# Synergistic Face Detection and Pose Estimation

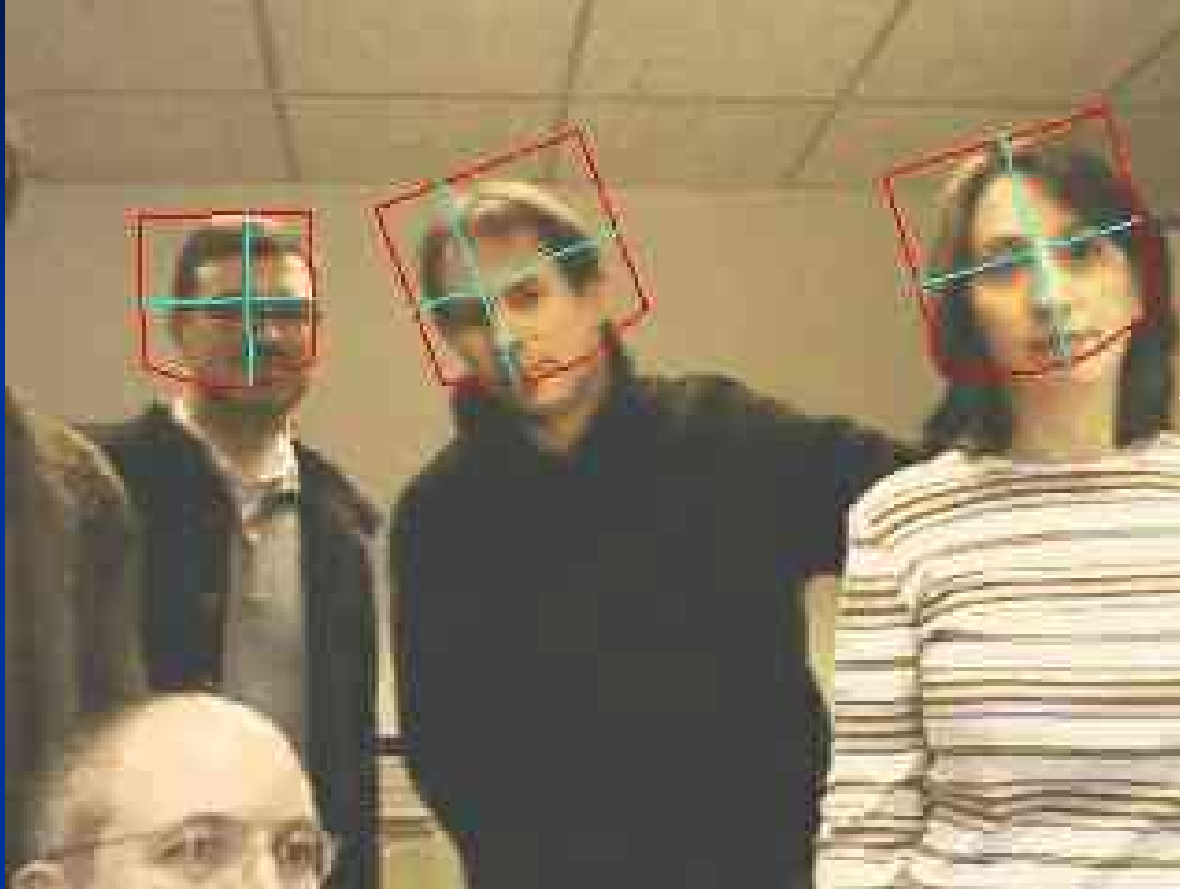M. Osadchy        M.Miller        Y. LeCun

Technion        NEC Labs        NYU

# Our System



- Detects faces independently of their poses.

- Estimates head poses.

