

Latent Variables

Yann LeCun

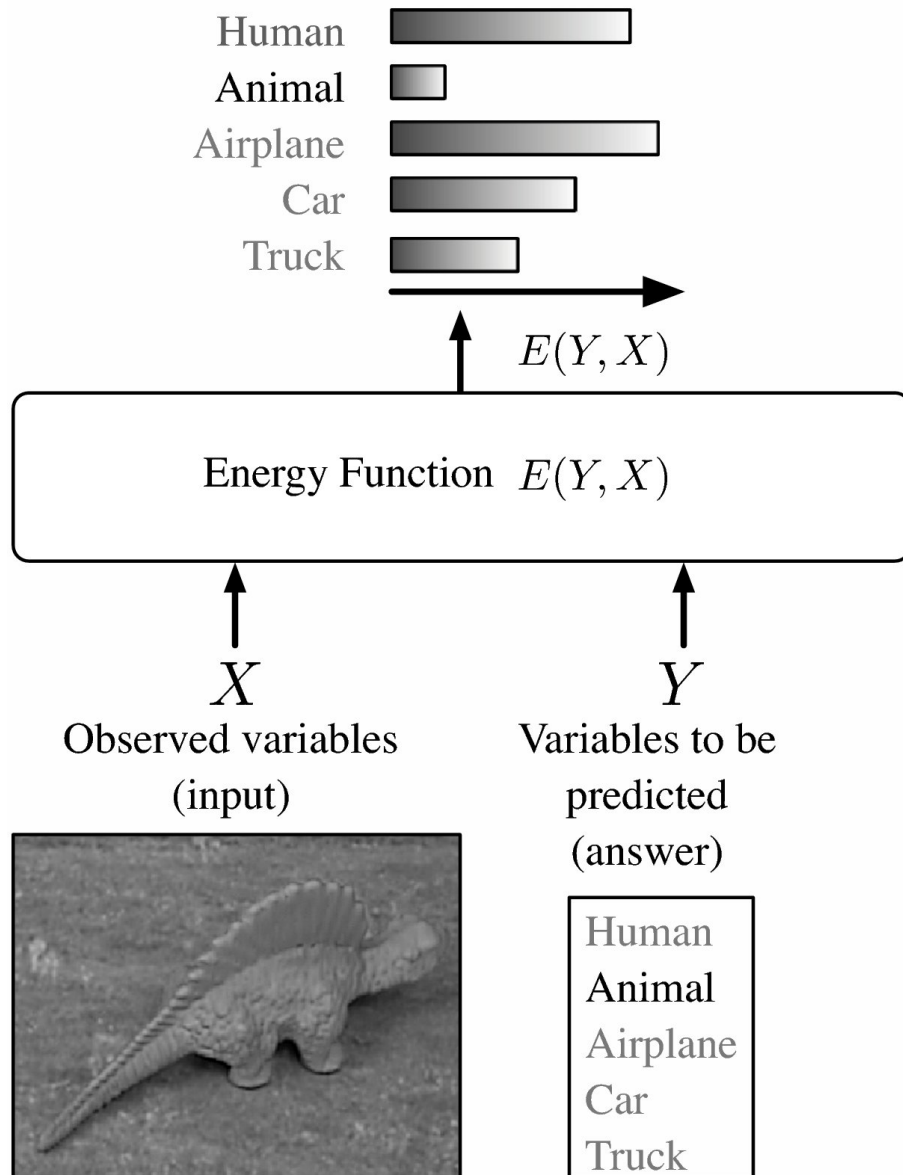
The Courant Institute of Mathematical Sciences

New York University

<http://yann.lecun.com>

<http://www.cs.nyu.edu/~yann>

Energy-Based Model for Decision-Making



• **Model:** Measures the compatibility between an observed variable X and a variable to be predicted Y through an energy function $E(Y, X)$.

$$Y^* = \operatorname{argmin}_{Y \in \mathcal{Y}} E(Y, X).$$

- **Inference:** Search for the Y that minimizes the energy within a set \mathcal{Y} .
- If the set has low cardinality, we can use exhaustive search.

Transforming Energies to Probabilities

• Energies are uncalibrated

- ▶ The energies of two separately-trained systems cannot be combined
- ▶ The energies are uncalibrated (measured in arbitrary units)

• How do we calibrate energies?

- ▶ We turn them into probabilities (positive numbers that sum to 1).
- ▶ Simplest way: Gibbs distribution
- ▶ Other ways can be reduced to Gibbs by a suitable redefinition of the energy.

$$P(Y|X) = \frac{e^{-\beta E(Y,X)}}{\int_{y \in \mathcal{Y}} e^{-\beta E(y,X)}},$$

Partition function

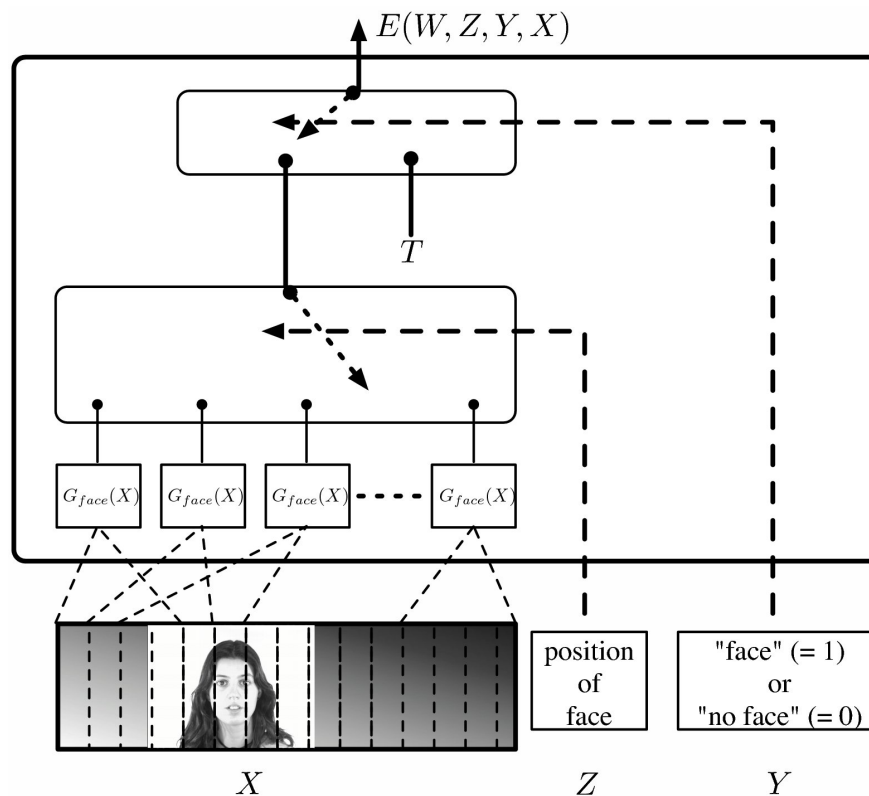
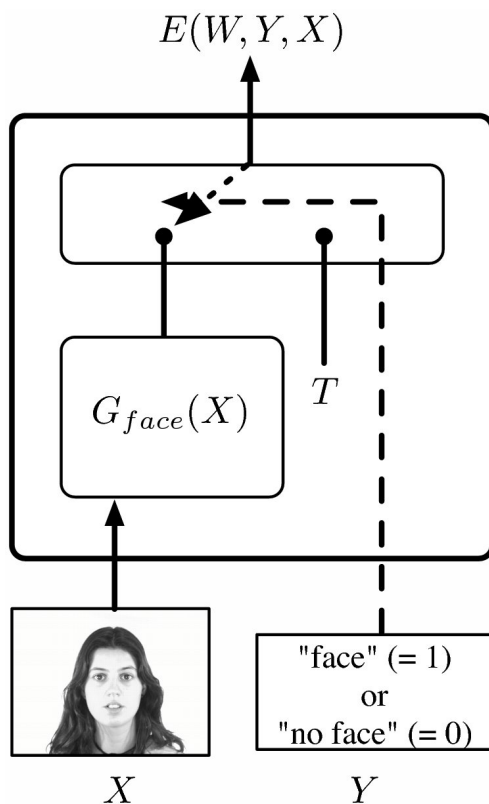
Inverse temperature

