

# Invariant Recognition

**Yann LeCun (Courant Institute, NYU)**

**Sumit Chopra, Raia Hadsell, Fu Jie Huang (Courant Institute, NYU)**

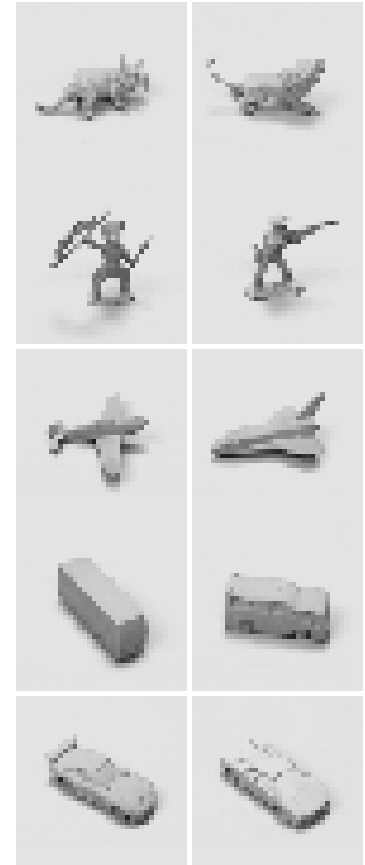
**Leon Bottou (NEC Labs)**

# Invariance

- The appearance of an object (in terms of pixels) changes considerably under changes of **pose, illumination, clutter, and occlusions**.
- Two instance of the same category may have widely differing shapes and appearances
  - ▶ An airliner and a fighter plane, a person standing and another one kneeling,...
- **Template-based methods are doomed** because the number of templates necessary to cover the space of variations grows **exponentially** with the number of dimensions of the variations.

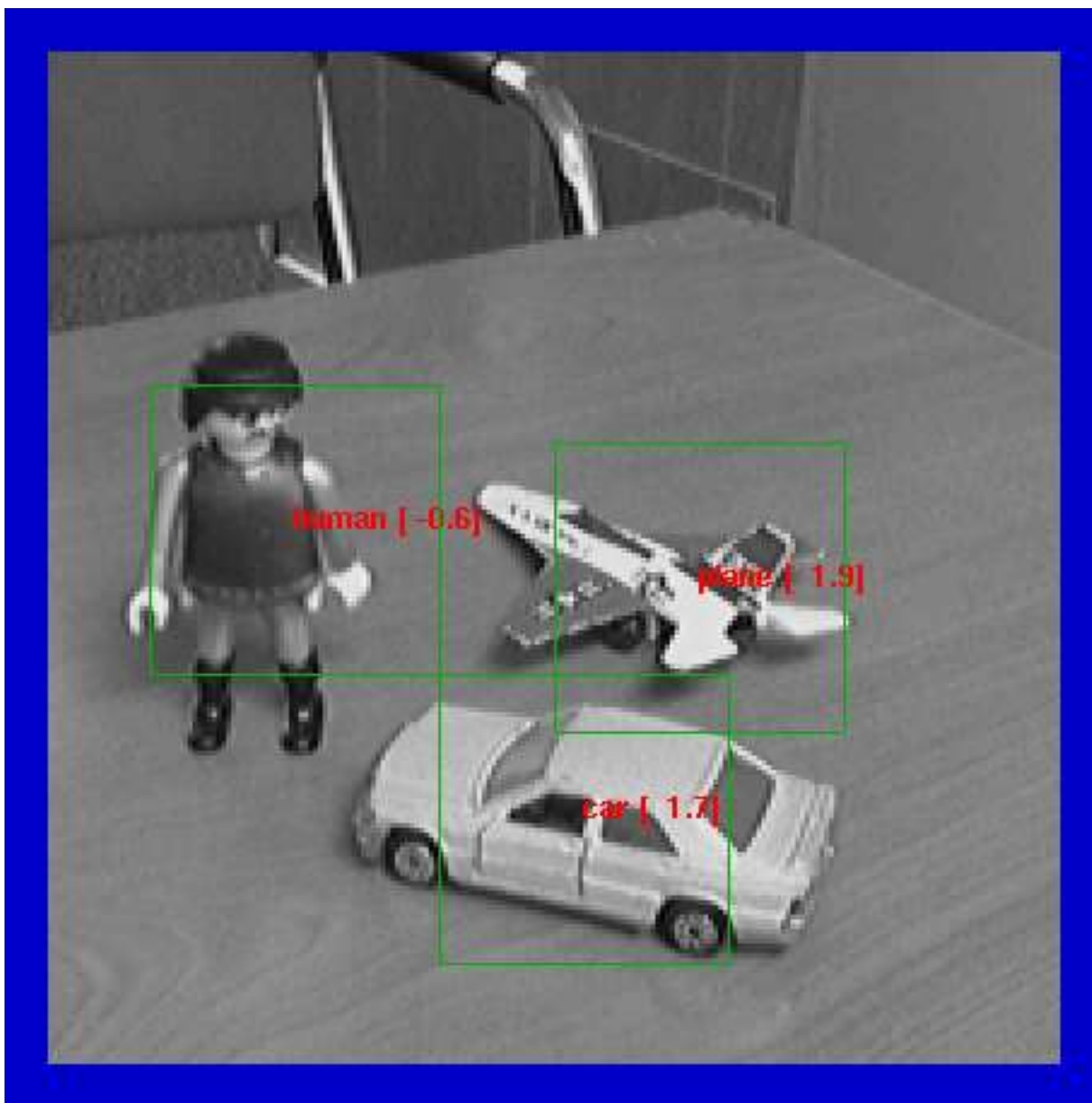
# Generic Object Recognition

- **Generic Object Recognition** is the problem of detecting and classifying objects into generic categories such as “cars”, “trucks”, “airplanes”, “animals”, or “human figures”
- **Appearances are highly variable within a category** because of shape variation, position in the visual field, scale, viewpoint, illumination, albedo, texture, background clutter, and occlusions.
- **Learning invariant representations is key.**
- **Understanding the neural mechanism behind invariant recognition is one of the main goals of Visual Neuroscience.**



# What we want to achieve

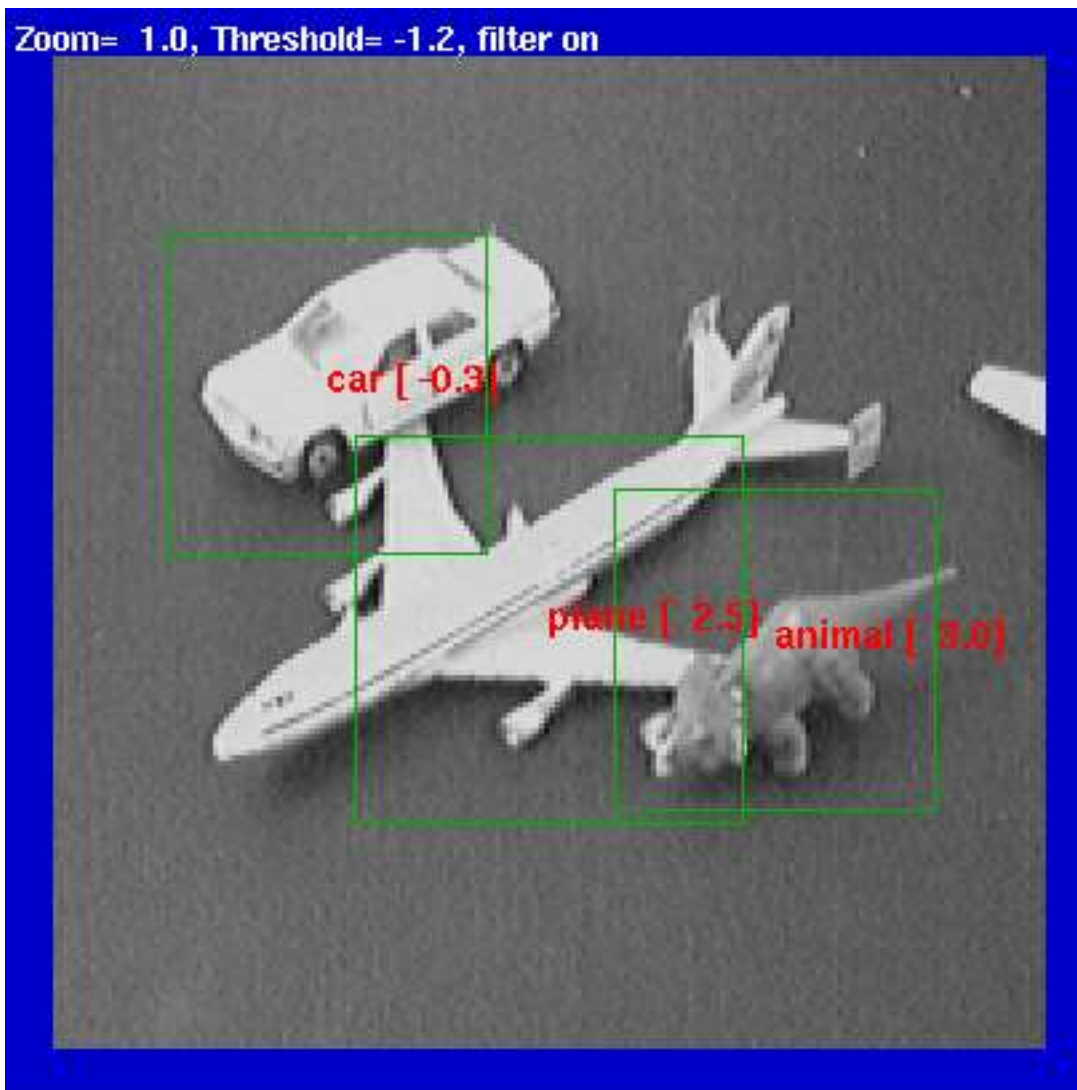
- **recognition from shape:**  
color, texture, and distinctive local features may be useful, but they merely allow us to sweep the real problems under the rug.
- **Full invariance to viewpoint, illumination, clutter, occlusions.**