# Our System

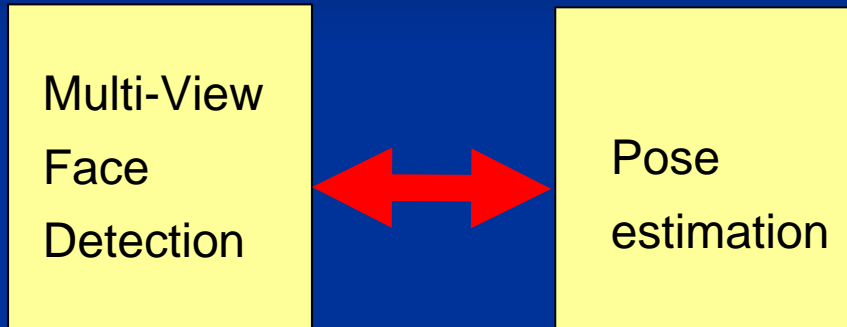Robust to: yaw (from left to right profile), roll (-45, 45), and pitch (-60, 60).

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.

# Synergy

closely related
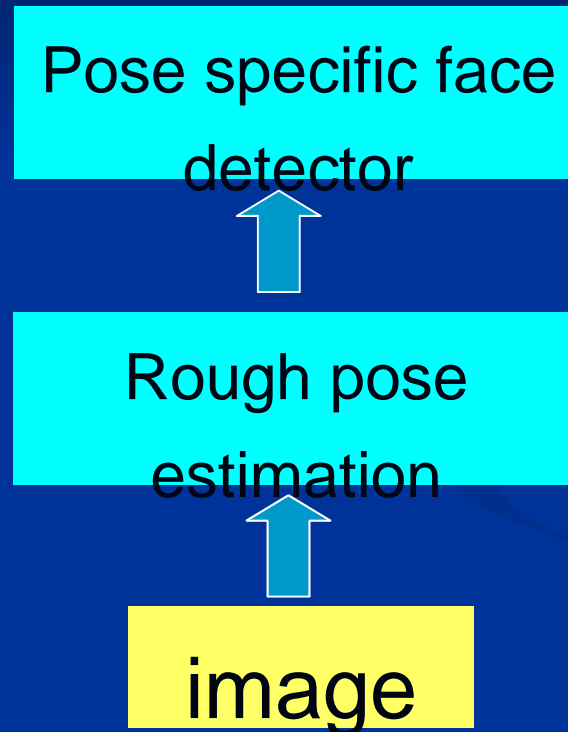
Multi-View Face Detection ⟷ Pose estimation

- Inner class variation (skin color, hair style, etc.)