Latent Variable Models

- The energy includes “hidden” variables Z whose value is never given to us

$$E(Y, X) = \min_{Z \in \mathcal{Z}} E(Z, Y, X).$$

$$Y^* = \operatorname{argmin}_{Y \in \mathcal{Y}, Z \in \mathcal{Z}} E(Z, Y, X).$$



What can the latent variables represent?

- **Variables that would make the task easier if they were known:**
 - ▶ **Face recognition:** the gender of the person, the orientation of the face.
 - ▶ **Object recognition:** the pose parameters of the object (location, orientation, scale), the lighting conditions.
 - ▶ **Parts of Speech Tagging:** the segmentation of the sentence into syntactic units, the parse tree.
 - ▶ **Speech Recognition:** the segmentation of the sentence into phonemes or phones.
 - ▶ **Handwriting Recognition:** the segmentation of the line into characters.
- **In general, we will search for the value of the latent variable that allows us to get an answer (Y) of smallest energy.**

Probabilistic Latent Variable Models

- Marginalizing over latent variables instead of minimizing.

$$P(Z, Y | X) = \frac{e^{-\beta E(Z, Y, X)}}{\int_{y \in \mathcal{Y}, z \in \mathcal{Z}} e^{-\beta E(y, z, X)}} \cdot$$

$$P(Y | X) = \frac{\int_{z \in \mathcal{Z}} e^{-\beta E(Z, Y, X)}}{\int_{y \in \mathcal{Y}, z \in \mathcal{Z}} e^{-\beta E(y, z, X)}} \cdot$$

- Equivalent to traditional energy-based inference with a redefined energy function:

$$Y^* = \operatorname{argmin}_{Y \in \mathcal{Y}} - \frac{1}{\beta} \log \int_{z \in \mathcal{Z}} e^{-\beta E(z, Y, X)}.$$

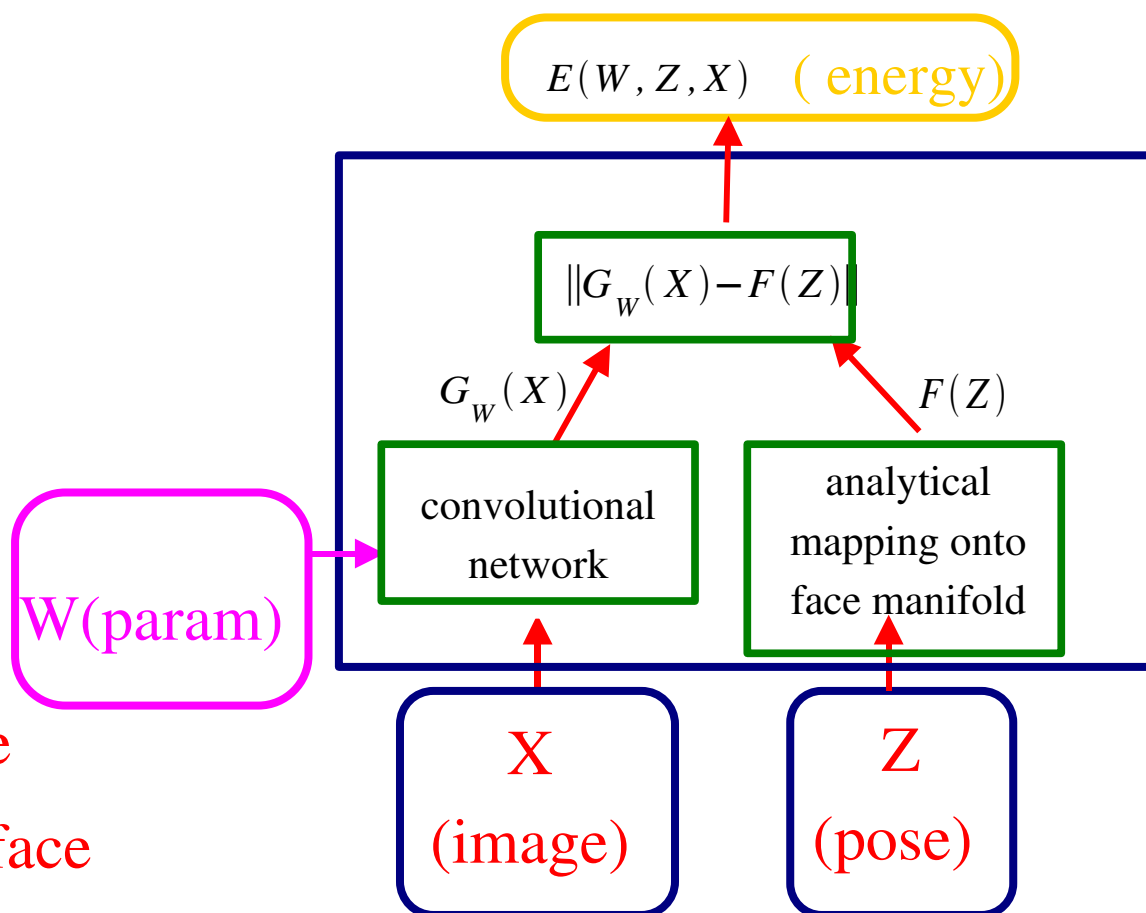
- Reduces to traditional minimization when Beta->infinity