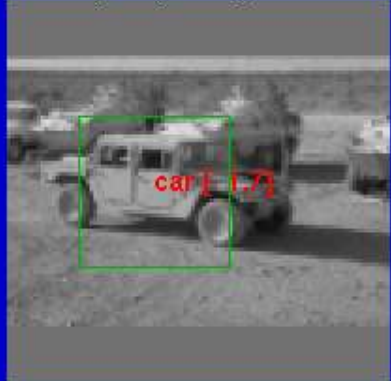
# Occlusions

Zoom= 1.0, Threshold= -1.2, filter on



# Clutter

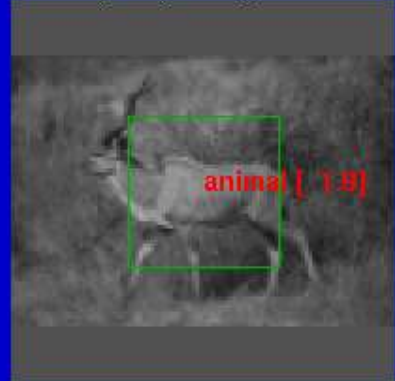
Thrs= 0.5, f on , os=40, nwin=23616



Thrs= 0.5, f on , os=40, nwin=23616



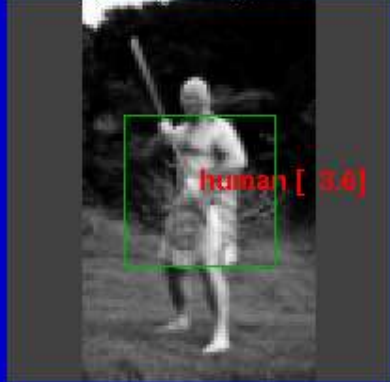
Thrs= 0.5, f on , os=40, nwin=23616



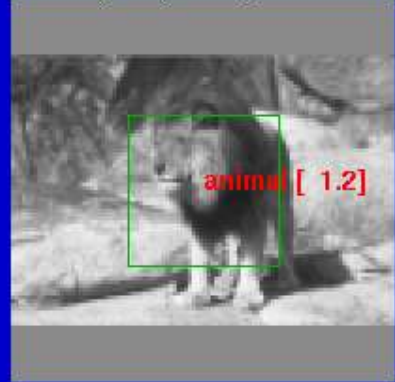
Thrs= 0.5, f on , os=40, nwin=23616



Thrs= 0.5, f on , os=40, nwin=23616



Thrs= 0.5, f on , os=40, nwin=23616



# The NYU Object Recognition Benchmark (NORB Dataset)

50 toys belonging to 5 categories: **animal**, **human figure**, **airplane**, **truck**, **car**

10 instance per category: **5 instances used for training**, **5 instances for testing**

**Raw dataset:** **972** stereo pair of each object instance. **48,600** image pairs total.

For each instance:

**18 azimuths**

0 to 350 degrees every 20 degrees

**9 elevations**

30 to 70 degrees from horizontal every 5 degrees

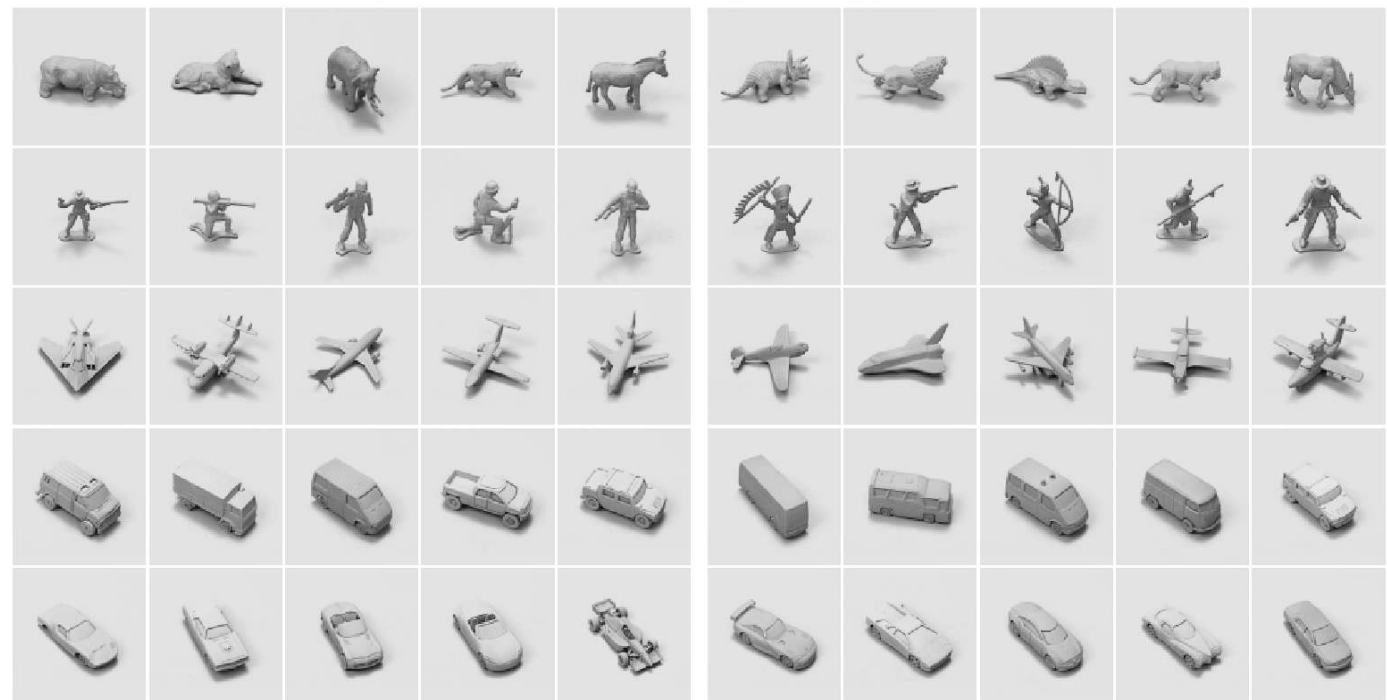
**6 illuminations**

on/off combinations of 4 lights

**2 cameras (stereo)**

7.5 cm apart

40 cm from the object



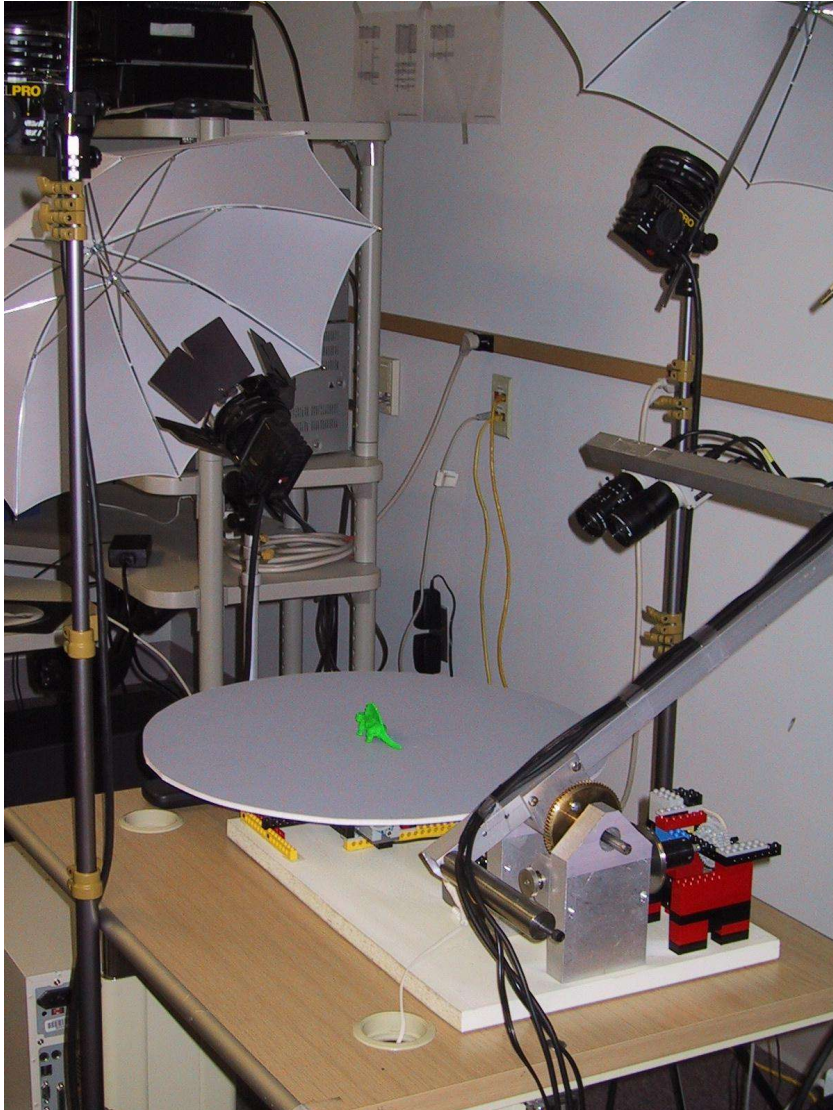
Training instances

Test instances



# Data Collection, Sample Generation

## Image capture setup

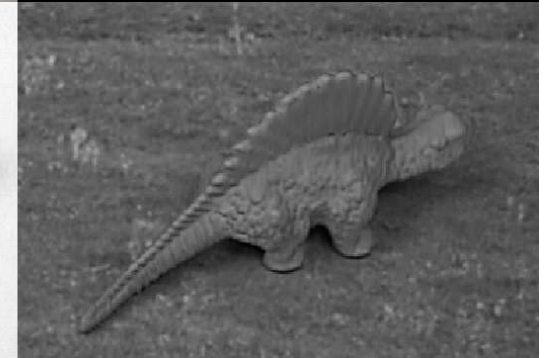
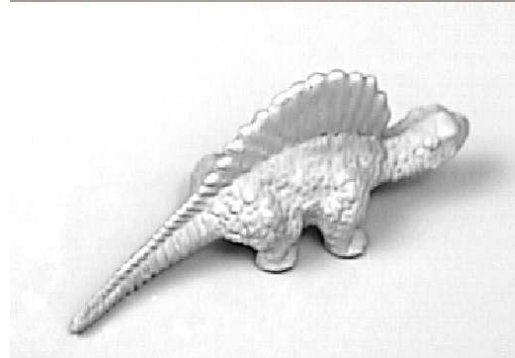
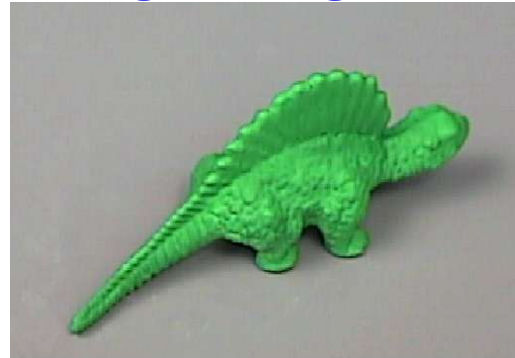


Objects are painted green so that:

- all features other than shape are removed
- objects can be segmented, transformed, and composited onto various backgrounds

**Original image**

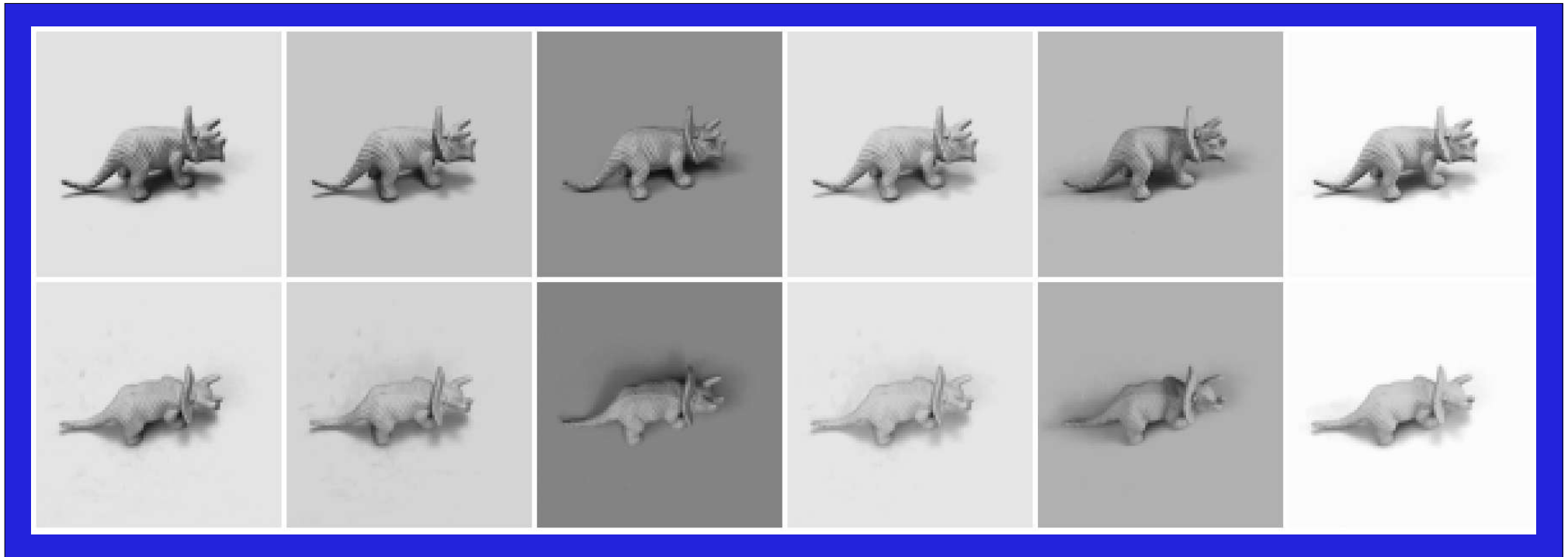
**Object mask**



**Shadow factor**

**Composite image**

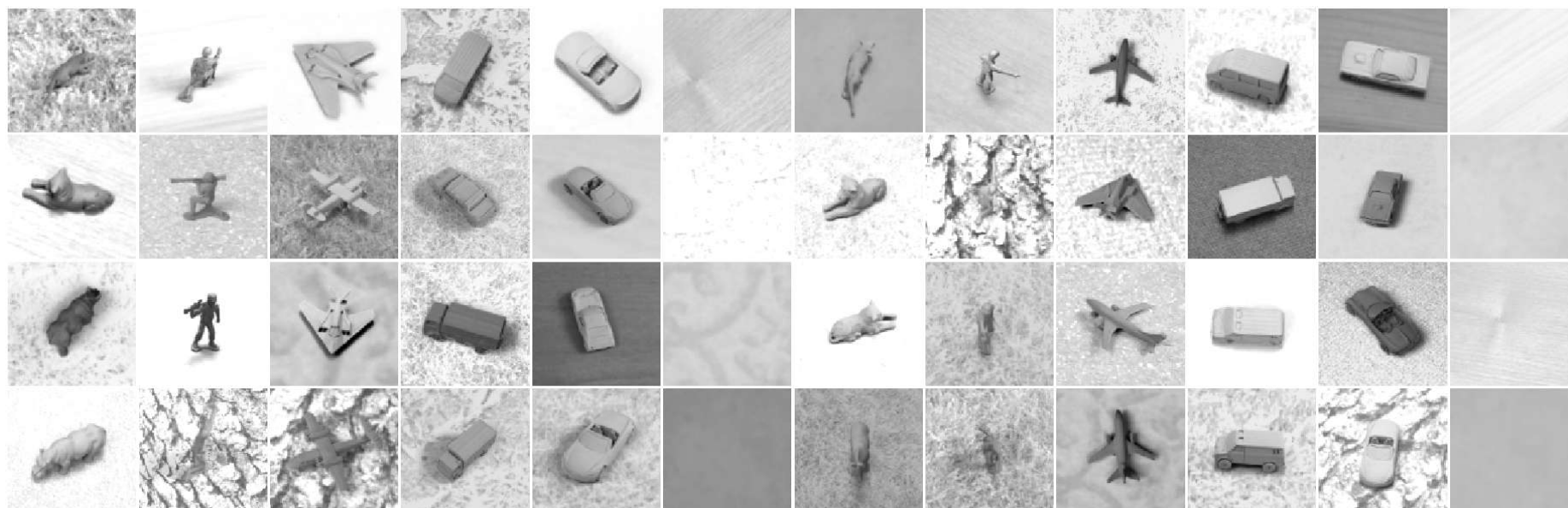
# Data Collection, Sample Generation



Samples showing the 6 different illuminations for 2 different elevations

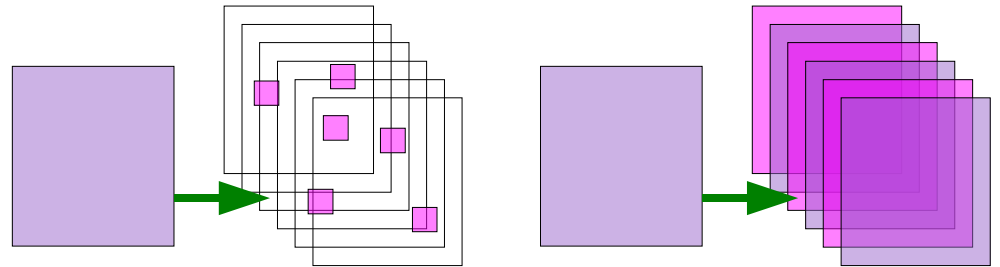


# Textured and Cluttered Datasets

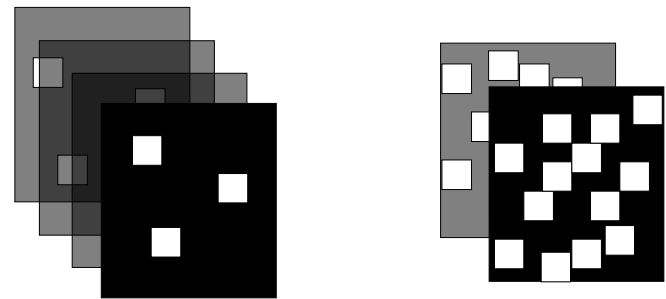


# Computational Models of Object Recognition

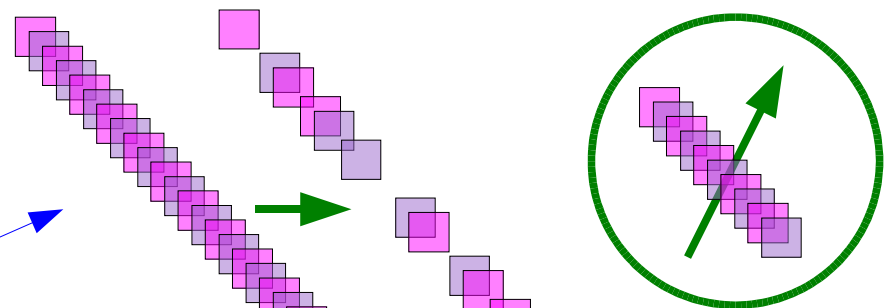
● Detecting features at interest points (Schmid, Perona, Ponce, Lowe) versus detecting them everywhere (LeCun, Ullman).