- Lighting Variations

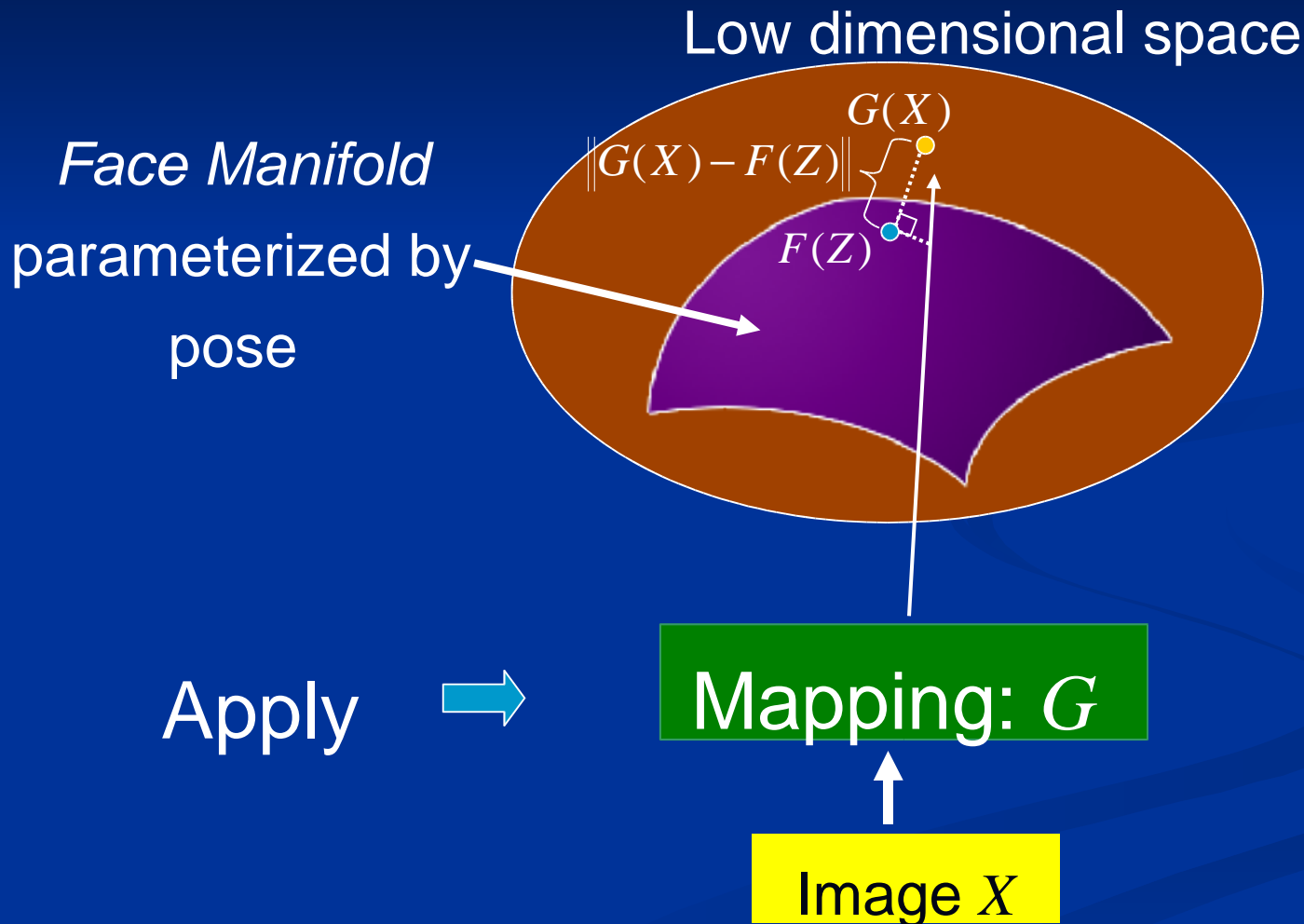- Scale  Variations

- Facial Expressions

- …

Train together ⟶ Better generalization

# Integrating Face Detection and Pose Estimation: Previous Methods

Pose specific face detector

⬆
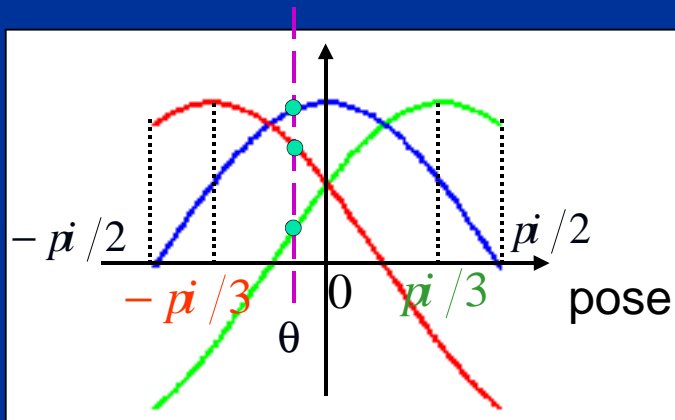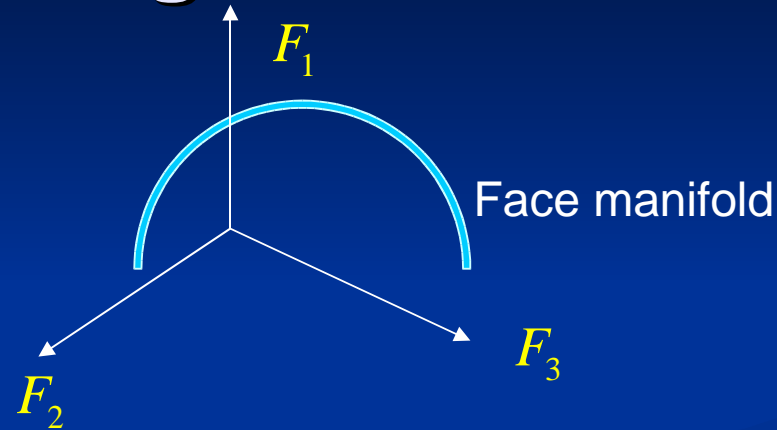
Rough pose estimation