Face Detection and Pose Estimation with a Convolutional EBM

- **Training:** 52,850, 32x32 grey-level images of faces, 52,850 selected non-faces.
- Each training image was used 5 times with random variation in scale, in-plane rotation, brightness and contrast.
- **2nd phase:** half of the initial negative set was replaced by false positives of the initial version of the detector .

$$E^*(W, X) = \min_Z \|G_W(X) - F(Z)\|$$

$$Z^* = \operatorname{argmin}_Z \|G_W(X) - F(Z)\|$$

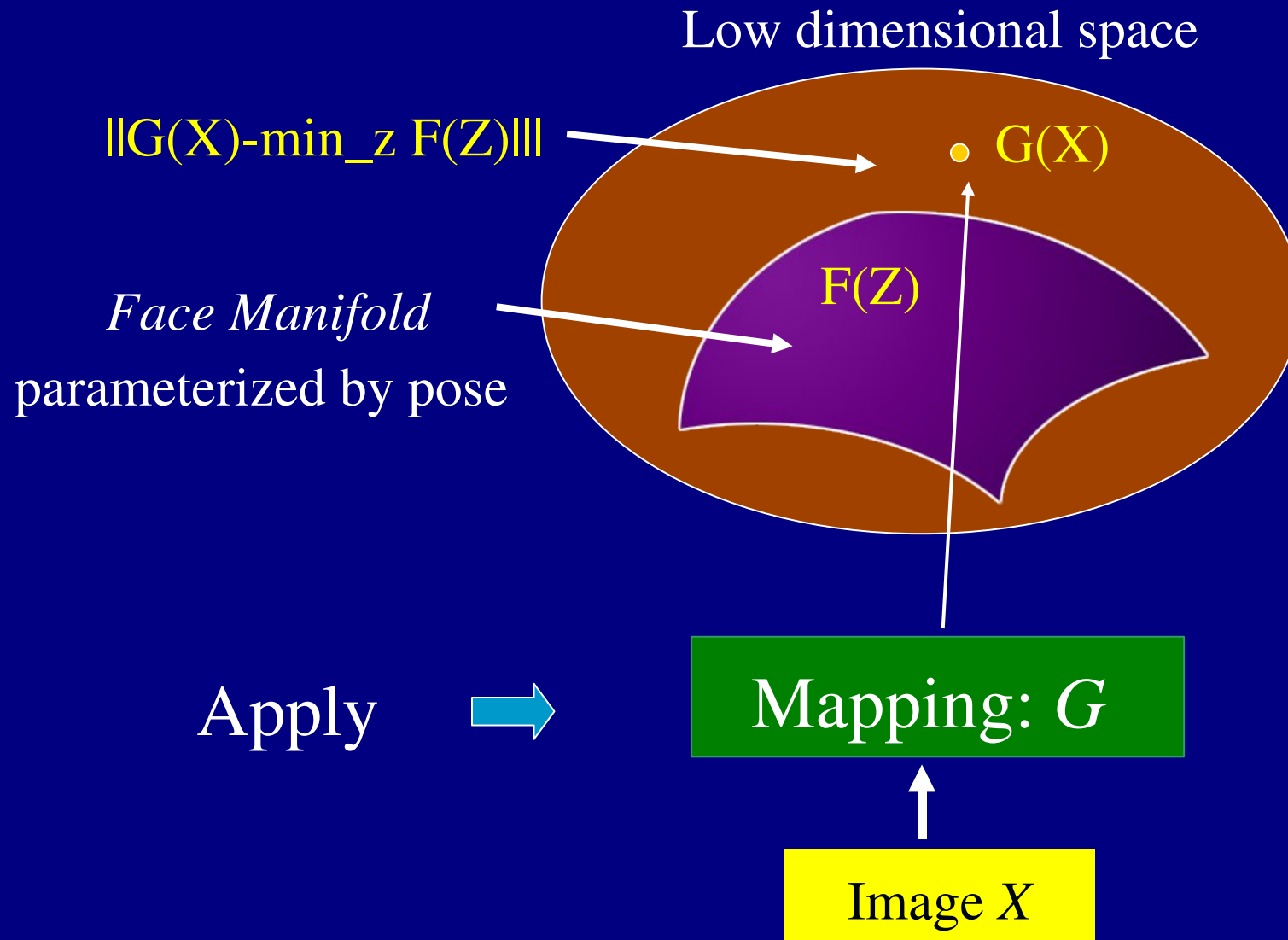


Small $E^*(W, X)$: face

Large $E^*(W, X)$: no face

[Osadchy, Miller, LeCun, NIPS 2004]

Face Manifold



Probabilistic Approach: Density model of joint P(face,pose)

Probability that image
X is a face with pose Z

$$P(X, Z) = \frac{\exp(-E(W, Z, X))}{\int_{X, Z \in \text{images, poses}} \exp(-E(W, Z, X))}$$

Given a training set of faces annotated with pose, find the W that maximizes the likelihood of the data under the model:

$$P(\text{faces} + \text{pose}) = \prod_{X, Z \in \text{faces} + \text{pose}} \frac{\exp(-E(W, Z, X))}{\int_{X, Z \in \text{images, poses}} \exp(-E(W, Z, X))}$$

Equivalently, minimize the negative log likelihood:

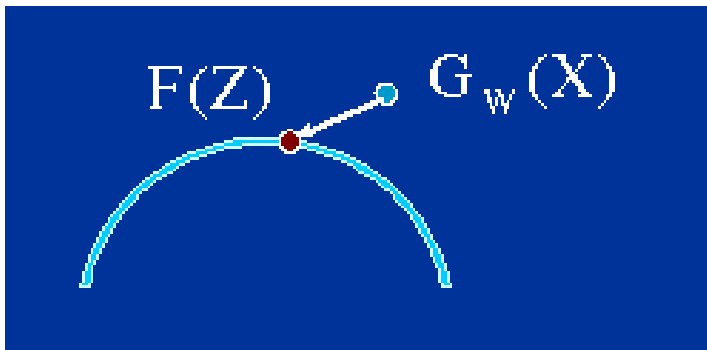
$$\mathcal{L}(W, \text{faces} + \text{pose}) = \sum_{X, Z \in \text{faces} + \text{pose}} E(W, Z, X) + \log \left[\int_{X, Z \in \text{images, poses}} \exp(-E(W, Z, X)) \right]$$