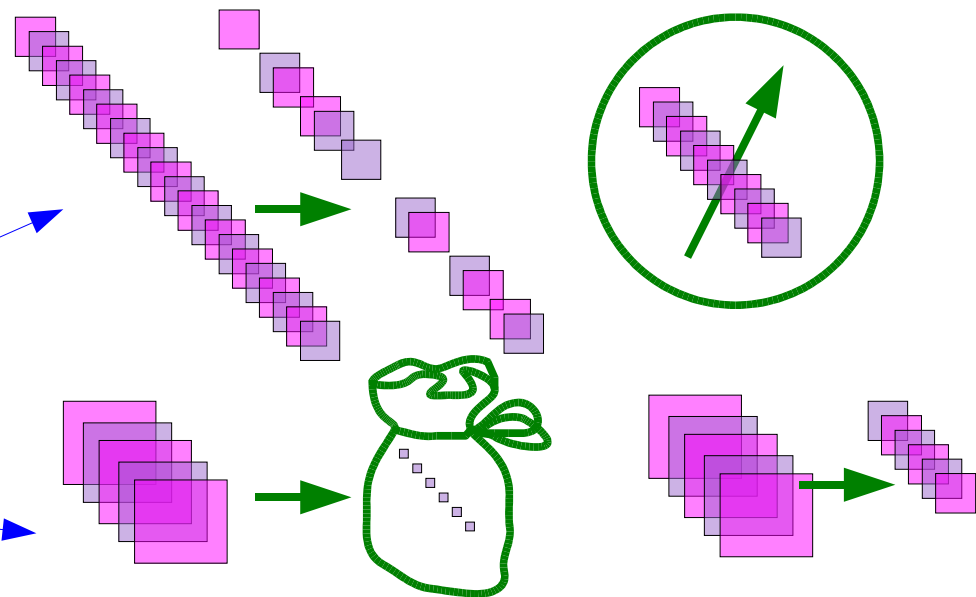
● Fixed features (Gabor, SIFT, Shape Context...), versus learned features



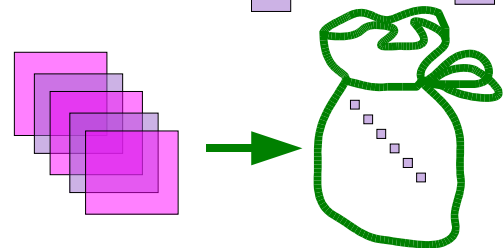
● Many sparse/selective features (Ullman's fragments) versus few dense/broad features (features that are "on" half the time).



● Selection from lots of simple features (Viola/Jones), vs tuning/optimization of a small number of features.

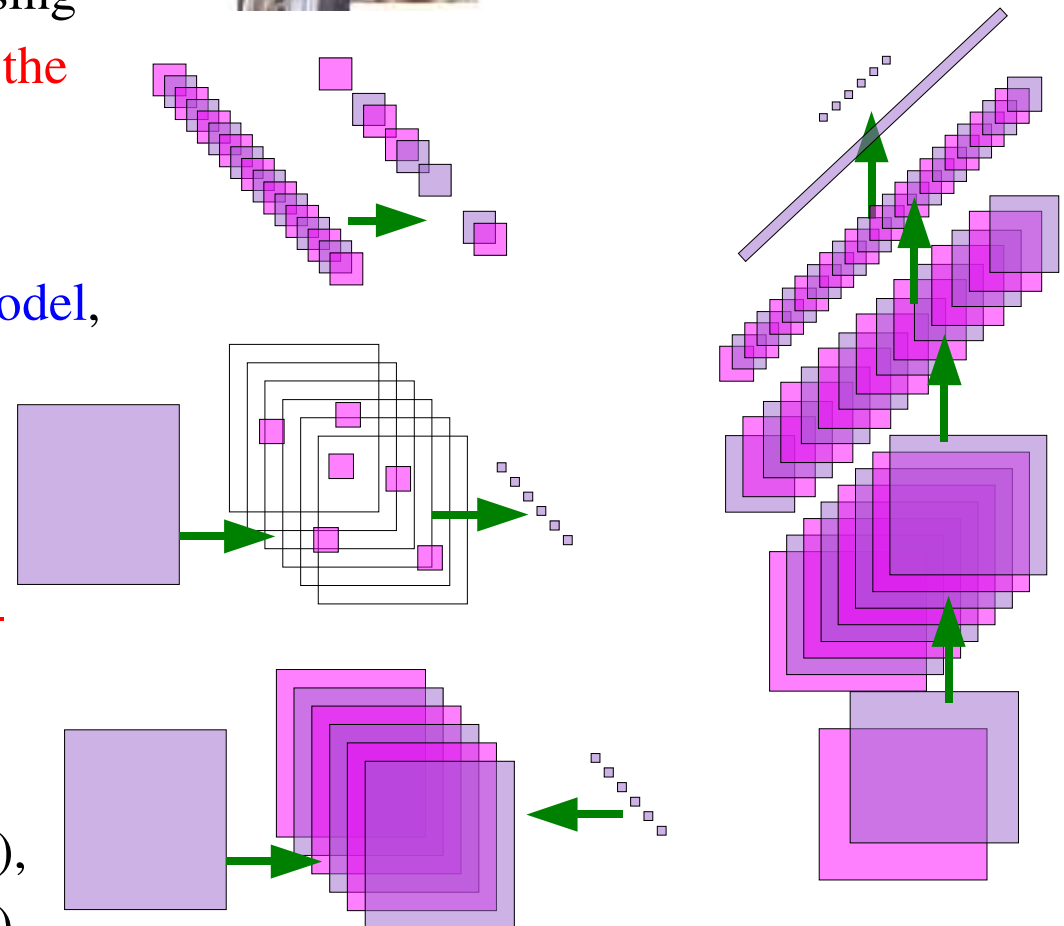
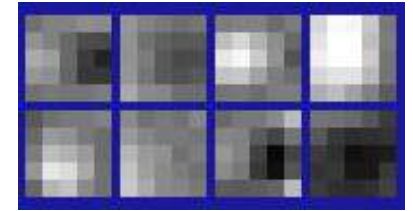


● Bag of features vs spatial relationships

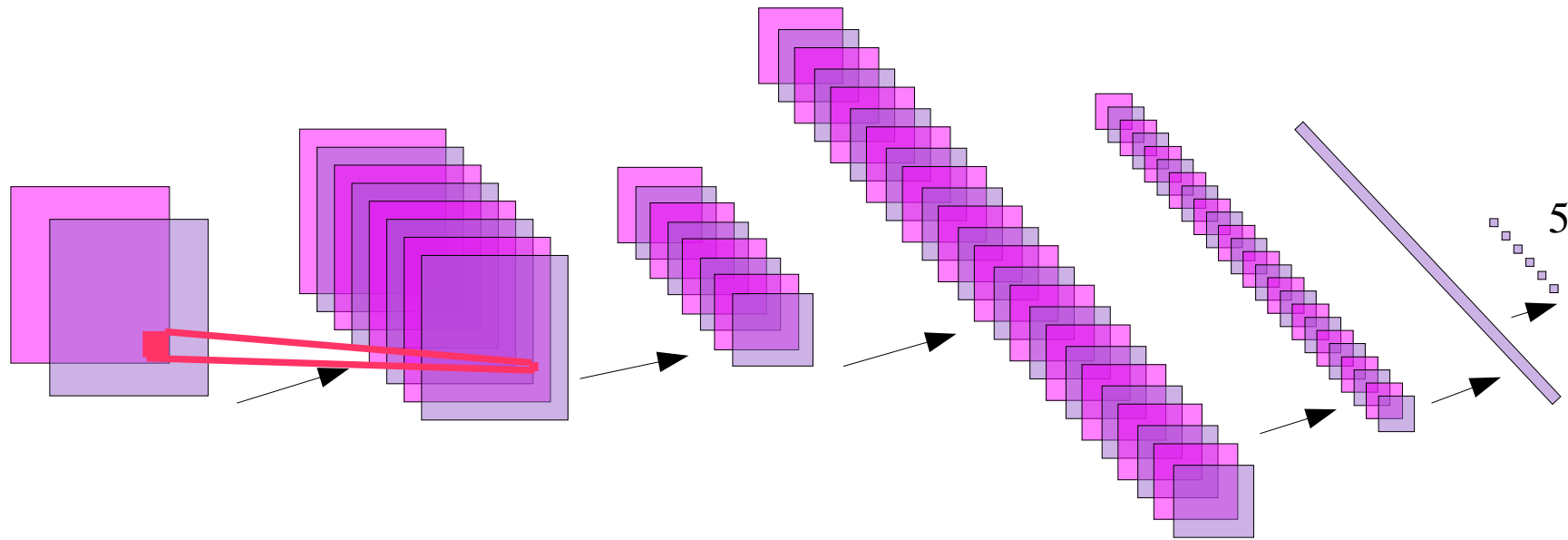


# What Architecture, what training?

- Selection of “patch” features (Schmid, Ullman, Ponce, Perona,.....), versus optimization of non-template features.
- “heuristic” feature selection (e.g. Using mutual information) versus learning the features by optimizing a global performance measure.
- Piecemeal training of feature and model, versus global training of the whole system
- 2-layer feature+model (almost everyone), versus hierarchical/multi-level (LeCun, Riesenhuber, Geman, Ullman)
- Generative (Perona, Amit, Freeman), versus discriminative (LeCun, Viola)



