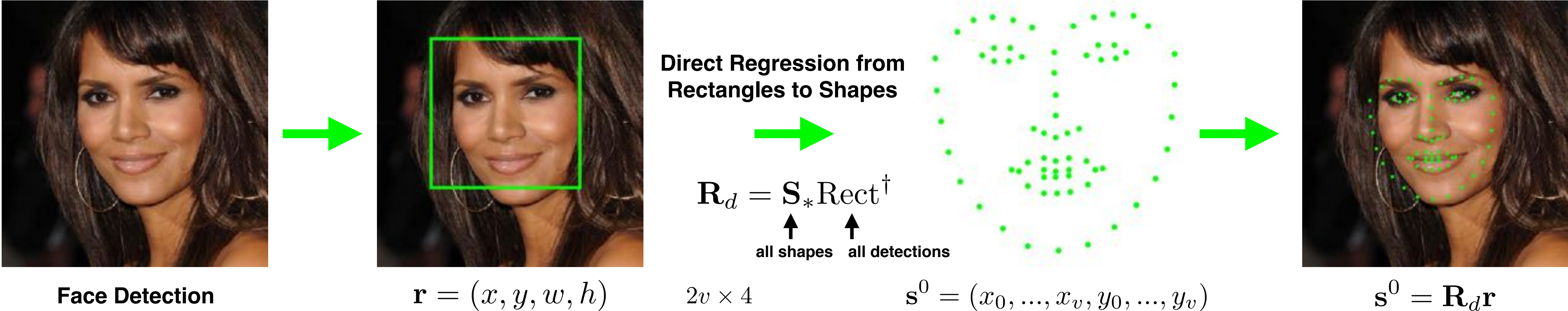
Normalized Image

Local Patches

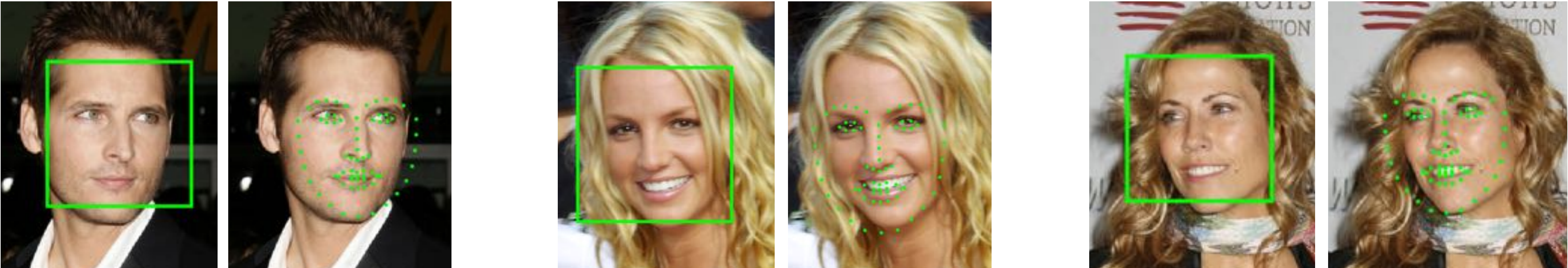
Sampled Local Patches



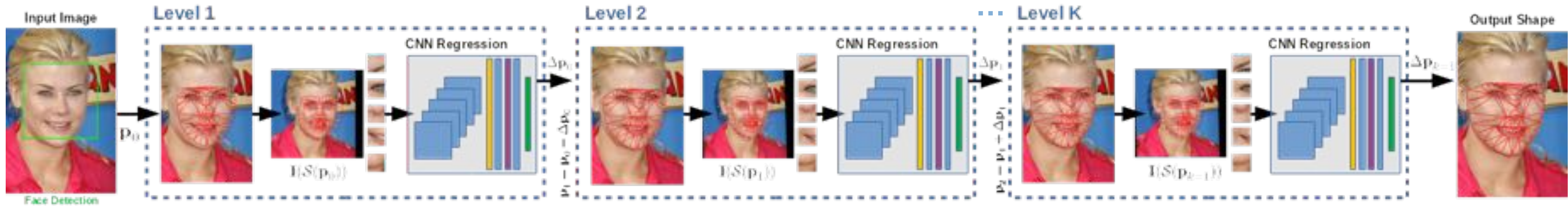
# Initial Shape Estimate ( $s^0$ )



Other Examples:



# Nonlinear Cascade Regression



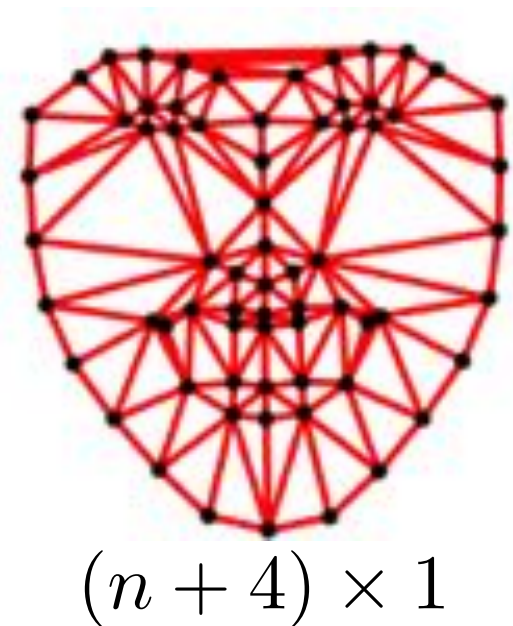
Combined shape + pose parameters

$$\mathbf{p} = \begin{bmatrix} \mathbf{b} \\ \mathbf{q} \end{bmatrix} \in \mathbb{R}^{n+4}$$

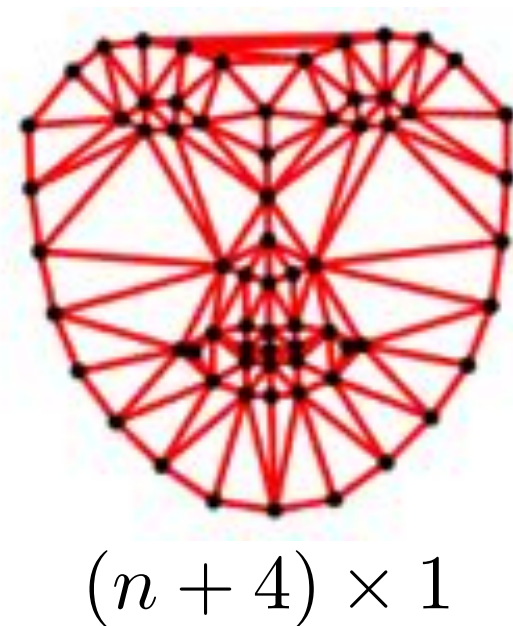
$$\mathbf{p}^k = \mathbf{p}^{k-1} + \gamma \mathcal{R}^{k-1} \{ \mathcal{L}(\mathbf{I}(\mathcal{S}(\mathbf{p}^{k-1}))) \}$$