COMPLICATED

Energy-Based Contrastive Loss Function

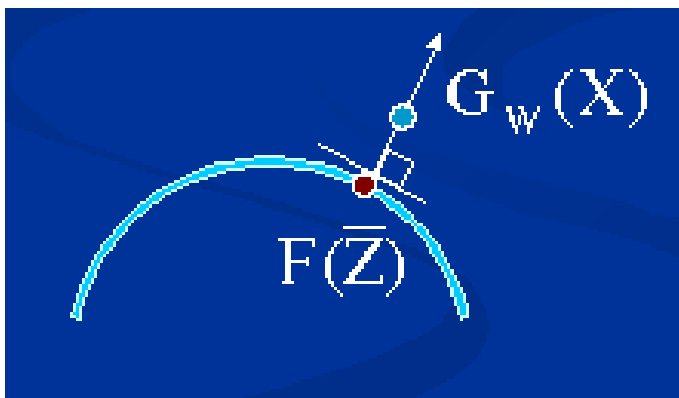
$$\mathcal{L}(W) = \frac{1}{|f + p|} \sum_{X, Z \in \text{faces} + \text{pose}} [L^+(E(W, Z, X))] + L^- \left(\min_{X, Z \in \text{bckgnd}, \text{poses}} E(W, Z, X) \right)$$

$$L^+(E(W, Z, X)) = E(W, Z, X)^2 = \|G_W(X) - F(Z)\|^2$$



Attract the network output $G_W(X)$ to the location of the desired pose $F(Z)$ on the manifold

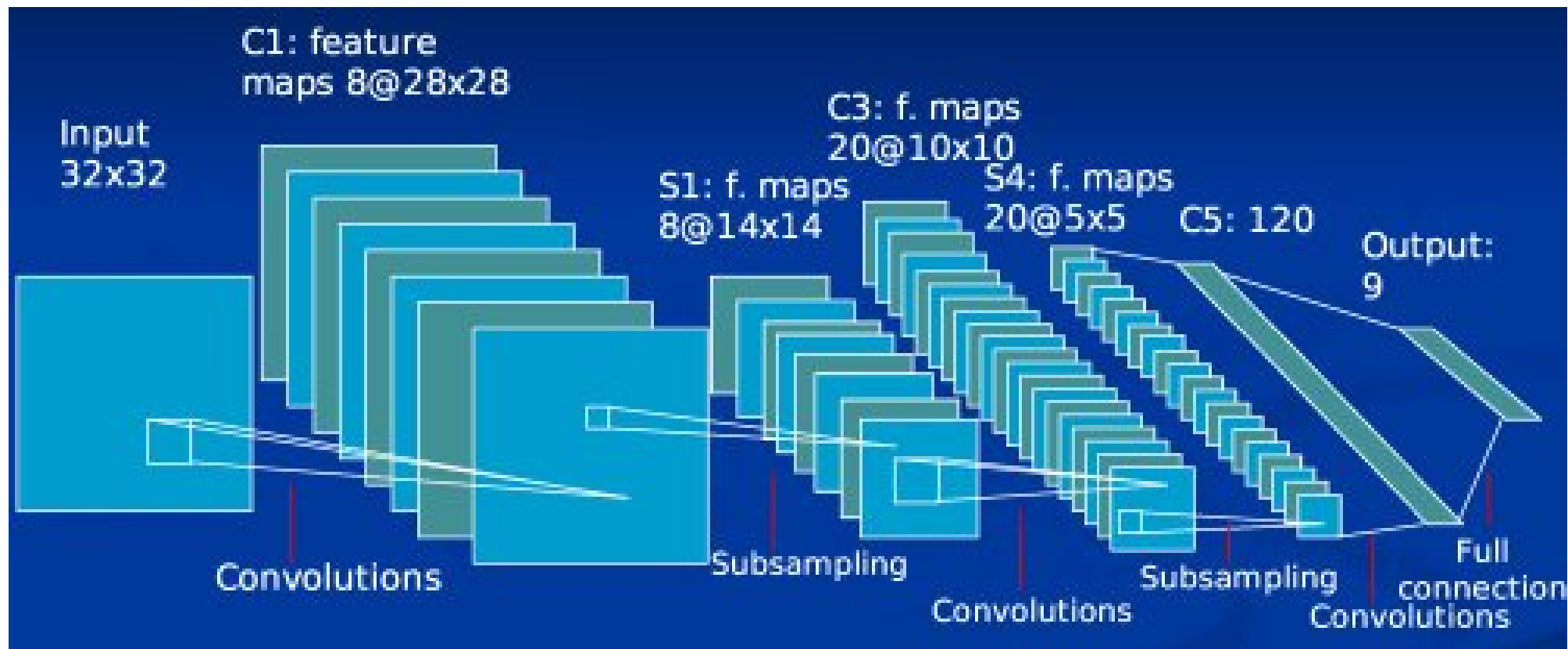
$$L^- \left(\min_{X, Z \in \text{bckgnd}, \text{poses}} E(W, Z, X) \right) = K \exp \left(-\min_{X, Z \in \text{bckgnd}, \text{poses}} \|G_W(X) - F(Z)\| \right)$$



Repel the network output $G_W(X)$ away from the face/pose manifold

Convolutional Network Architecture

[LeCun et al. 1988, 1989, 1998, 2005]