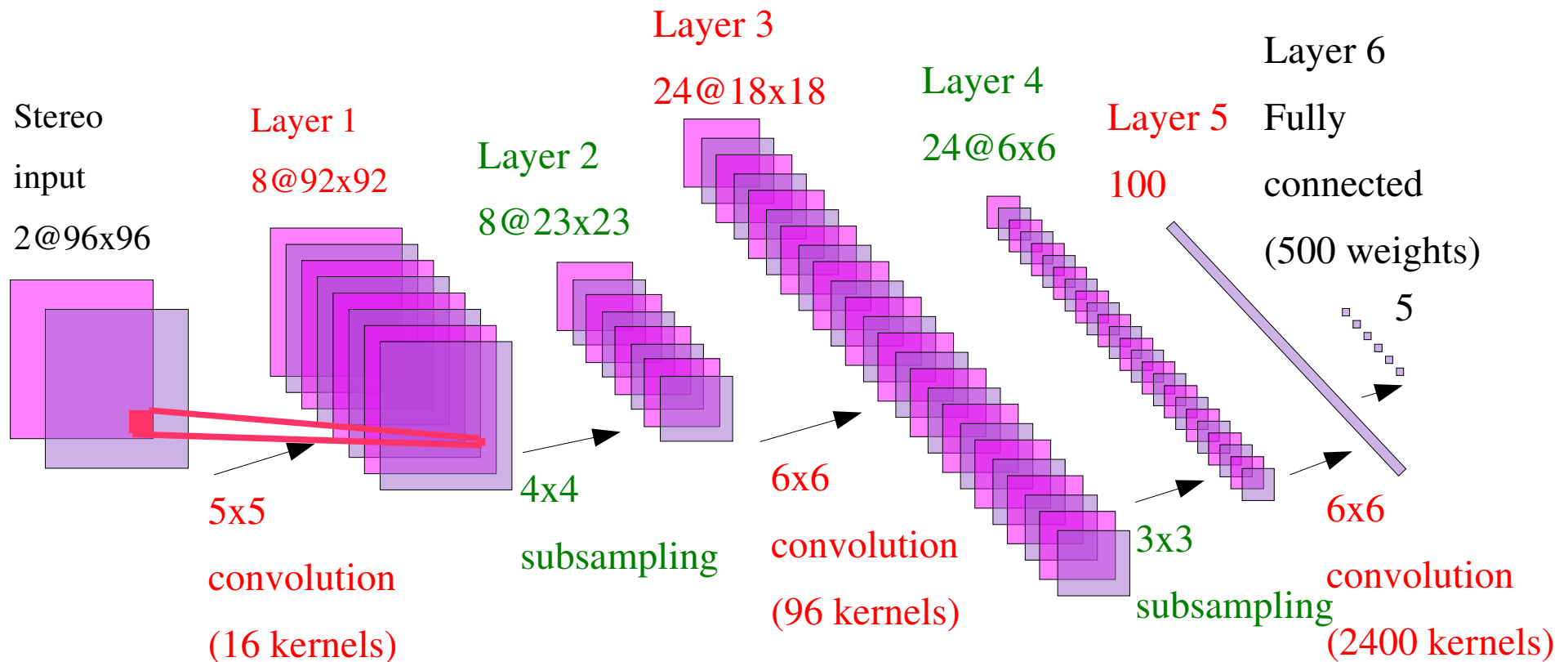
# Convolutional Network



- **Hierarchical/multilayer:** features get progressively more global, invariant, and numerous
- **dense features:** features detectors applied everywhere (no interest point)
- **broadly tuned (possibly invariant) features:** sigmoid units are on half the time.
- **Global discriminative training:** The whole system is trained “end-to-end” with a gradient-based method to minimize a global loss function
- **Integrates segmentation, feature extraction, and invariant classification in one fell swoop.**



# Convolutional Network



90,857 free parameters, 3,901,162 connections.

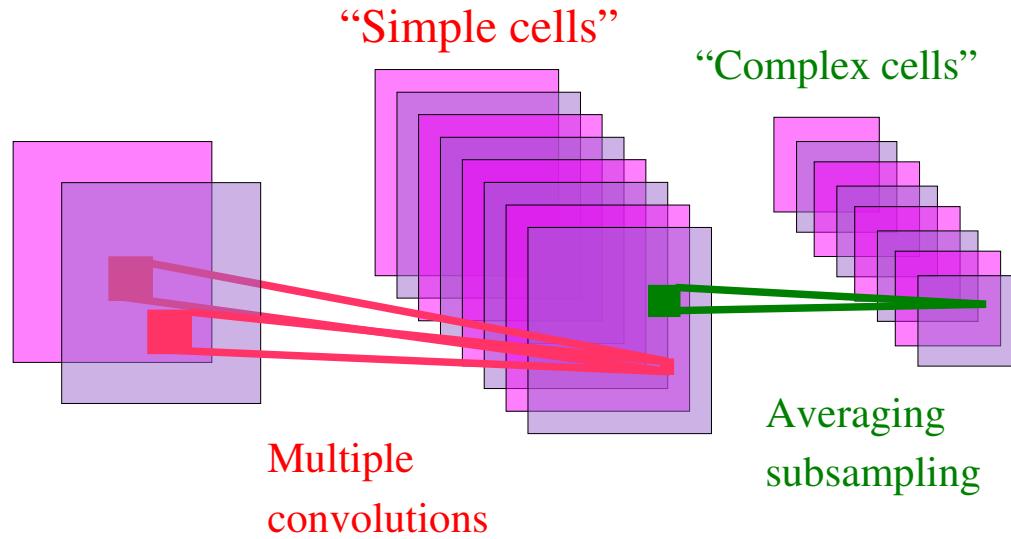
The architecture alternates **convolutional layers** (feature detectors) and **subsampling layers** (local feature pooling for invariance to small distortions).

**The entire network is trained end-to-end** (all the layers are trained simultaneously).

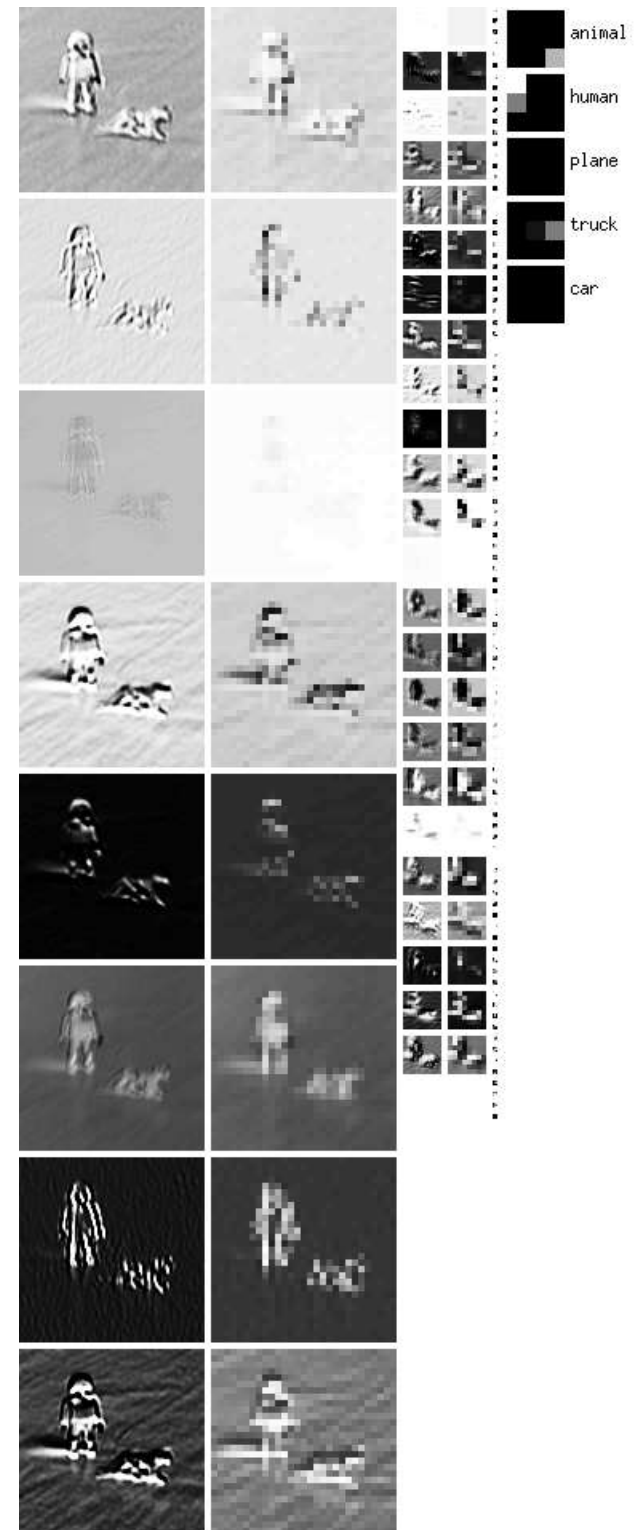
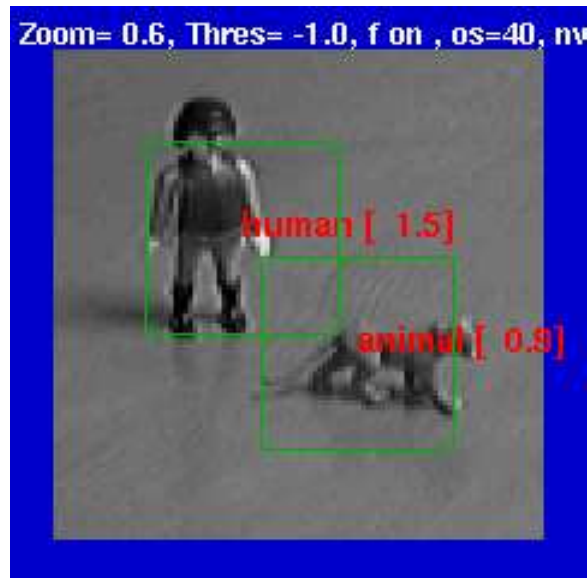
A gradient-based algorithm is used to minimize a supervised loss function.



