

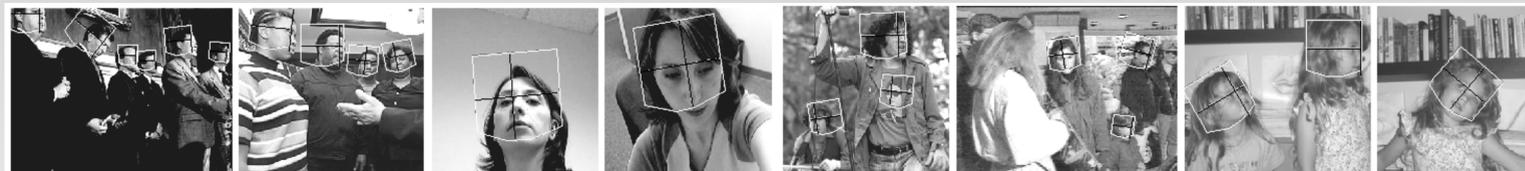
# Synergistic Face Detection and Pose Estimation with Energy-Based Models

M. Osadchy *Technion*   M. Miller *NEC Labs*   Y. LeCun *NYU*

We developed a method for

Simultaneous face detection and pose estimation.  
Robust to: yaw (from left to right profile), roll (-45, 45), and pitch (-60, 60).

Single Detector is applied to all poses.  
Pose estimation: Within 15° error about 90% of poses are estimated correctly.  
Near real-time: 5 frames per second on standard hardware.

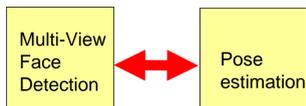


## Integrating face detection and pose estimation

### Synergy

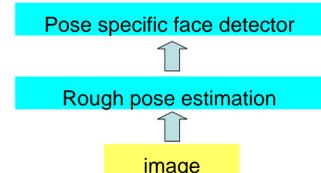
Common Problems:

Inner class variation (skin color, hair style, etc.)  
Lighting Variations  
Scale Variations  
Facial Expressions  
...



Train together → Better generalization

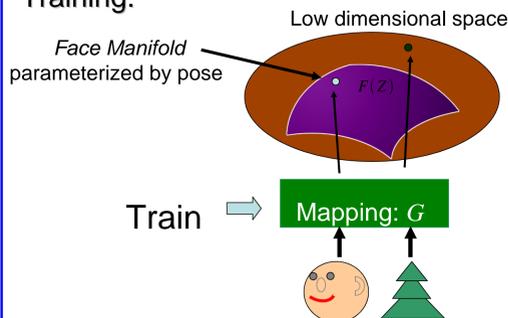
### Previous Methods



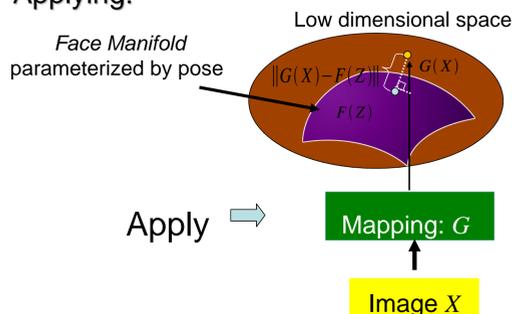
Unmanageable in real problems

### Our Approach

#### Training:



#### Applying:



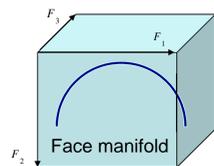
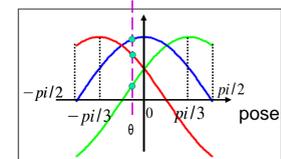
### Parameterization of the face space

#### Single pose parameter

Yaw:  $Z = \theta = [-\pi/2, \pi/2]$

$$F_i(\theta) = \cos(\theta - \alpha_i) \quad i=1,2,3$$

$$\alpha = [-\pi/3, 0, \pi/3]$$



$$\theta = \arctan \frac{\sum_{i=1}^3 G_i \cos \alpha_i}{\sum_{i=1}^3 G_i \sin \alpha_i}$$

#### More Pose Parameters

Yaw and roll  $Z = (\theta, \phi)$ :

$\theta = [-\pi/2, \pi/2]$   
 $\phi = [-\pi/4, \pi/4]$  } a portion of the surface of a sphere

$$F_{ij}(\theta, \phi) = \cos(\theta - \alpha_i) \cos(\phi - \beta_j); \quad i, j = 1, 2, 3$$

$$\alpha, \beta = [-\pi/3, 0, \pi/3]$$

$$\theta = 0.5 (\text{atan2}(cs + sc, cc - ss) + \text{atan2}(sc - cs, cc + ss))$$

$$\phi = 0.5 (\text{atan2}(cs + sc, cc - ss) - \text{atan2}(sc - cs, cc + ss))$$

where

$$cc = \sum_{ij} G_{ij}(x) \cos(\alpha_i) \cos(\beta_j) \quad cs = \sum_{ij} G_{ij}(x) \cos(\alpha_i) \sin(\beta_j)$$

$$ss = \sum_{ij} G_{ij}(x) \sin(\alpha_i) \sin(\beta_j) \quad sc = \sum_{ij} G_{ij}(x) \sin(\alpha_i) \cos(\beta_j)$$

### Minimum Energy Machine

Energy function:  $E_w(Y, Z, X)$     $[Z] = [-90, 90] \times [-45, 45]$   
parameters:  $Y$  (label),  $Z$  (pose),  $X$  (image)  
measures compatibility between  $X, Z, Y$ .