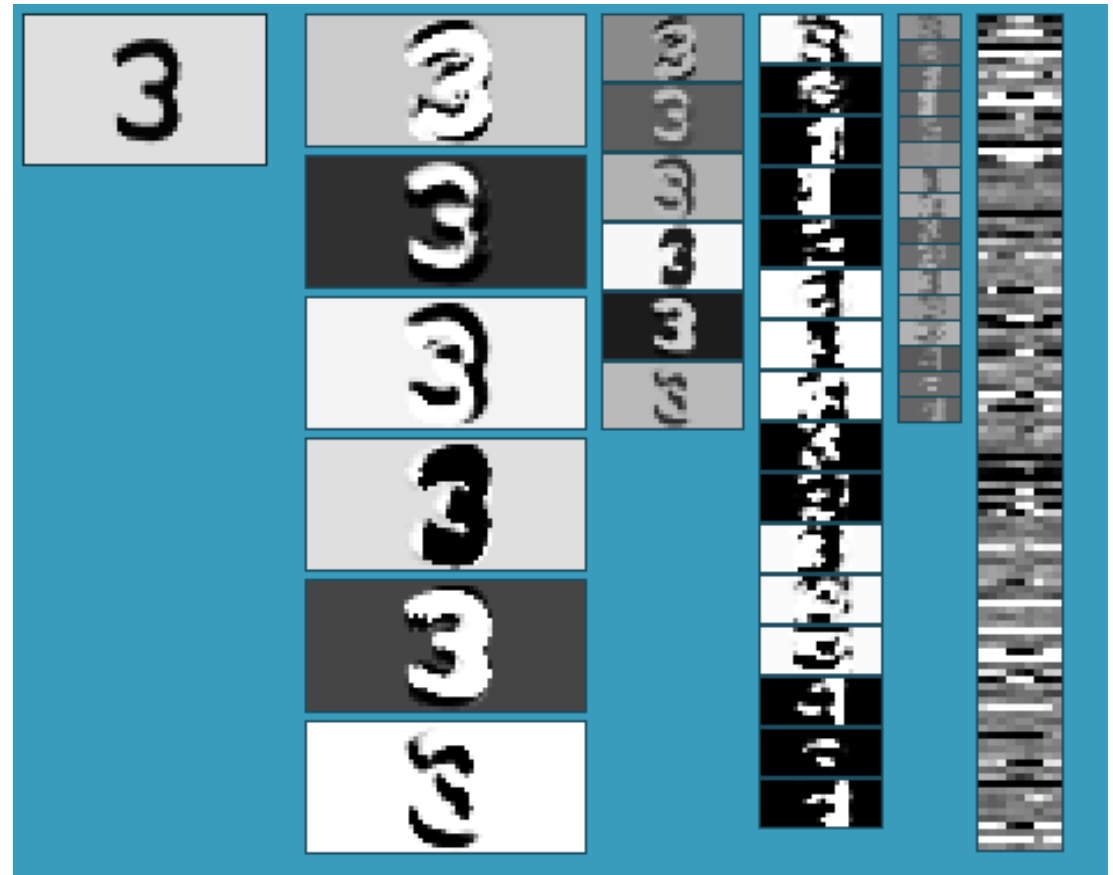
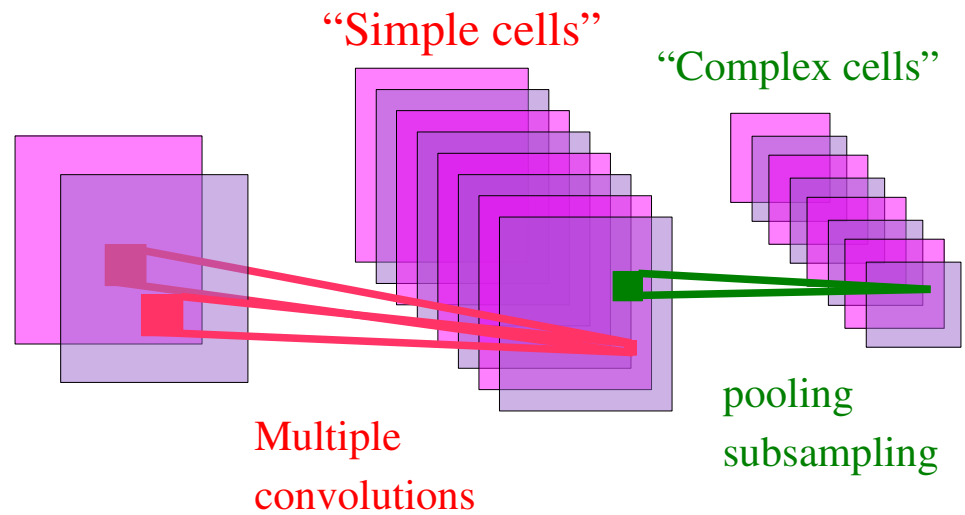
Hierarchy of local filters (convolution kernels),

sigmoid pointwise non-linearities, and spatial subsampling

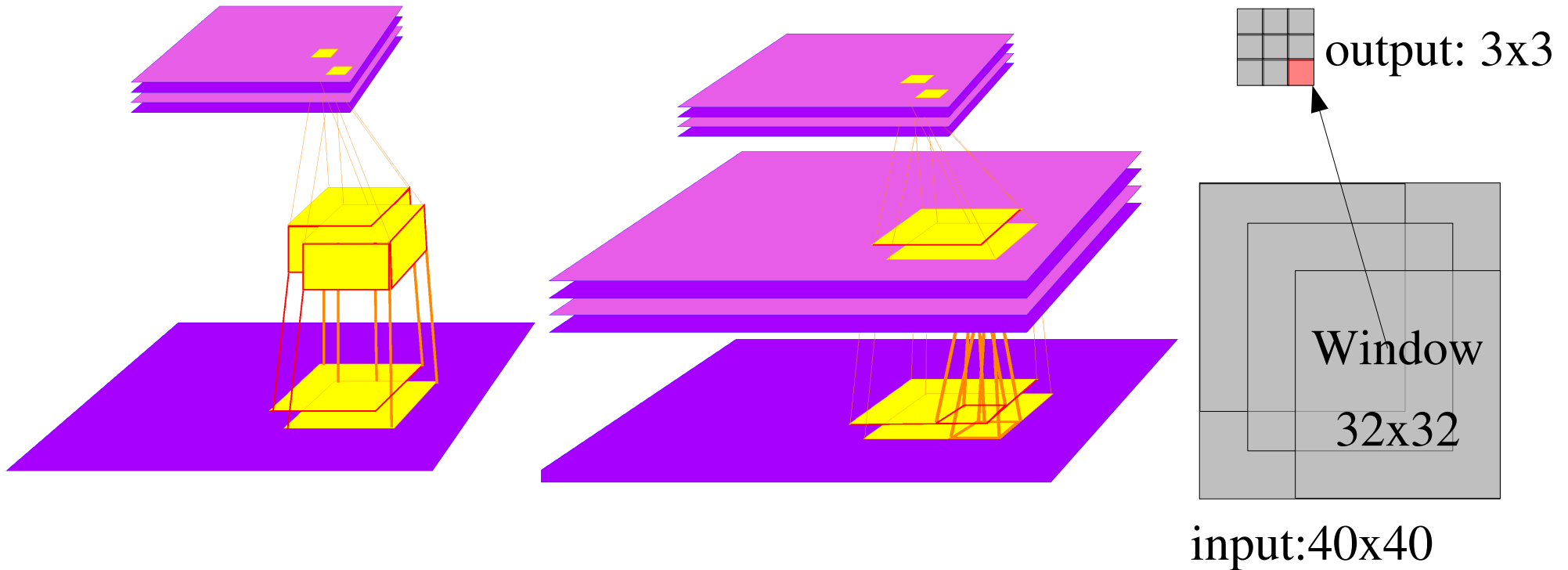
All the filter coefficients are learned with gradient descent (back-prop)

Alternated Convolutions and Pooling/Subsampling

- Local features are extracted everywhere.
- pooling/subsampling layer builds robustness to variations in feature locations.
- Long history in neuroscience and computer vision:
 - Hubel/Wiesel 1962,
 - Fukushima 1971-82,
 - LeCun 1988-06
 - Poggio, Riesenhuber, Serre 02-06
 - Ullman 2002-06
 - Triggs, Lowe,....