⬆

image

Unmanageable in real problems

# Integrating Face Detection and Pose Estimation: Our Approach

# Parameterization of the Face Manifold – Single Parameter

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

$F_1$

Face manifold

$F_3$

$F_2$

$F_i(\theta) = \cos(\theta - \alpha_i)$ $\qquad \alpha = [-pi/3, 0, pi/3]$

$i = 1, 2, 3$

$-pi/2$ $\qquad$ $pi/2$

$-pi/3$ $\quad 0 \quad pi/3$ $\qquad$ pose

$\theta$

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

# Parameterization of the Face Manifold – Two Parameters

Yaw and roll  $Z = (\theta, \varphi)$ ;

$$\left. \begin{array}{l} \theta = [-pi/2, \; pi/2] \\ \varphi = [-pi/4, \; pi/4] \end{array} \right\} \cong \quad \text{a portion of the surface of a sphere}$$

$$F_{ij}(\theta, \varphi) = \cos(\theta - \alpha_i)\cos(\varphi - \beta_j); \quad i, j = 1,2,3 \quad \alpha, \beta = [-pi/3, 0, pi/3]$$

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

$$\overline{\varphi} = 0.5(\text{atan2}\left(cs + sc, cc - ss\right) - \text{atan2}\left(sc - cs, cc + ss\right)$$

where

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

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

# Minimum Energy Machine

parameters

- Energy function:

$$E_W(Y, Z, X)$$

label    pose    image

$$\{Z\} = [-90, 90] \times [-45, 45]$$

$$Y = \begin{cases} 1 & \text{face} \\ 0 & \text{non face} \end{cases}$$

- $E_W(Y, Z, X)$ 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:
$$\left(\overline{Y}, \overline{Z}\right) = \underset{Y \in \{Y\}, Z \in \{Z\}}{\arg \min} \ E_W\left(Y, \ Z, \ X\right)$$

- Complete energy:
$$E_W\left(Y, \ Z, \ X\right) = Y \cdot \left\| G_W(X) - F(Z) \right\| + \left(1 - X\right) \cdot T$$

| X is a face | $\Rightarrow$ | Y=1 |

# Operating the Machine

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

- Find $Z$ and $Y$ such that:

$$\left(\overline{Y}, \overline{Z}\right) = \arg \min_{Y \in \{Y\}, Z \in \{Z\}} E_W\left(Y, Z, X\right)$$

- Complete energy:
$$E_W\left(Y, Z, X\right) = Y \cdot \left\|G_W(X) - F(Z)\right\| + \left(1 - Y\right) \cdot T$$