$k$  - cascade level

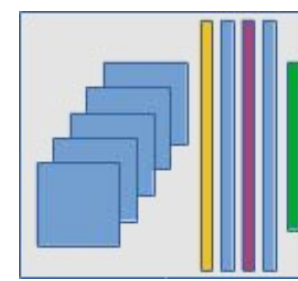
Updated shape instance



Previous shape instance



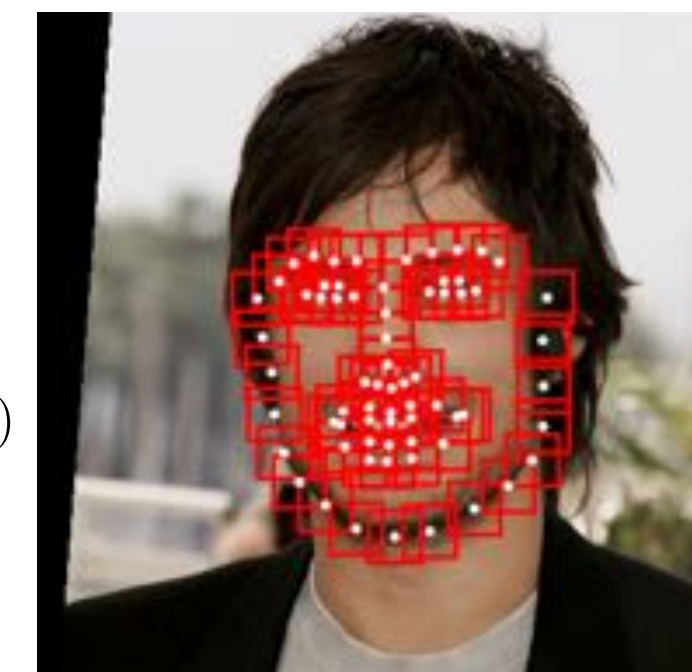
Nonlinear Mapping



Local Feature Extraction at Normalized Frame



Similarity Warp  
 $\mathbf{p}(n + 1 : n + 4)$



Sampled 3D Array

$$P \times P \times v$$



$$\mathcal{L}(\mathbf{I}(\mathcal{S}(\mathbf{p})))$$

# CNN Regression Architecture

## Nonlinear Regression

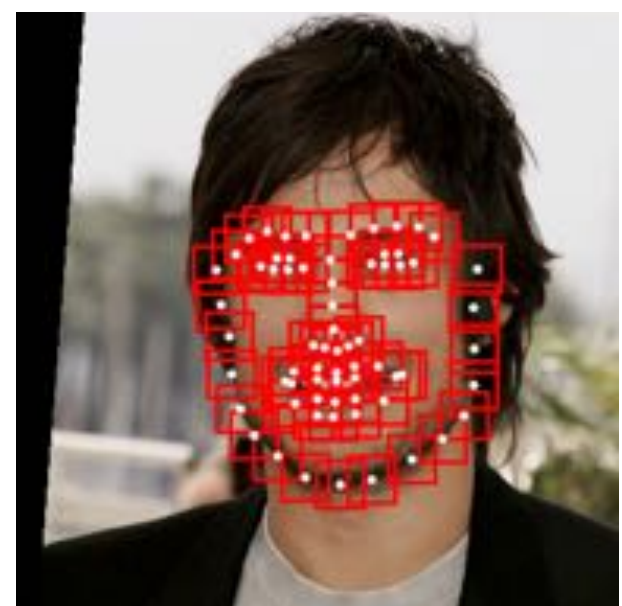
$$\arg \min_{\mathcal{R}^k} \sum_{i=1}^N \sum_{j=1}^M \|\Delta \mathbf{p}_{ij}^k - r_L(\dots r_1(\mathcal{L}(\mathbf{I}_i(\mathcal{S}(\mathbf{p}_j^k)))))\|_{\Sigma^k}^2$$

↓ CNN Topology

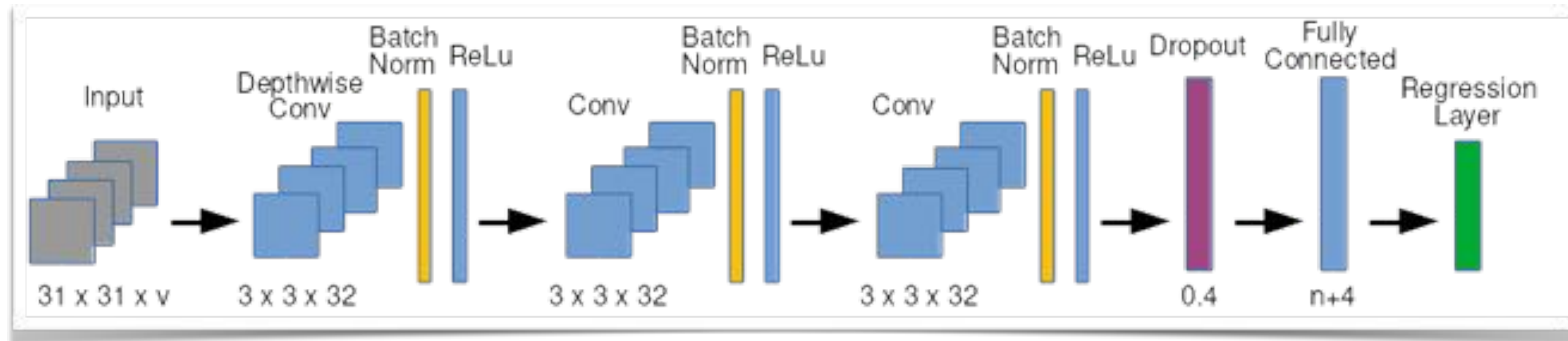
Input



$\mathcal{L}(\mathbf{I}(\mathcal{S}(\mathbf{p})))$



Pose Normalized Image



**Depthwise Convolution**  
32 filters (3x3) for each local patch

**Convolution**  
32 filters (3x3)

**Convolution**  
32 filters (3x3)

**Dropout**  
0.4

**Regression Layer**  
(n+4)

Output

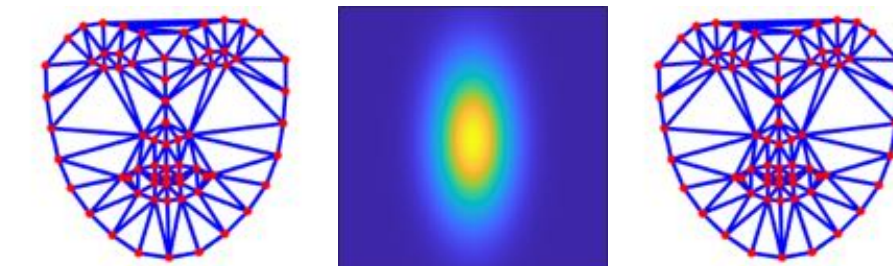


shape parameters update  
 $\Delta \mathbf{p}$

## Loss function

$$L_r = \frac{1}{N} \sum_{j=1}^N \Delta \mathbf{p}_j^T \Sigma_{\mathbf{p}}^{-1} \Delta \mathbf{p}_j$$

Mahalanobis Distance



# CNN Learning - Data Collection

$$\arg \min_{\mathcal{R}^k} \sum_{i=1}^N \sum_{j=1}^M \|\Delta \mathbf{p}_{ij}^k - \text{CNN}^k(\mathcal{L}(\mathbf{I}_i(\mathcal{S}(\mathbf{p}_j^k))))\|_{\Sigma^k}^2$$

$k$  - cascade level  
 $i$  - training image  
 $j$  - virtual sample

Estimate noise

$$\Sigma^k = \text{cov}(\mathbf{p}_* - \mathbf{p}_{ij})$$

Deviation from Ground Truth

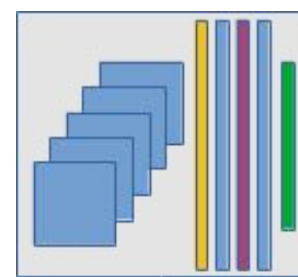
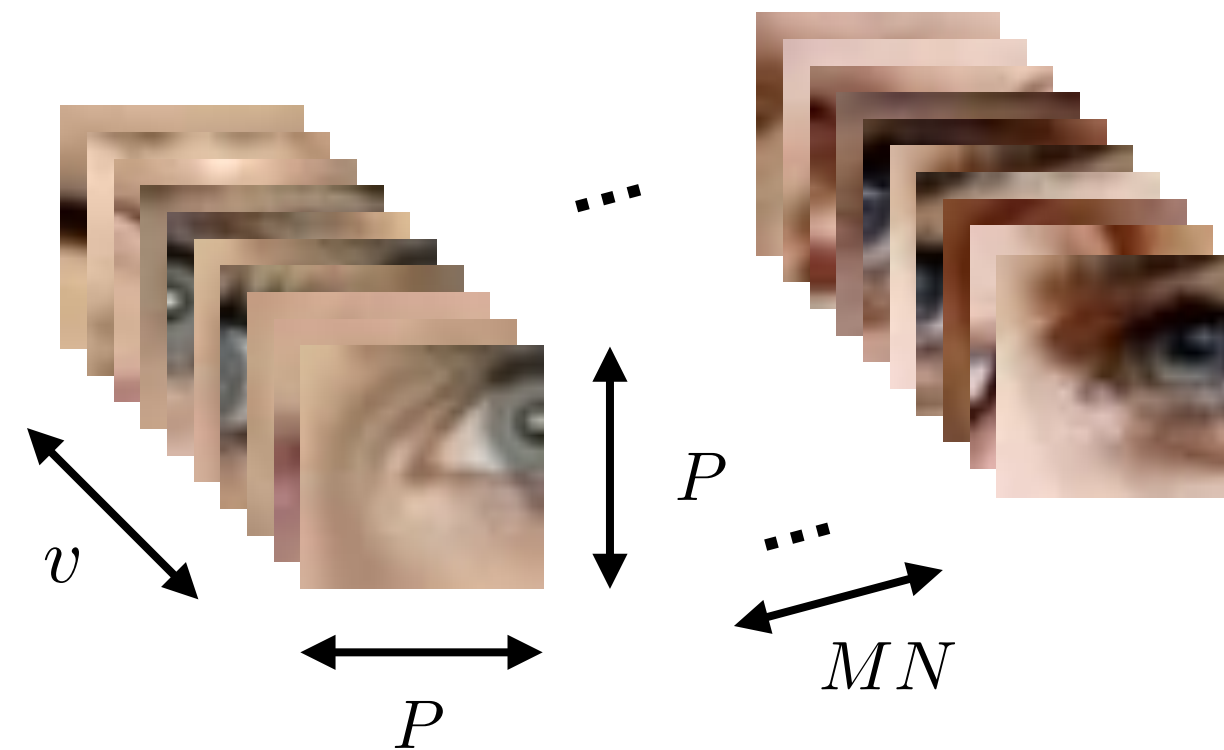
$$\Delta \mathbf{p}_{ij} = \mathbf{p}_* - \mathbf{p}_{ij}$$