Building a Detector/Recognizer: Replicated Conv. Nets



- Traditional Detectors/Classifiers must be applied to every location on a large input image, at multiple scales.
- Convolutional nets can be replicated over large images very cheaply.
- The network is applied to multiple scales spaced by $\sqrt{2}$
- Non-maximum suppression with exclusion window

Building a Detector/Recognizer: Replicated Convolutional Nets

● Computational cost for replicated convolutional net:

● 96x96 -> 4.6 million multiply-accumulate operations

● 120x120 -> 8.3 million multiply-accumulate operations

● 240x240 -> 47.5 million multiply-accumulate operations

● 480x480 -> 232 million multiply-accumulate operations

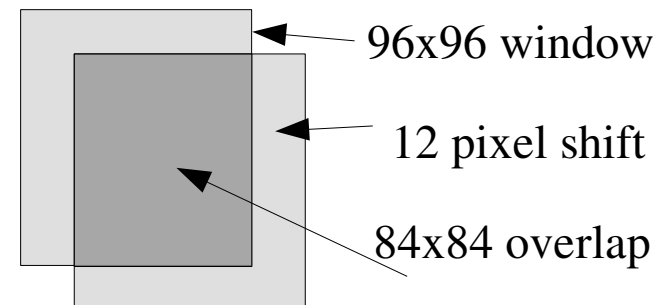
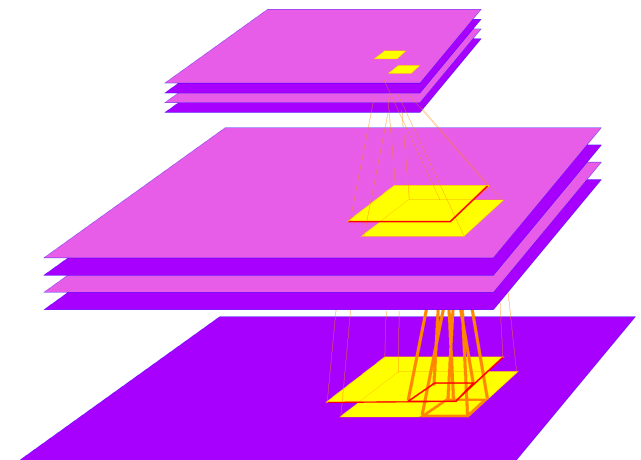
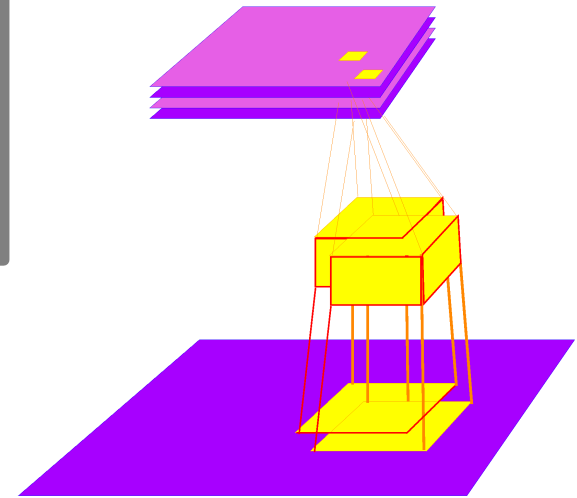
● Computational cost for a non-convolutional detector of the same size, applied every 12 pixels:

● 96x96 -> 4.6 million multiply-accumulate operations

● 120x120 -> 42.0 million multiply-accumulate operations

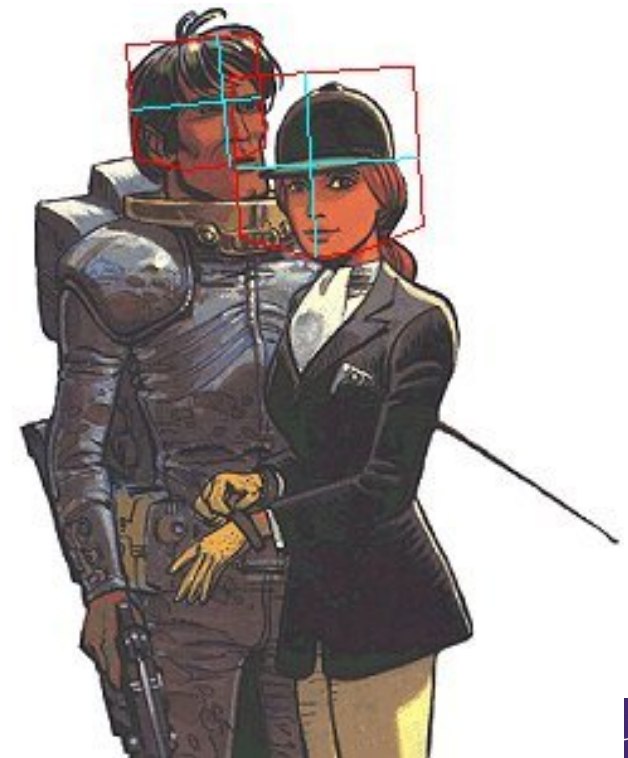
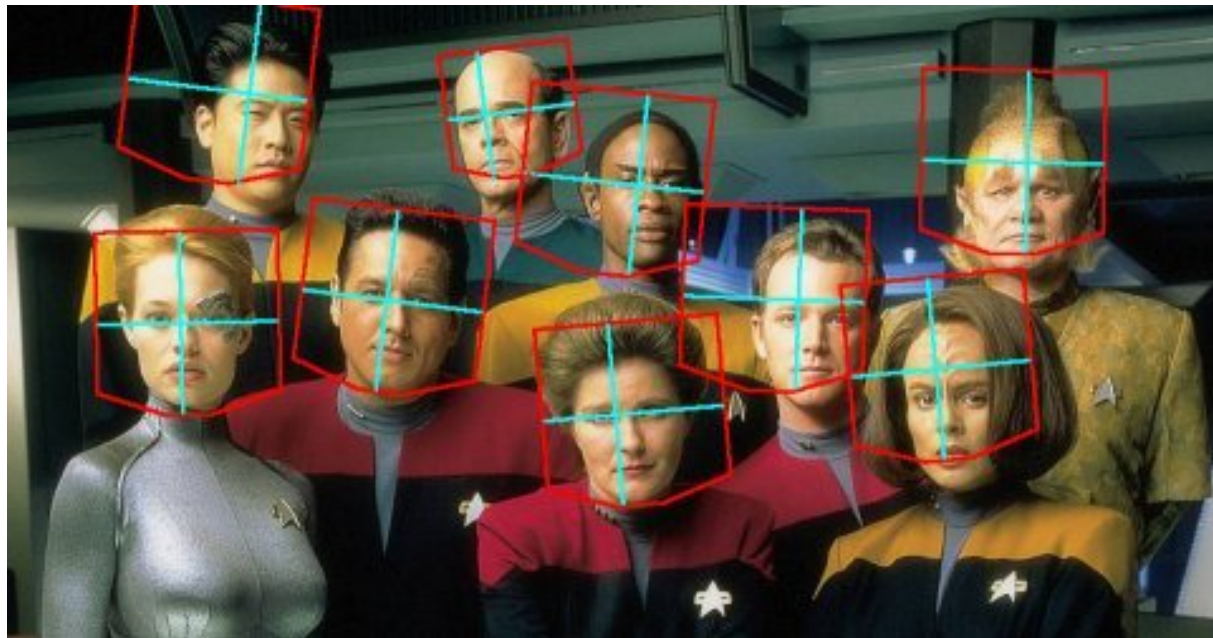
● 240x240 -> 788.0 million multiply-accumulate operations