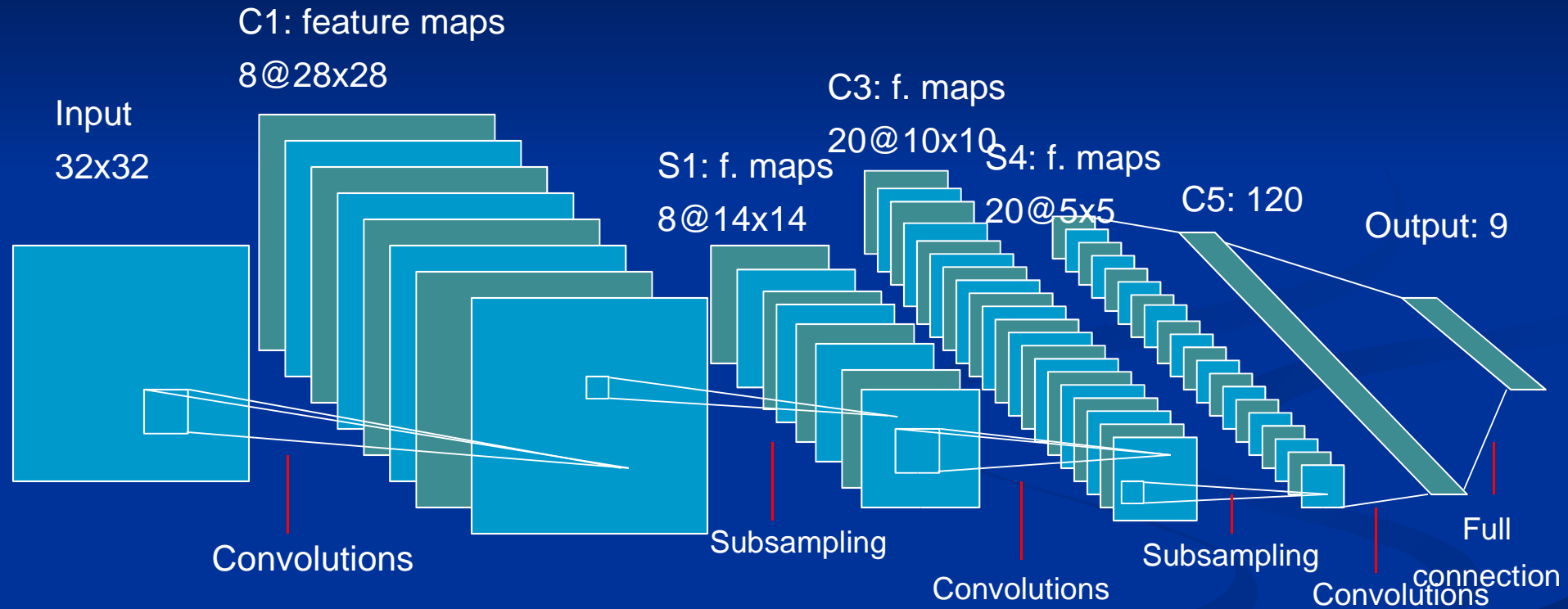| X is not a face | $\Rightarrow$ | Y=0 |
|---|---|---|

# Architecture

# 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.

Considerable advantage for real-time systems!

# Similar to LeNet5, with more maps:

# Training with Discriminative Loss Function
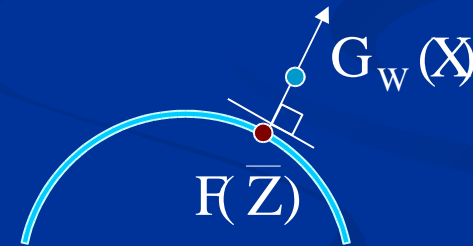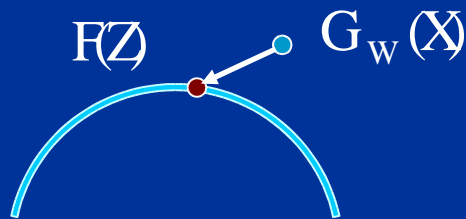
loss for face sample with known pose

loss for non-face sample

Minimize:
$$L(W) = \frac{1}{|S_1|} \sum_{i \in S_1} L_1\left(W, Z^i, X^i\right) + \frac{1}{|S_0|} \sum_{i \in S_0} L_0\left(W, X^i\right)$$

training faces

training non-faces

$$L_1(W, 1, Z, X) = E_W(1, Z, X)^2$$

$$L_0(W, 0, X) = K \exp\left[-E(1, \overline{Z}, X)\right]$$

$F(Z)$   $G_W(X)$

$G_W(X)$   $F(\overline{Z})$

We showed that this loss function causes the machine to exhibit proper behavior:

$$E\left(Y^{\text{desired}}, \ldots\right) < E\left(Y^{\text{undesired}}, \ldots\right) + \text{margin}$$