Data Matrix (local normalized patches)

$$\mathbf{D}_{ij} = \mathcal{L}(\mathbf{I}_i(\mathcal{S}(\mathbf{p}_j)))$$

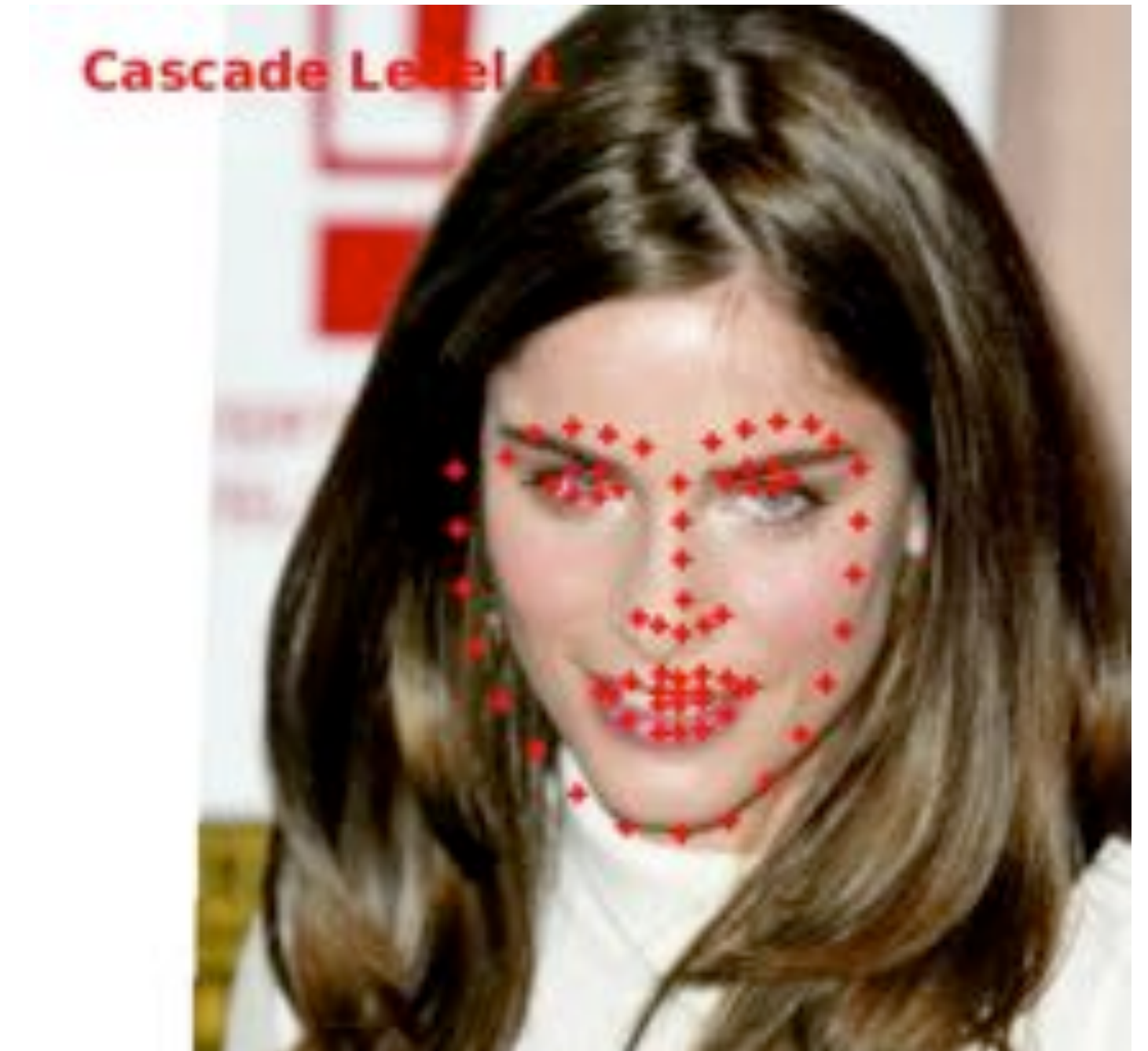
Full Data: 4D Array  
( $P \times P \times v \times N.M$ )

Regression Labels



CNN<sup>k</sup>

M - augmented examples  
N - real examples

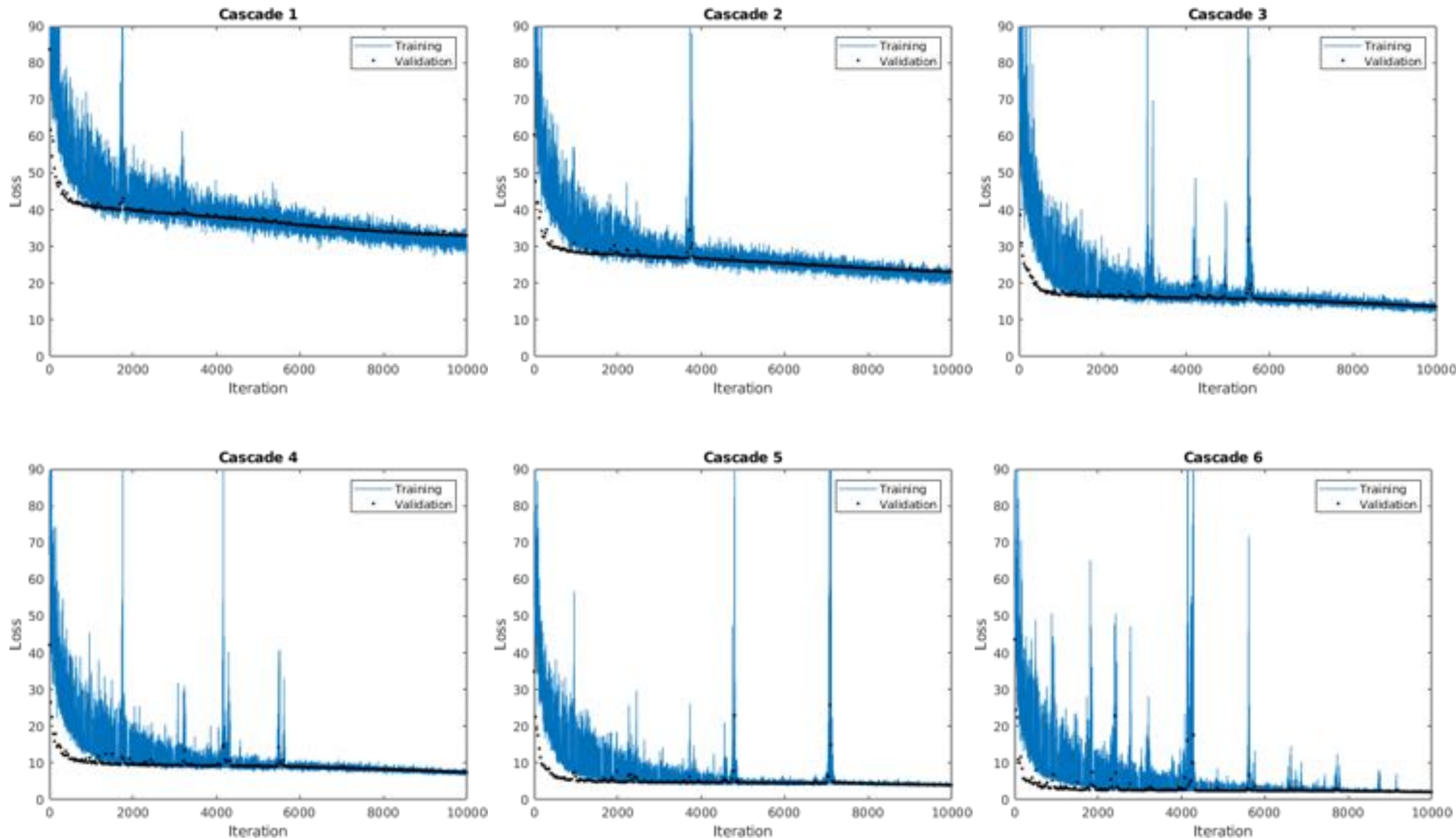


Data Collection (D matrix)