● 480x480 -> 5,083 million multiply-accumulate operations

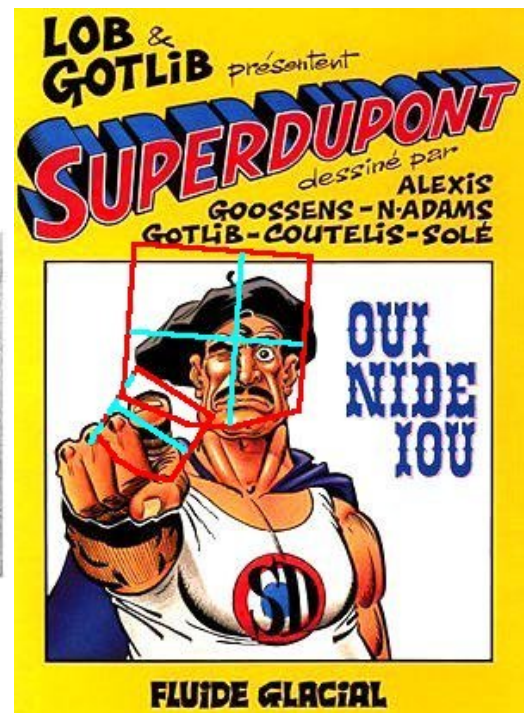
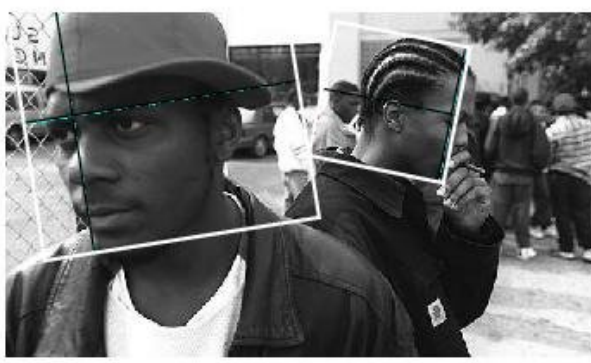
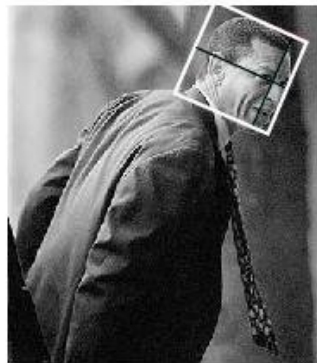
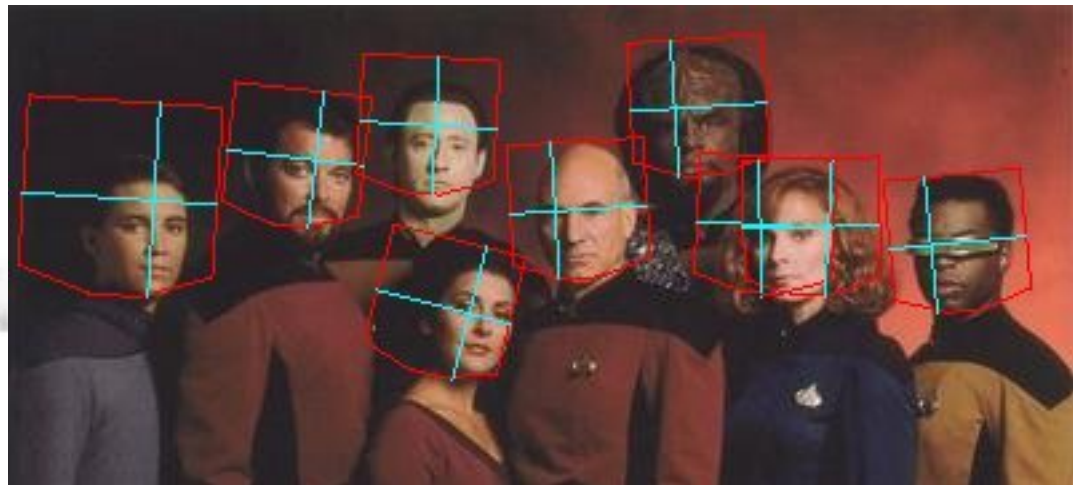


Face Detection: Results

<i>Data Set-></i>	TILTED		PROFILE		MIT+CMU	
	<i>False positives per image-></i>					
	4.42	26.9	0.47	3.36	0.5	1.28
Our Detector	90%	97%	67%	83%	83%	88%
Jones & Viola (tilted)	90%	95%	x		x	
Jones & Viola (profile)	x		70%	83%	x	