# 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

- Our system is robust to yaw $\pm 90$, in-plane rotation $\pm 60$, and pitch $\pm 45$

# 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.
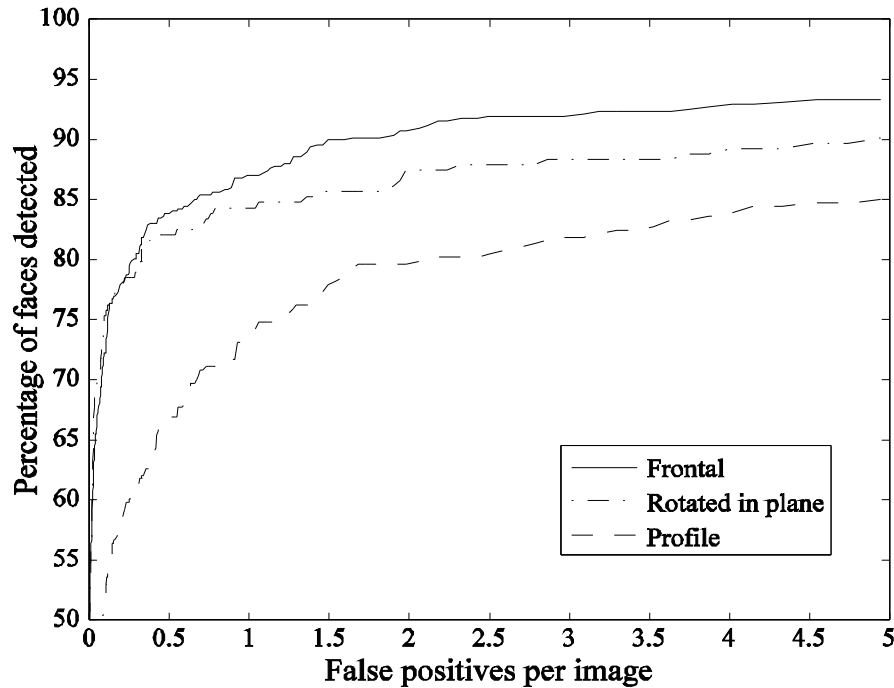
# Test on Standard Data Sets

- No standard set tests all poses, that our system is designed to detect.

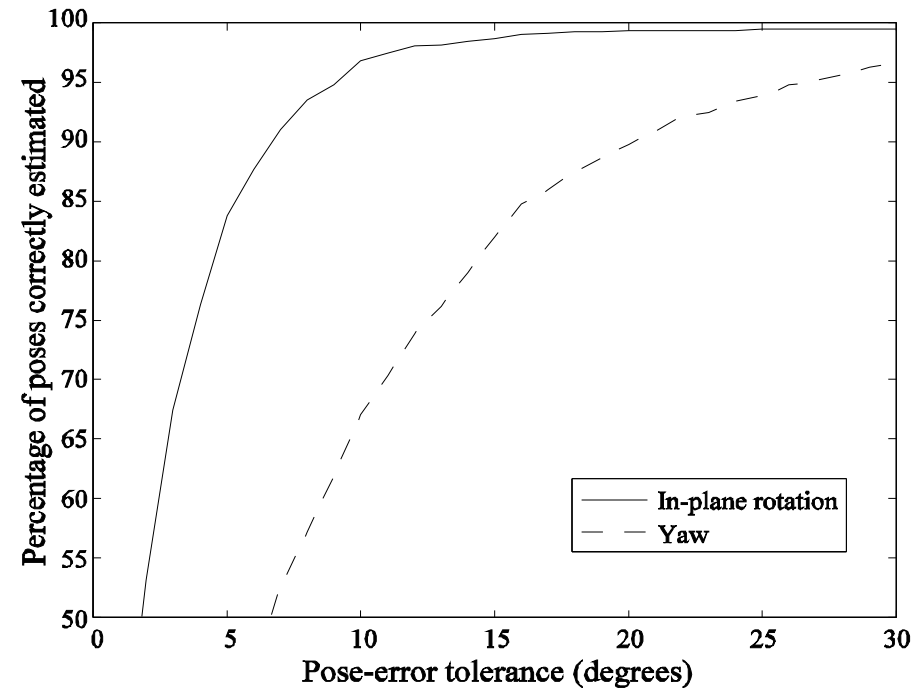- 3 standard sets focusing on particular pose variation: tilted, profile, and frontal.