# Alternated Convolutions and Subsampling



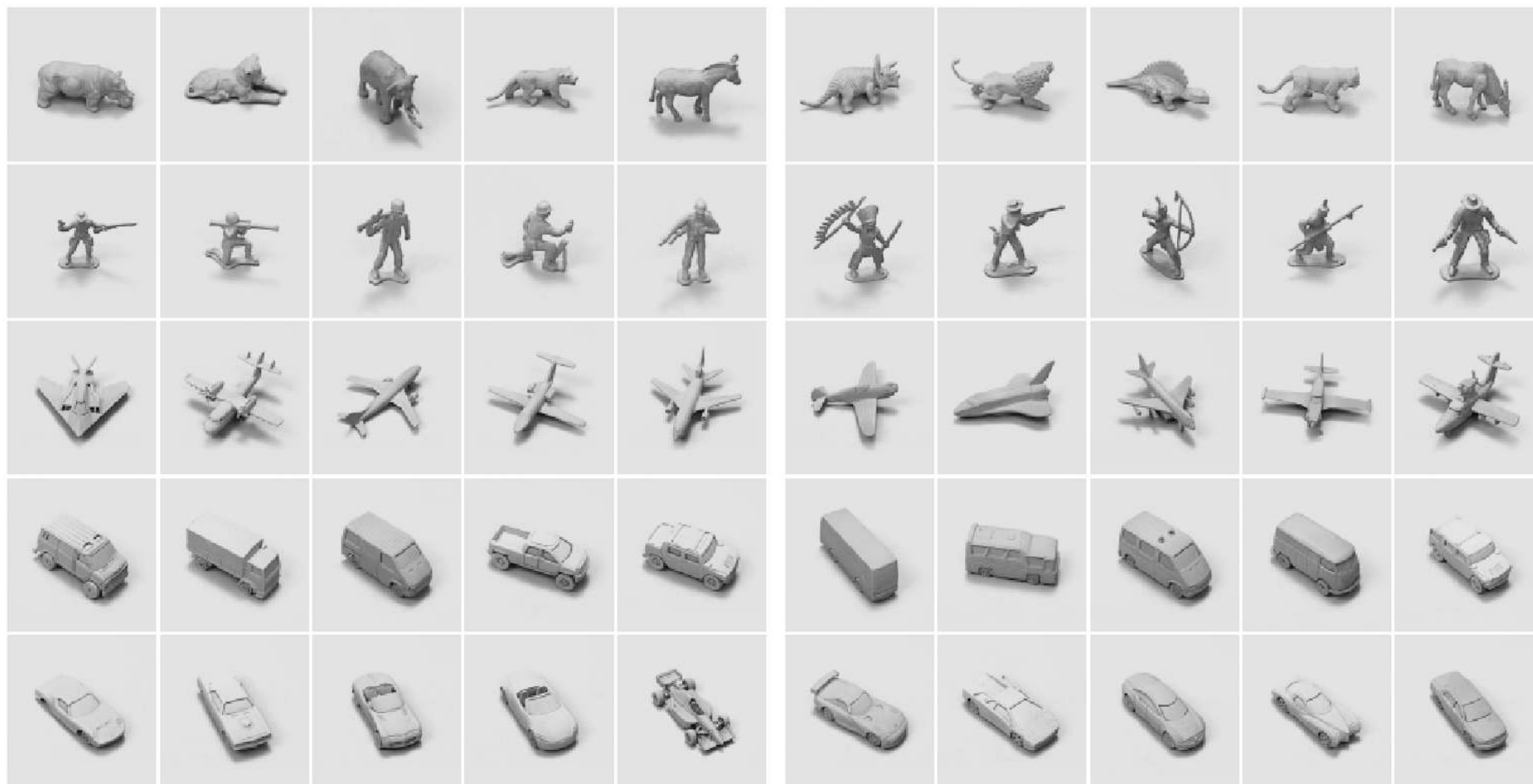
- Local features are extracted everywhere.
- averaging/subsampling layer builds robustness to variations in feature locations.
- Hubel/Wiesel'62, Fukushima'71, LeCun'89, Riesenhuber & Poggio'02, Ullman'02,....



# Experiment 1: Normalized-Uniform Dataset

- Normalized-Uniform Dataset: **972 stereo pair of each object instance** (18 azimuths X 9 elevations X 6 illuminations).
- 5 categories. 5 instances/category for training, 5 instances/category for testing**
- 24,300 stereo pairs for training, 24,300 for testing**
- Objects are centered and size-normalized so all the views of each object instance fits in an 80x80 pixel window.
- Objects are placed on uniform backgrounds (one for each of the 6 illuminations) of size 96x96 pixels
- Each sample is composed of two 96x96 images