$$\mathbf{p}_{ij} \sim \mathcal{N}(\mathbf{p}_i, \Sigma^k)$$

↑  
virtual shape sample

# CNN Learning - Optimization Details



| Name   | Type                | Activations | Learnables                            |
|--|---------------------|-------------|---------------------------------------|
| imageinput<br>31x31x68 images  | Image Input         | 31×31×68    | -                                     |
| grouped_conv<br>68 groups of 32 3x3x1 convolutions with stride [1 1] and padding [0 0 0] | Grouped Convolution | 29×29×2176  | Weights 3×3×1×32×68<br>Bias 1×1×32×68 |
| batchnorm_1<br>Batch normalization with 2176 channels                                    | Batch Normalization | 29×29×2176  | Offset 1×1×2176<br>Scale 1×1×2176     |
| relu_1<br>ReLU   | ReLU                | 29×29×2176  | -                                     |
| conv_1<br>32 3x3x2176 convolutions with stride [1 1] and padding [0 0 0]                 | Convolution         | 27×27×32    | Weights 3×3×2176×32<br>Bias 1×1×32    |
| batchnorm_2<br>Batch normalization with 32 channels                                      | Batch Normalization | 27×27×32    | Offset 1×1×32<br>Scale 1×1×32         |
| relu_2<br>ReLU   | ReLU                | 27×27×32    | -                                     |
| conv_2<br>32 3x3x32 convolutions with stride [1 1] and padding [0 0 0]                   | Convolution         | 25×25×32    | Weights 3×3×32×32<br>Bias 1×1×32      |
| batchnorm_3<br>Batch normalization with 32 channels                                      | Batch Normalization | 25×25×32    | Offset 1×1×32<br>Scale 1×1×32         |
| relu_3<br>ReLU   | ReLU                | 25×25×32    | -                                     |
| dropout<br>40% dropout   | Dropout             | 25×25×32    | -                                     |
| fc<br>31 fully connected layer   | Fully Connected     | 1×1×31      | Weights 31×28068<br>Bias 31×1         |
| matl<br>Mandelbrot sq loss v1  | Regression Output   | -           | -                                     |