| | Data Set-> | TILTED | | PROFILE | | MIT+CMU | |
|---|---|---|---|---|---|---|---|
| | *False positives per image->* | 4.42 | 26.9 | 0.47 | 3.36 | 0.5 | 1.28 |
| **Real time** | **Our Detector** | **90%** | **97%** | **67%** | **83%** | **83%** | **88%** |
| | **Jones & Viola (tilted)** | **90%** | **95%** | x | | x | |
| | **Jones & Viola (profile)** | x | | **70%** | **83%** | x | |
| | Rowley *et al* | 89% | 96% | | | | x |
| | Schneiderman & Kanade | | | 86% | 93% | | x |

# Standard Sets

## Detection
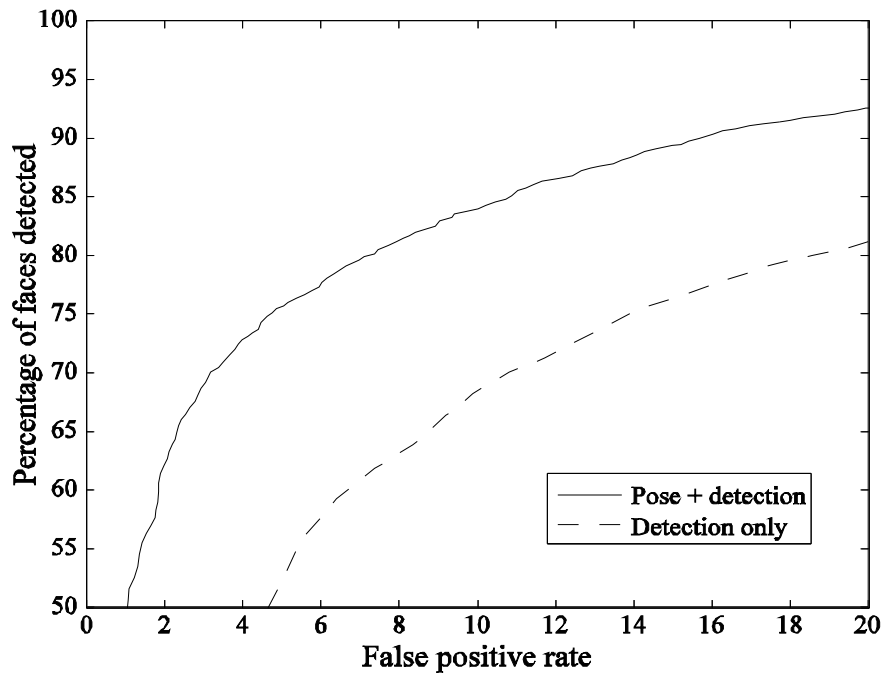
## Pose Estimation of the detected faces



Note: typical pose estimation systems input centered faces; when we hand localize this faces we get: 89% of yaw and 100% of in-plane rotations within 15 degrees.

# Synergy Test

## Detection       Pose Estimation