# Experiment 1: Normalized-Uniform Dataset

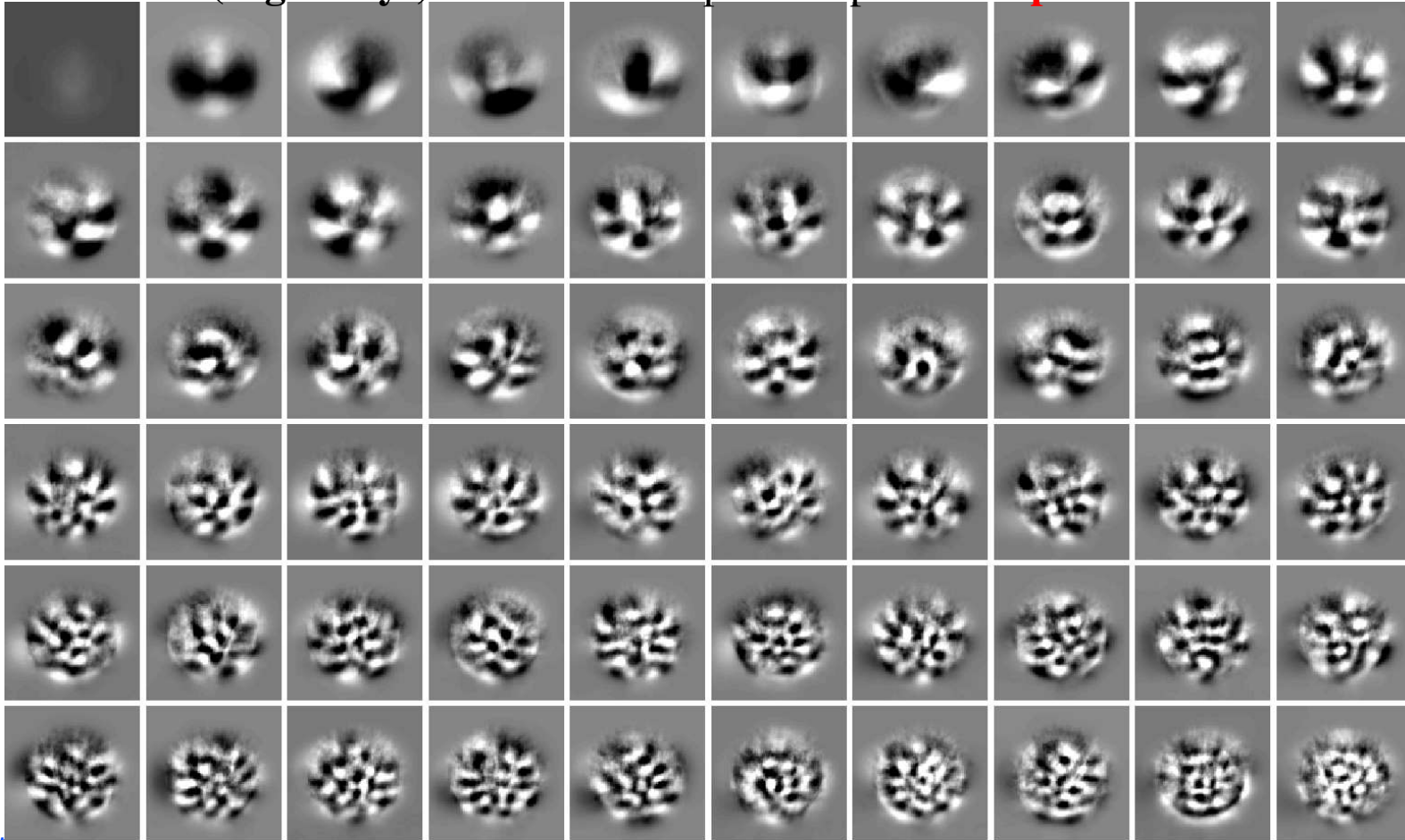


**Training instances**

**Test instances**

# Experiment 1: Normalized-Uniform Set: Representations

- 1 - Raw Stereo Input: 2 images 96x96 pixels **input dim. = 18432**
- 2 - Raw Monocular Input: 1 image, 96x96 pixels **input dim. = 9216**
- 3 - Subsampled Mono Input: 1 image, 32x32 pixels **input dim = 1024**
- 4 - PCA-95 (EigenToys): First 95 Principal Components **input dim. = 95**



First 60 eigenvectors (EigenToys)

## Experiment 1: Normalized-Uniform Set: Error Rates

- Linear Classifier on raw stereo images: **30.2% error.**
- K-Nearest-Neighbors on raw stereo images: **18.4% error.**
- K-Nearest-Neighbors on PCA-95: **16.6% error.**
- Pairwise SVM on 96x96 stereo images: **14.1% error**
- Pairwise SVM on 48x48 stereo images: **12.5% error**
- Pairwise SVM on 32x32 stereo images: **11.8% error.**
- Pairwise SVM on 48x48 monocular images: **13.9% error.**
- Pairwise SVM on 32x32 monocular images: **12.6% error.**
- Pairwise SVM on 95 Principal Components: **13.3% error.**
- Convolutional Net on 32x32 stereo images: **11.3% error.**
- Convolutional Net on 48x48 stereo images: **8.7% error.**
- Convolutional Net on 96x96 stereo images: **6.6% error.**

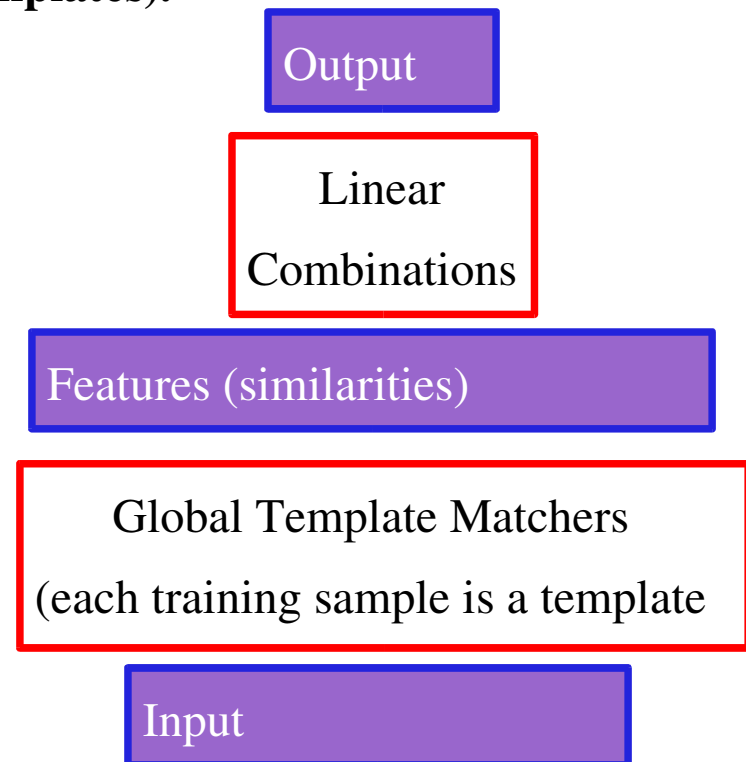


# What's wrong with K-NN and SVMs?

- K-NN and SVM with Gaussian kernels are based on **matching global templates**
- Both are “shallow” architectures
- There is now way to learn invariant recognition tasks with such naïve architectures (unless we use an impractically large number of templates).

● The number of necessary templates grows **exponentially** with the number of dimensions of variations.

● Global templates are in trouble when the variations include: category, instance shape, configuration (for articulated object), position, azimuth, elevation, scale, illumination, texture, albedo, in-plane rotation, background luminance, background texture, background clutter, .....

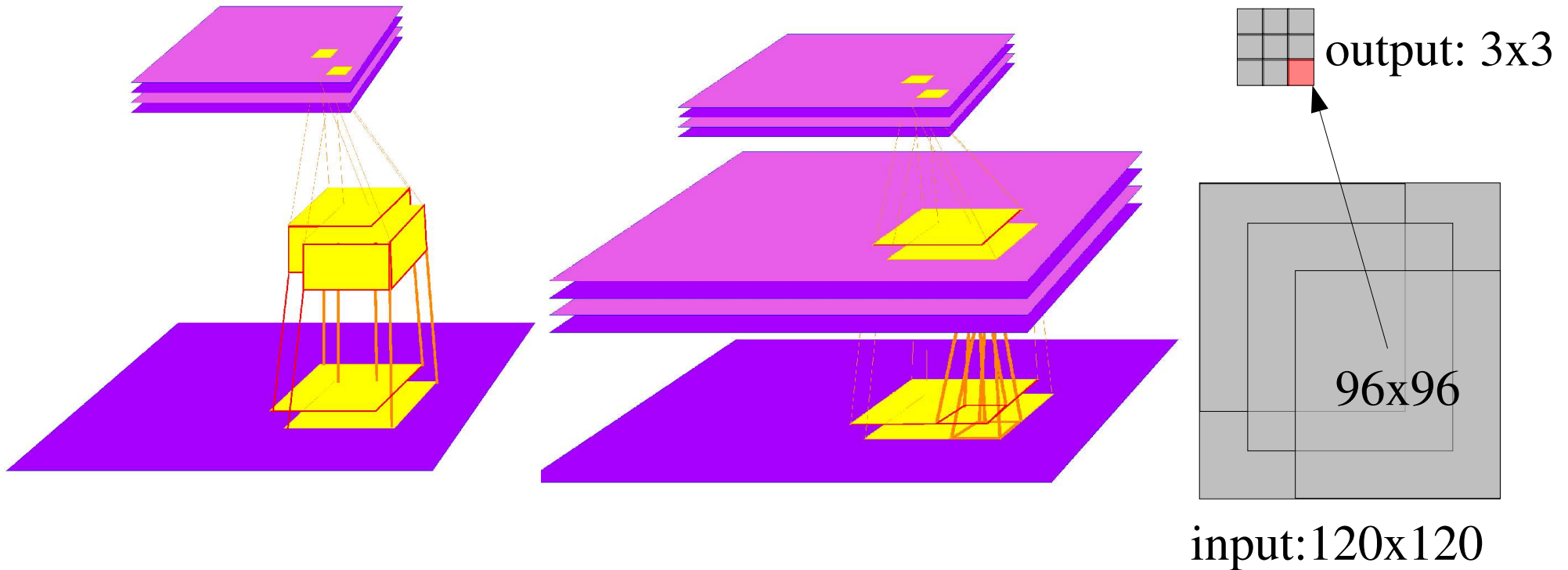


## Experiment 2: Jittered-Cluttered Dataset



- 291,600 training samples, 58,320 test samples
- Convolutional Net with **binocular** input: **7.8% error**
- Convolutional Net + SVM on top: **5.8% error**
- Convolutional Net with **monocular** input: **20.8% error**
- Smaller **mono** net (DEMO): **26.0% error**
- Dataset available from <http://www.cs.nyu.edu/~yann>