CNN topology details

## CNN Optimization:

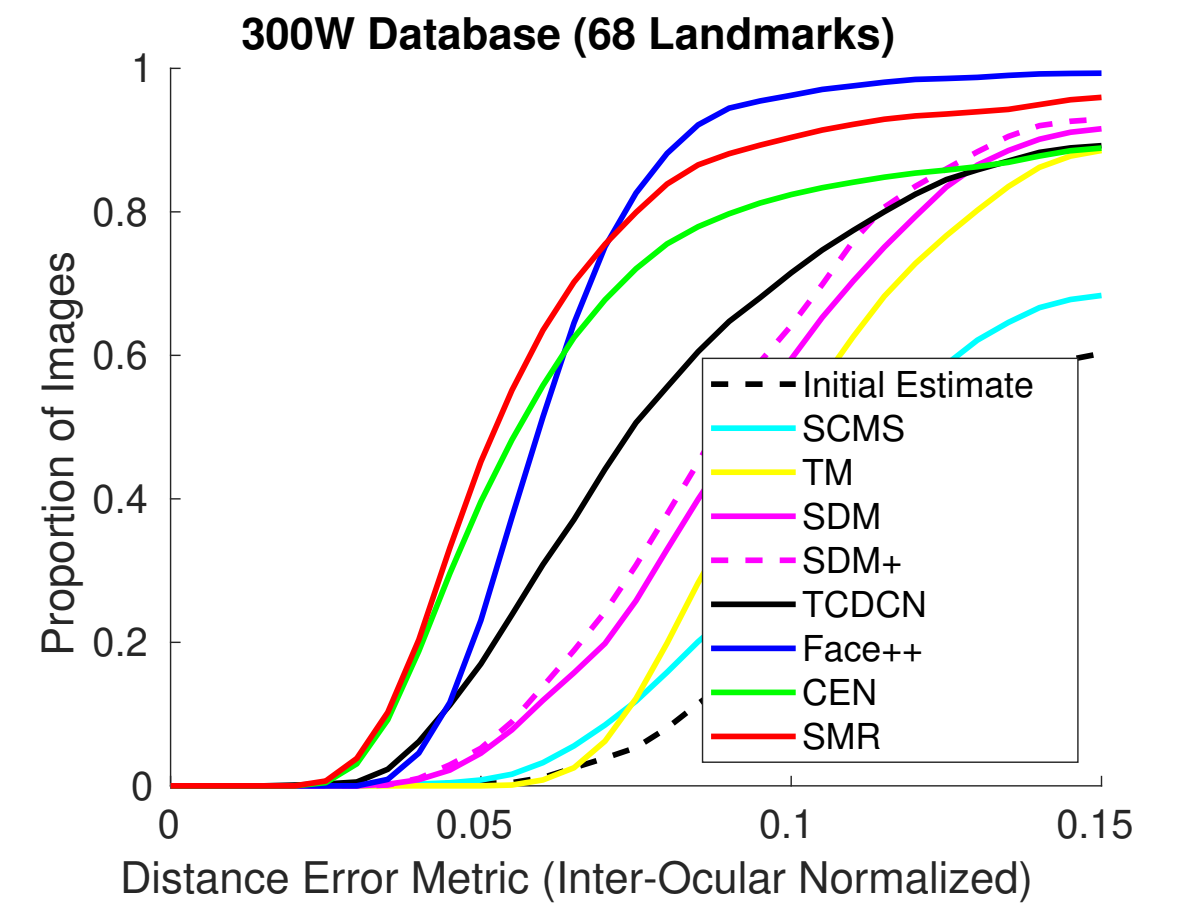
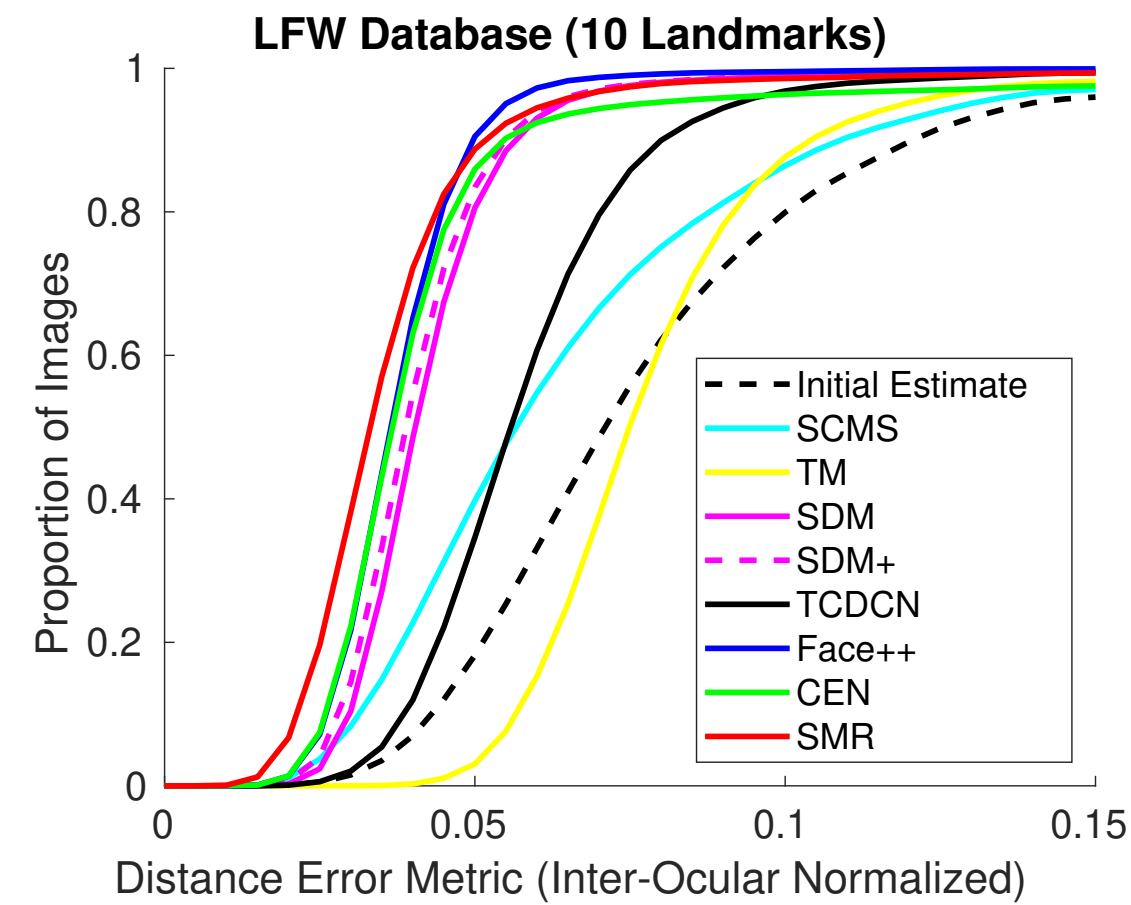
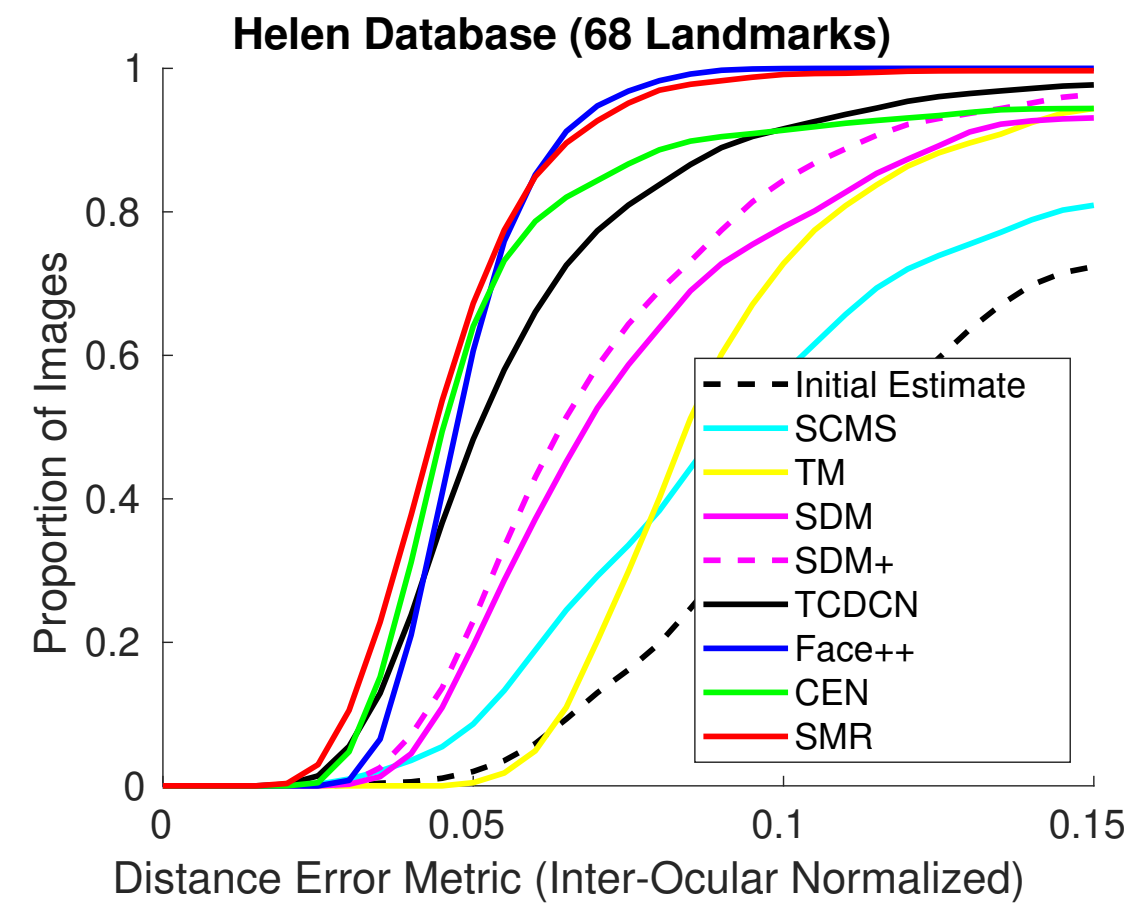
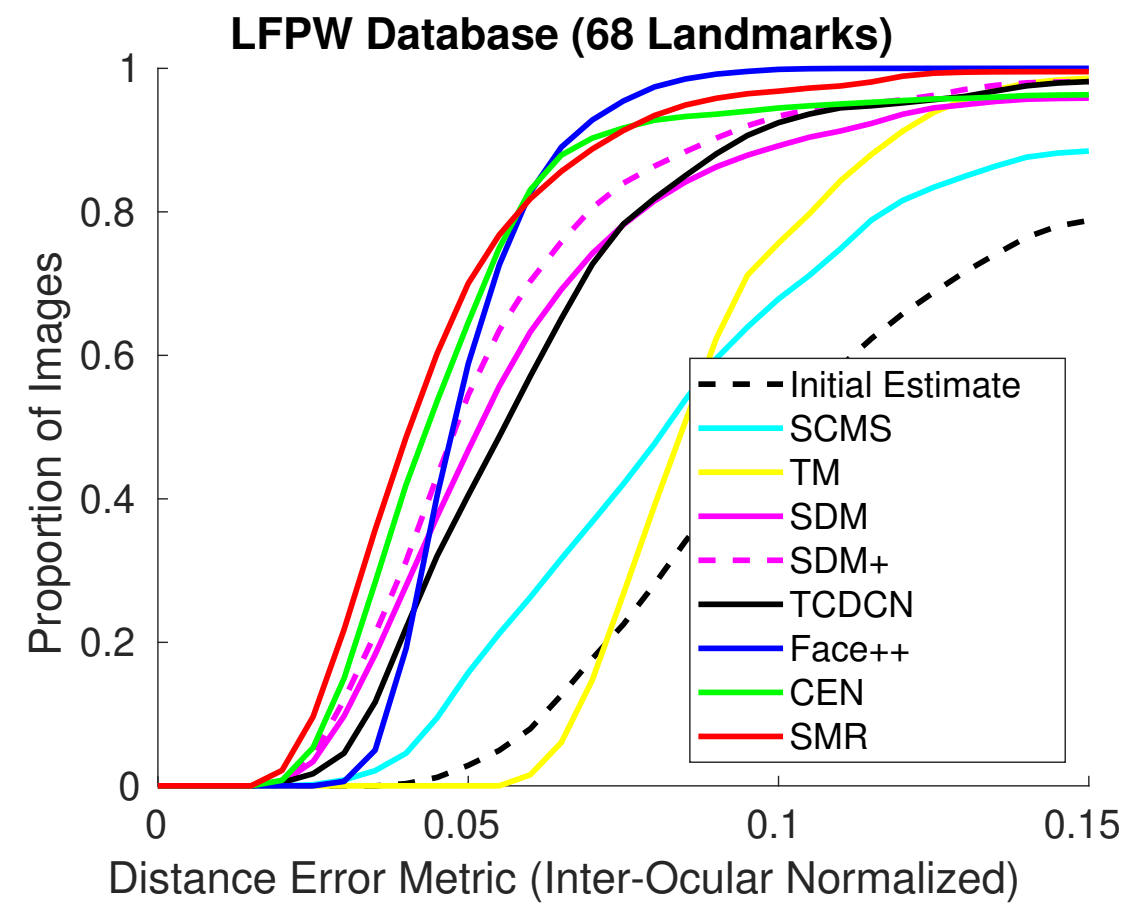
- Adam solver (default hyperparameters)
- Initial learning rate  $10^{-4}$ , exp. decay 0.9 every 5000 iterations
- Mini-batch w/ 64 examples
- Max cascade levels  $K = 6$