If  $X$  is a face with pose  $Z$  then we want:

$$E_w(1, Z, X) < E_w(0, Z, X) \quad \forall Z;$$

$$E_w(1, Z, X) < E_w(1, Z', X) \quad \forall Z' \neq Z$$

### Operating the Machine

Clamp  $X$  to the observed value (the image)

Find  $Z$  and  $Y$  such that:

$$(Y, Z) = \underset{Y \in \{1, 0\}, Z \in [Z]}{\text{argmin}} E_w(Y, Z, X)$$

Complete energy:

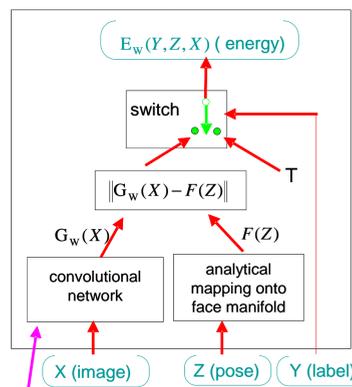
$$E_w(Y, Z, X) = Y \cdot \|G_w(X) - F(Z)\| + (1 - Y) \cdot T$$

$$X \text{ is a face, } Y=1 \rightarrow E_w(1, Z, X) = \|G_w(X) - F(Z)\|$$

$$X \text{ is not a face, } Y=0 \rightarrow E_w(0, Z, X) = T$$

## Learning Machine

### Architecture

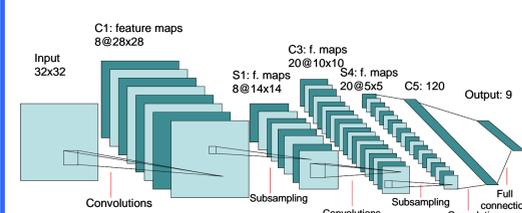


#### Operating the machine:

$$\bar{Z} = \underset{Z \in [Z]}{\text{argmin}} \|G_w(X) - F(Z)\|$$

$$\bar{Y} = \begin{cases} 1 & \|G_w(X) - F(\bar{Z})\| < T \\ 0 & \text{otherwise} \end{cases}$$

### Convolutional Network:



"end-to-end" trainable systems from low-level features to high-level representations.

Easily learn the type of shift-invariant features, relevant to object recognition.

Can be replicated over large images much more efficiently than traditional classifiers.

### Discriminative Loss Function

Minimize:  $L(W) = \frac{1}{|S_1|} \sum_{i \in S_1} L_1(W, Z^i, X^i) + \frac{1}{|S_0|} \sum_{i \in S_0} L_0(W, X^i)$   
loss for face sample with known pose   loss for non-face sample  
training faces   training non-faces

$$L_1(W, 1, Z, X) = E_w(1, Z, X)^2 \quad L_0(W, 0, X) = K \exp[-E(1, Z, X)]$$

We showed that this loss function causes the machine to exhibit proper behavior:  $E(Y^{\text{desired}}, \dots) < E(Y^{\text{undesired}}, \dots) + \text{margin}$

### Running the Machine

Works on grey-level images.

Applied at range of scales stepping by a factor of  $\sqrt{2}$

The network is replicated over the image at each scale, stepping by 4 pixels in  $x$  and  $y$ .

Overlapping detections are replaced by the strongest.

## Results

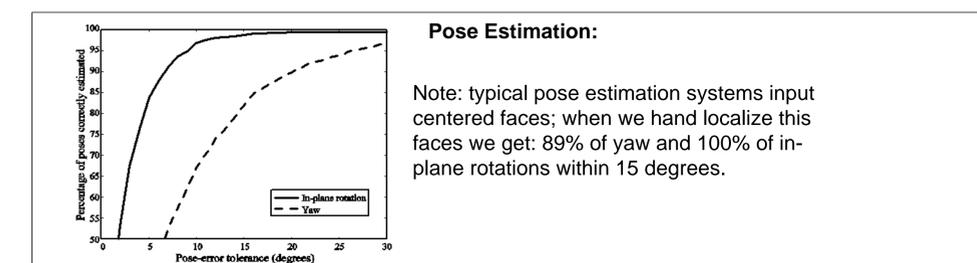
### Training

"52,850, 32x32 grey-level images of faces (NEC Labs hand annotated set) with uniform distribution of poses.  
Initial negative set: 52,850 random non-face natural images.  
Second phase: half of the initial negative set was replaced by false positives of the initial version of the detector.  
Each training image was used 5 times with random variation in scale, in-plane rotation, brightness and contrast.  
9 passes on the data: 26 hours on 2Ghz Pentium 4.  
The system converged to an EER of 5% on training set and 6% on test set of 90,000 images.

### Test on Standard Data Sets

No standard set tests all poses, that our system is designed to detect.  
3 standard sets focusing on particular pose variation: tilted, profile, and frontal.

Data Set →	Detection:					
	TILTED	PROFILE	MIT+CMU			
False positives per image →	4.42	26.9	0.47	3.36	0.5	1.28
<b>Our Detector</b>	<b>90%</b>	<b>97%</b>	<b>67%</b>	<b>83%</b>	<b>83%</b>	<b>88%</b>
Jones & Viola (tilted)	90%	95%	x	x	x	x
Jones & Viola (profile)	x	x	70%	83%	x	x
Rowley <i>et al</i>	89%	96%			x	x
Schneiderman & Kanade			86%	93%	x	x



### Synergy Test