Face Detection and Pose Estimation: Results

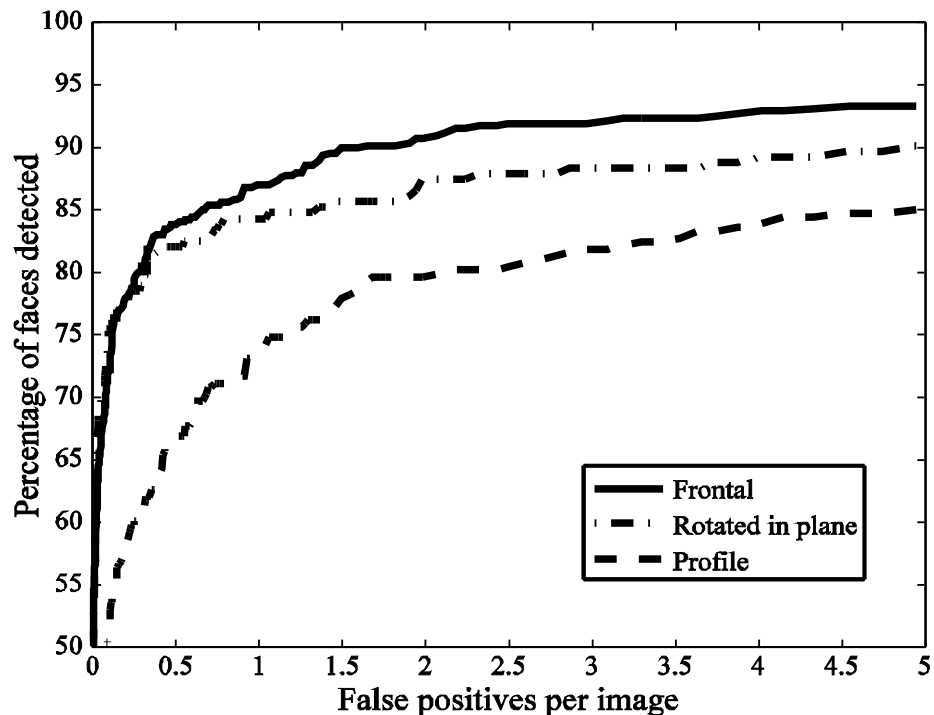


Face Detection with a Convolutional Net

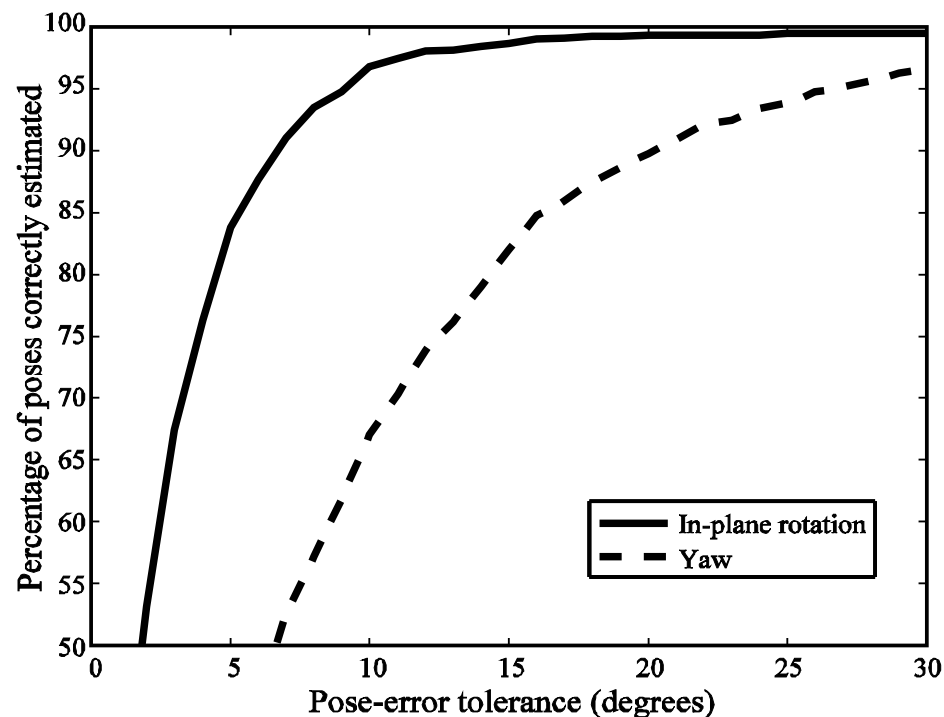


Performance on standard dataset

Detection



Pose estimation

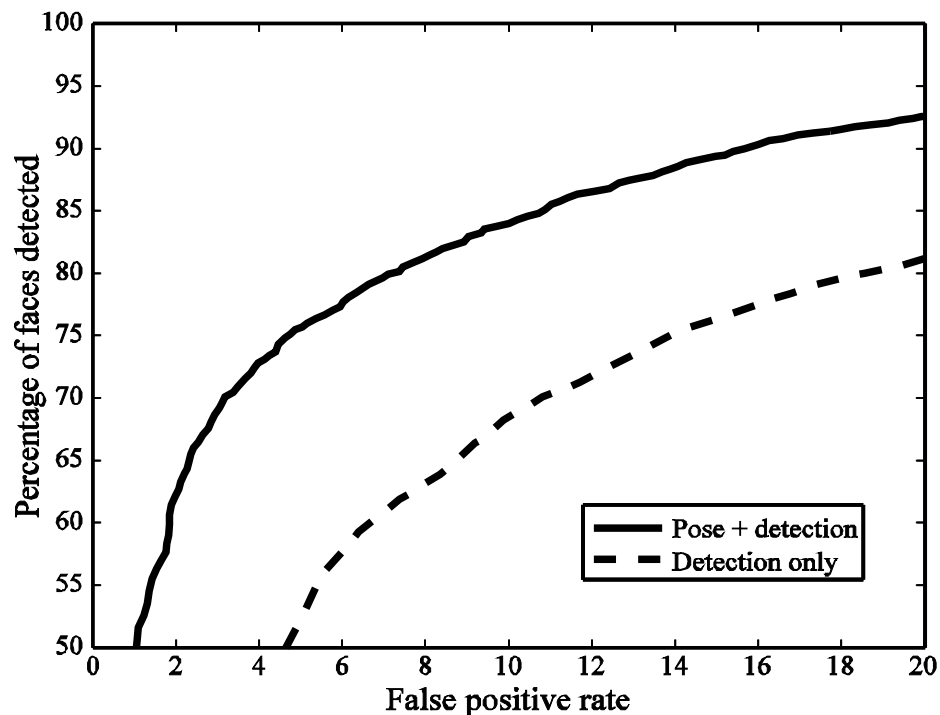


Pose estimation is performed on faces located automatically by the system

when the faces are localized by hand we get: 89% of yaw and 100% of in-plane rotations within 15 degrees.

Synergy Between Detection and Pose Estimation

**Pose Estimation Improves
Detection**



**Detection improves
pose estimation**