**Demo:**

**Nonlinear Cascaded SMR  
Fitting in the LFPW Database**

# Evaluation Results

## Cumulative error distribution function (CDF)

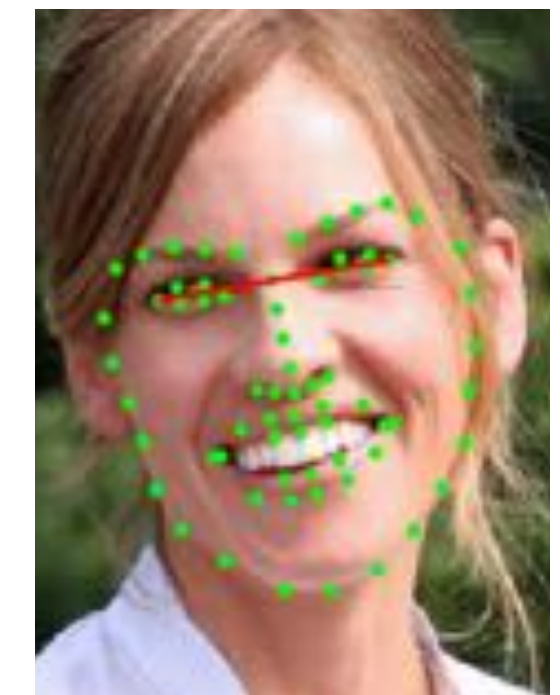


| Method / AUC      | LFPW        | HELEN       | LFW         | 300W        |
|-------------------|-------------|-------------|-------------|-------------|
| Initial Estimate  | 30.5        | 25.4        | 49          | 17,6        |
| <b>SCMS</b>       | <b>42.4</b> | <b>36.0</b> | <b>57.3</b> | <b>23.2</b> |
| <b>Tree-Model</b> | <b>40.8</b> | <b>39.7</b> | <b>47.7</b> | <b>30.6</b> |
| <b>SDM</b>        | <b>60.2</b> | <b>48.6</b> | <b>71.5</b> | <b>36.5</b> |
| <b>SDM+</b>       | <b>63.6</b> | <b>52.1</b> | <b>72.4</b> | <b>39.0</b> |
| <b>TCDCN</b>      | <b>59.5</b> | <b>61.1</b> | <b>61.0</b> | <b>44.6</b> |
| <b>Face++</b>     | <b>66.7</b> | <b>67.4</b> | <b>74.8</b> | <b>58.4</b> |
| <b>CEN</b>        | <b>67.2</b> | <b>63.6</b> | <b>72.3</b> | <b>54.0</b> |
| <b>SMR</b>        | <b>69.6</b> | <b>69.1</b> | <b>75.9</b> | <b>59.5</b> |

This work →

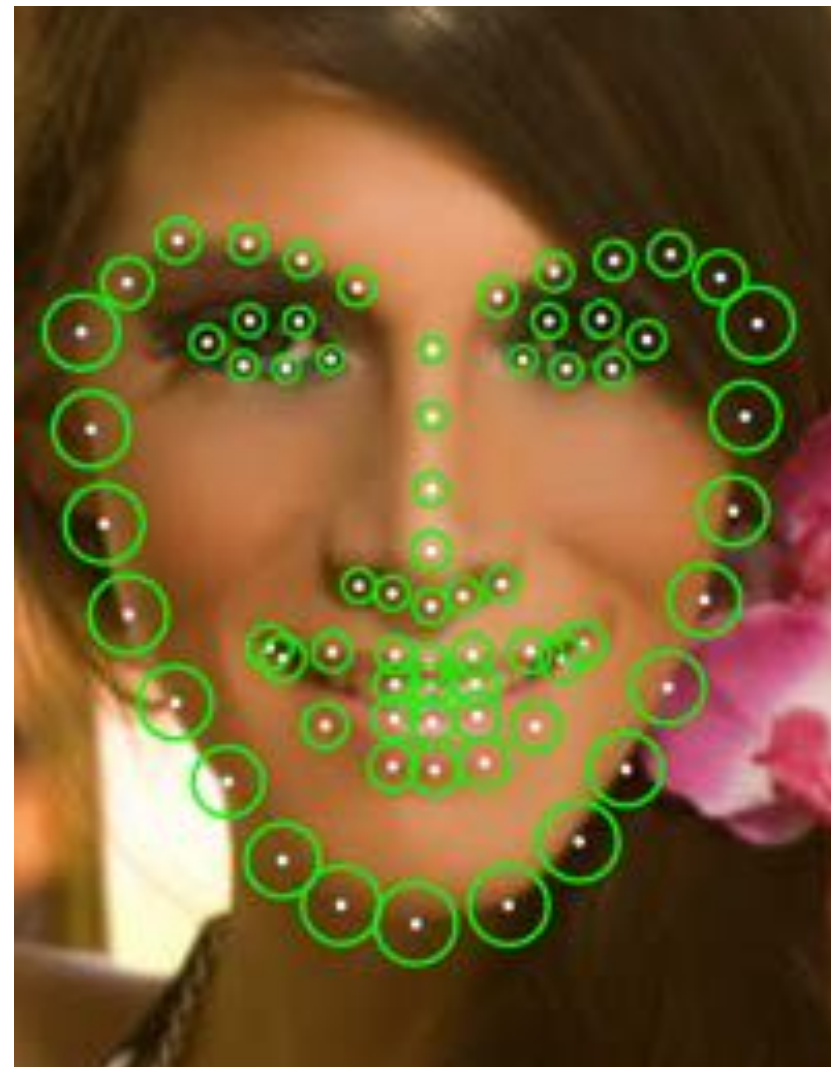
AUC - Area Under Curve

## Inter-ocular normalized error

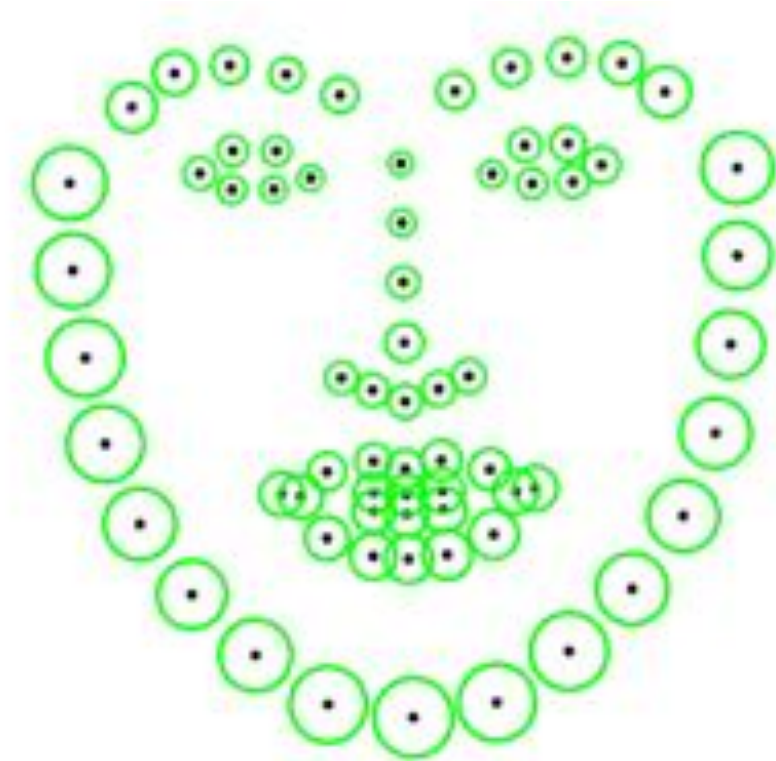


$$e_m(\mathbf{s}) = \frac{1}{v d_{\text{eyes}}} \sum_{i=1}^v \|\mathbf{s}_i - \mathbf{s}_i^*\| \quad 15$$

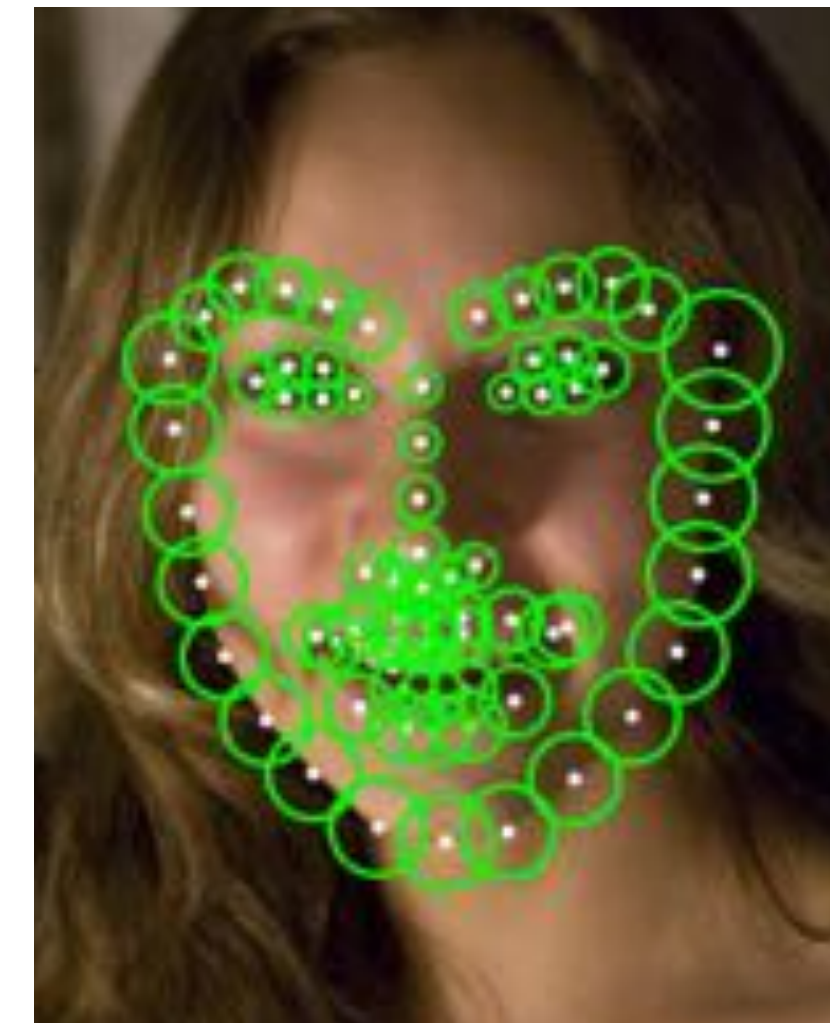
# Evaluation Results - Fitting Error Standard Deviation



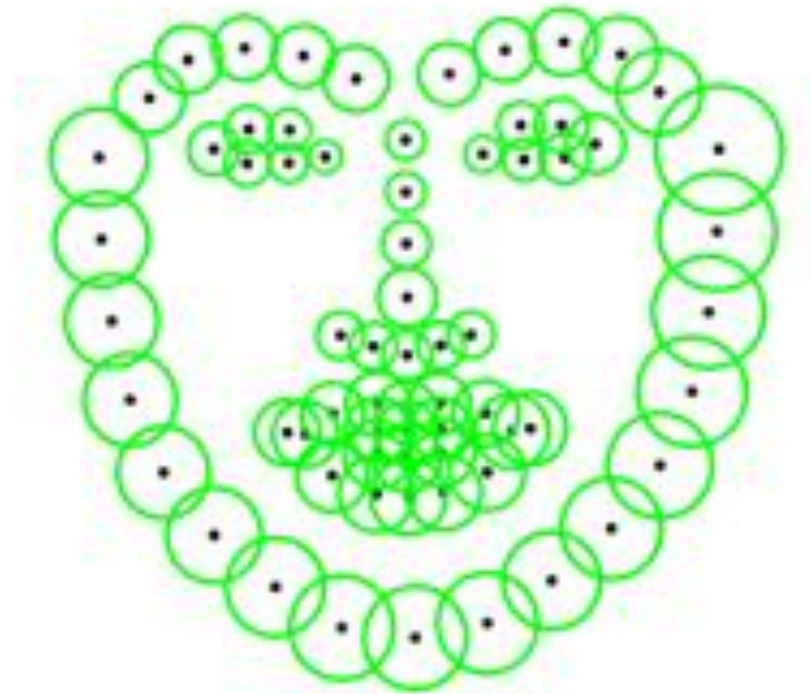
LFPW Database



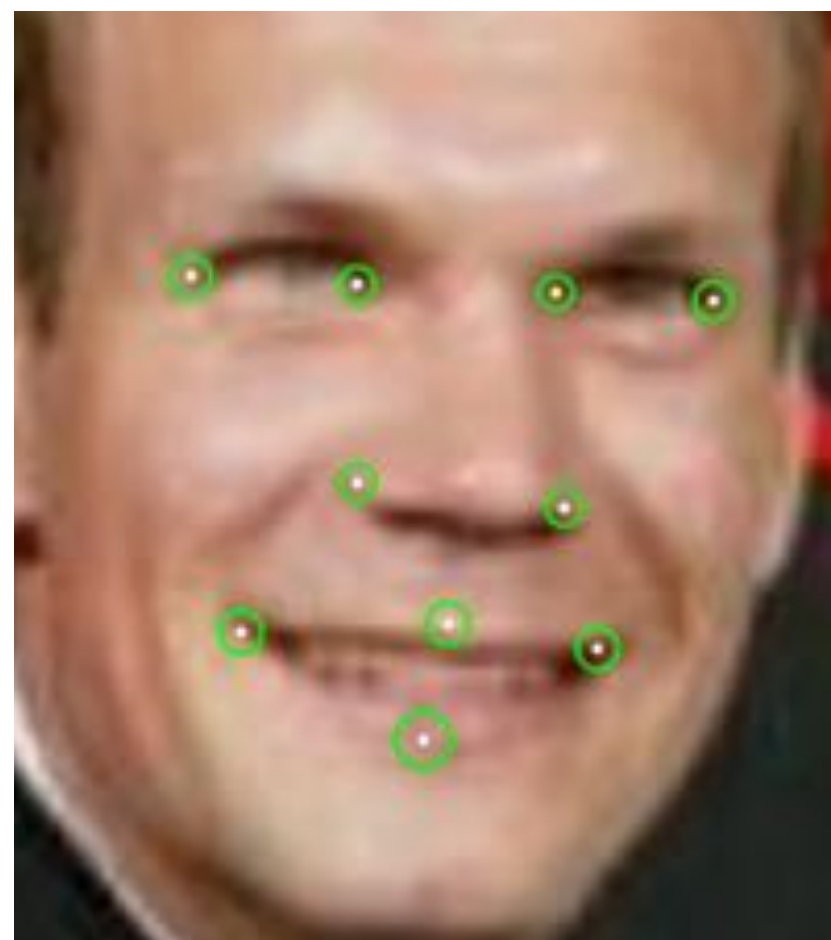
811 train  
244 test



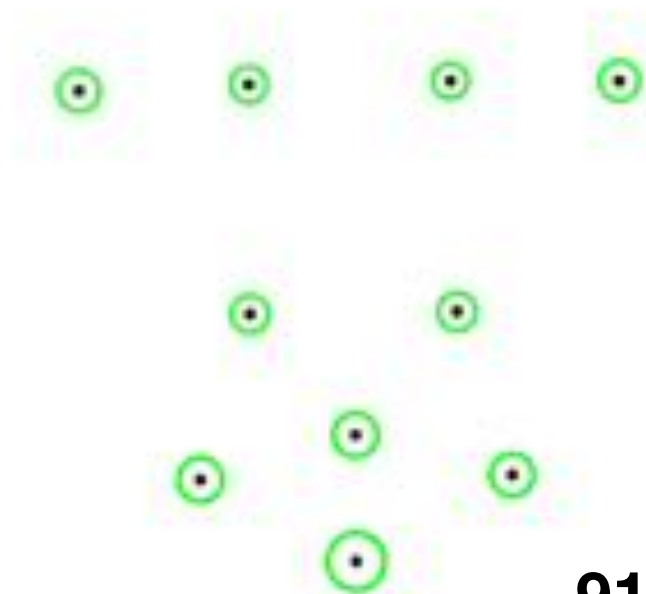
Helen Database



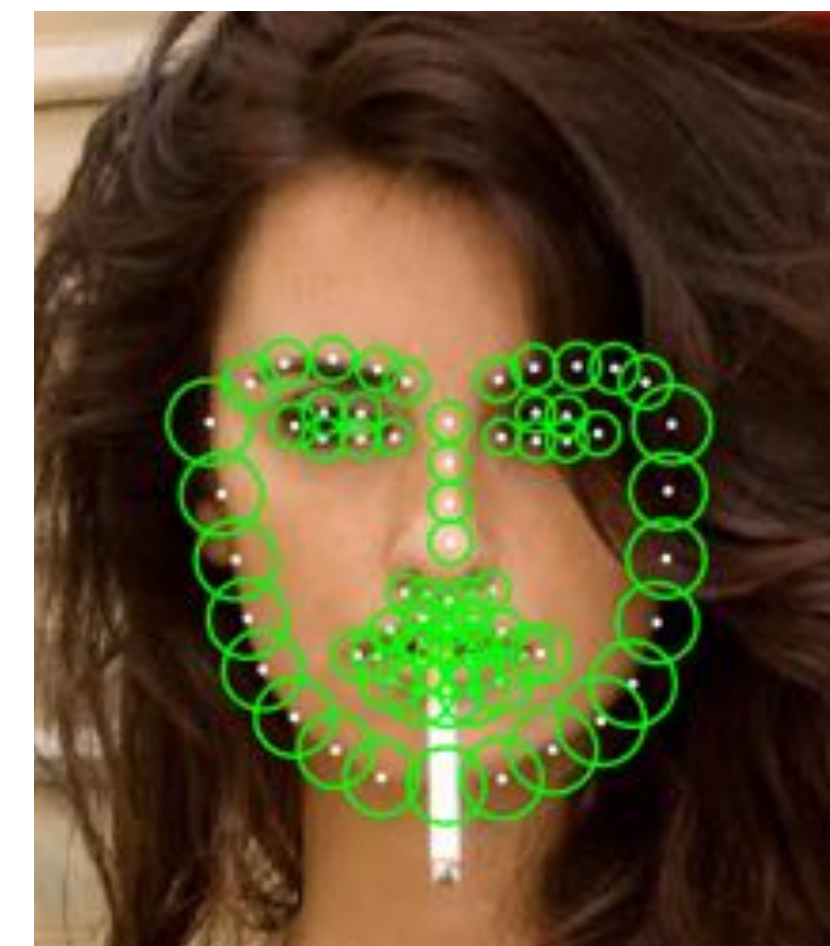
2000 train  
300 test



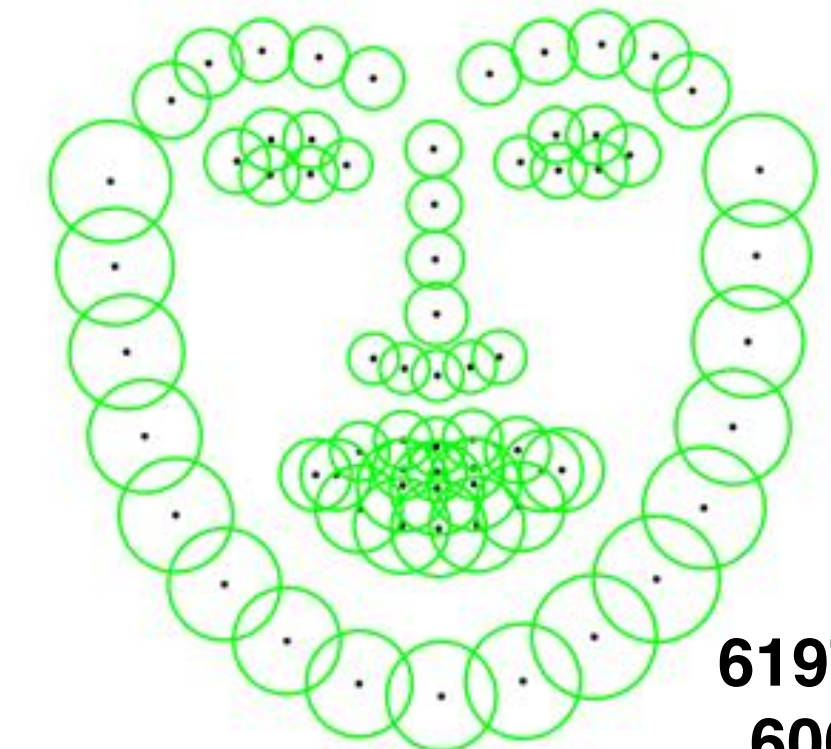
LFW Database



9100 train  
3900 test



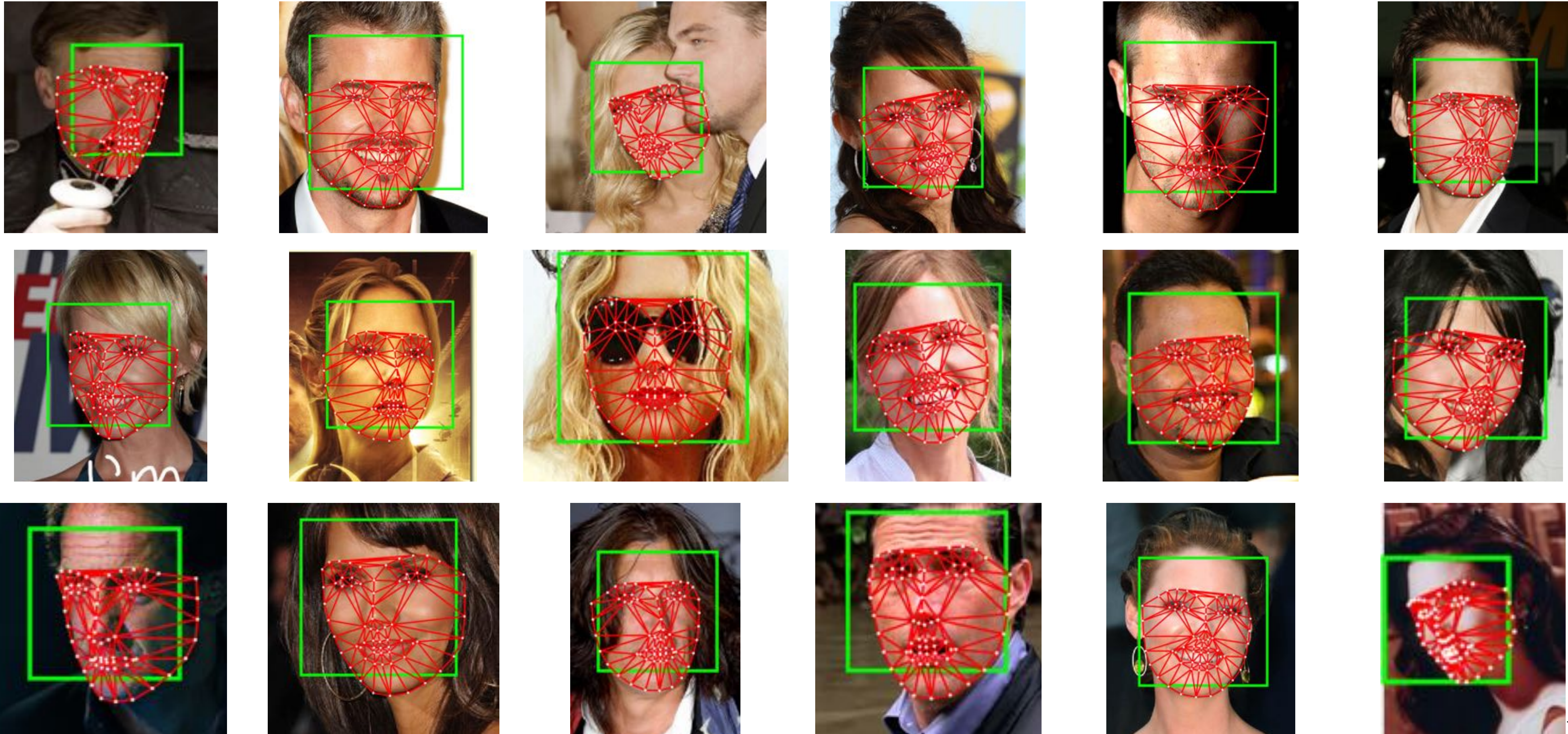
300W Database



6197 train  
600 test



# Qualitative Results - LFPW Database



# Qualitative Results - HELEN Database



# Conclusions

---

- Proposed an improved face alignment approach (facial landmark localization) w/ deformable face model.
- Nonlinear Cascaded Regression Extension.
  - Compactness of shape model.
  - Enforced shape consistency in Regression (reduced regression effort).
  - Bootstrap model to augment training data (CNN).
  - Loss function, w/ shape aware weighting.
- Demonstrated results in LFPW, HELEN, LFW and 300W datasets.
  - Improvement of  $\sim 1.5\%$  AUC (average across all datasets).

# Thank you

---

<https://www.isr.uc.pt/~pedromartins>  
[pedromartins@isr.uc.pt](mailto:pedromartins@isr.uc.pt)