# 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 1.5.

# 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

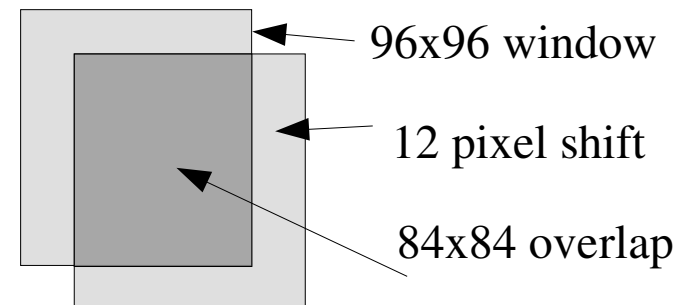
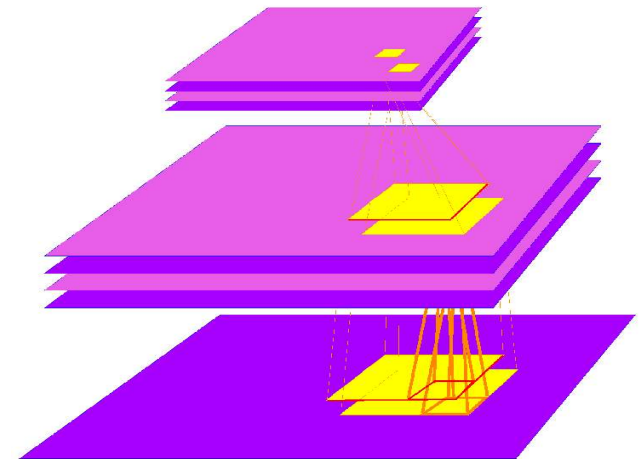
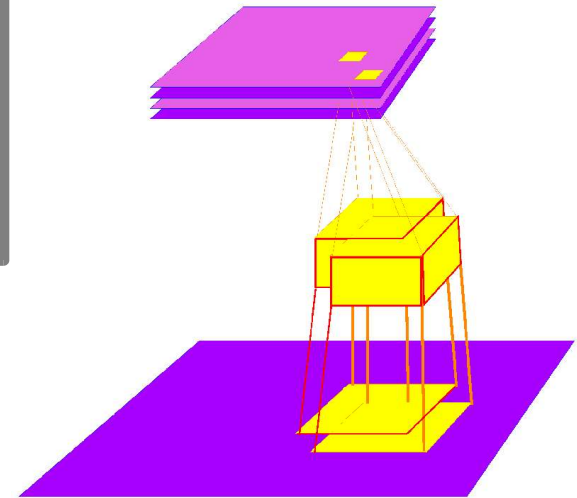
● Computational cost for a non-convolutional detector of the same size, applied every 12 pixels:

● 96x96 -> 4.6 million multiply-accumulate operations

● 120x120 -> 42.0 million multiply-accumulate operations

● 240x240 -> 788.0 million multiply-accumulate operations

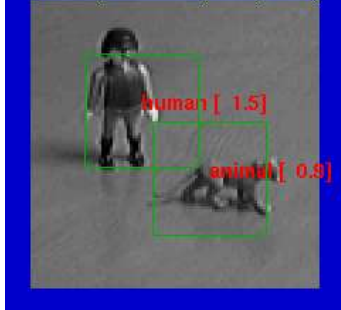
● 480x480 -> 5,083 million multiply-accumulate operations



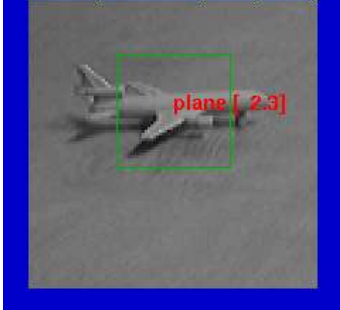


# Examples (Monocular Mode)

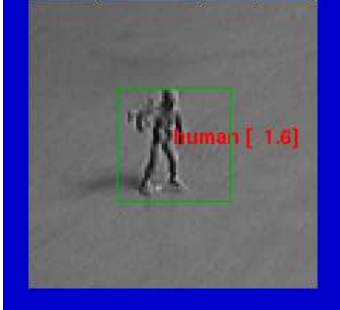
Zoom= 0.6, Thres= -1.0, f on , os=40, nv



Zoom= 0.6, Thres= -1.0, f on , os=40, nv



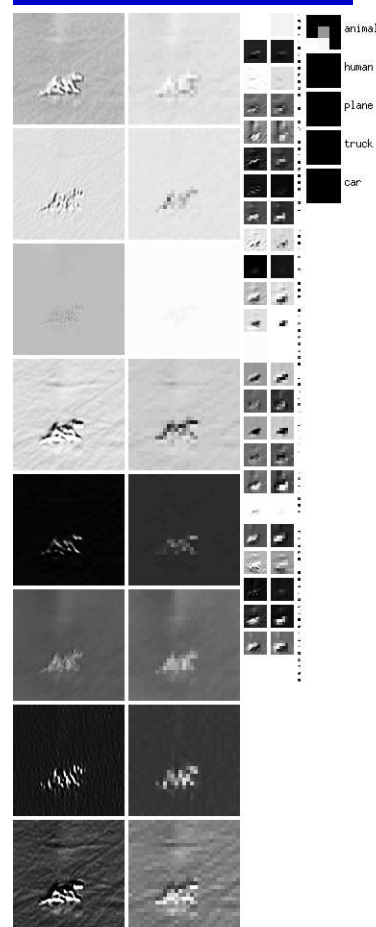
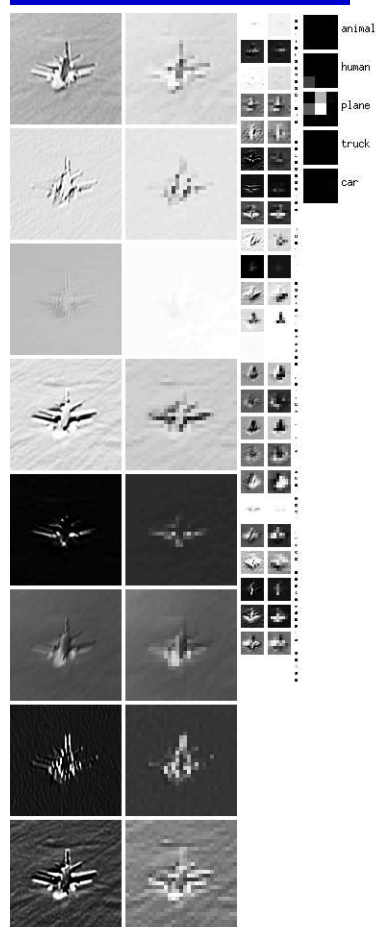
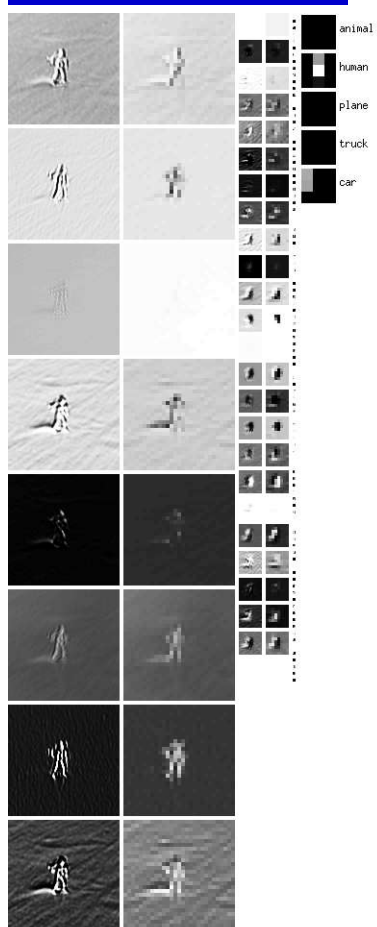
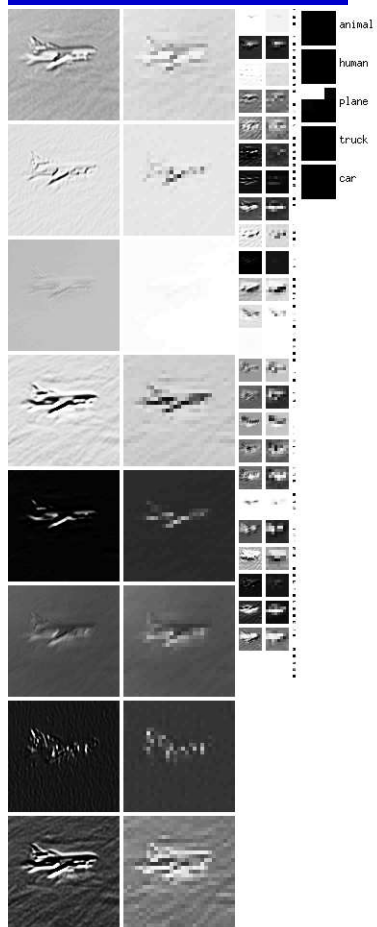
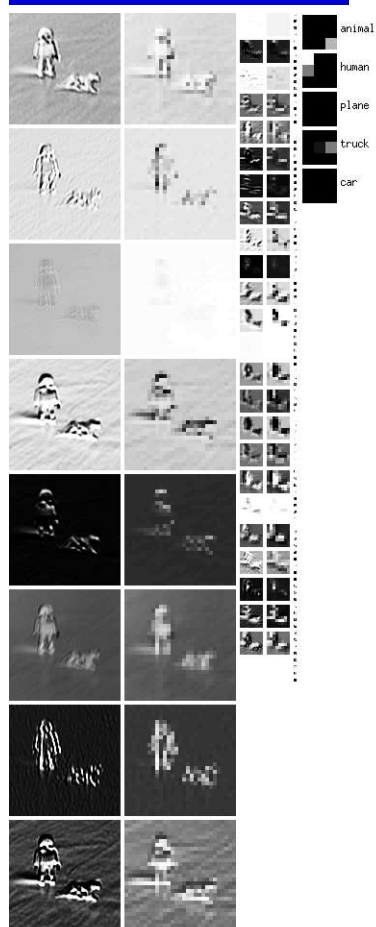
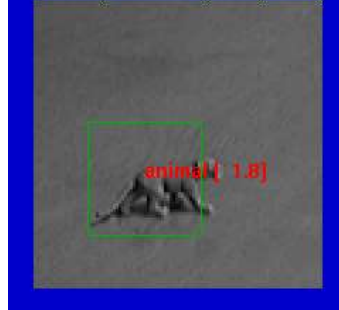
Zoom= 0.6, Thres= -1.0, f on , os=40, nv



Zoom= 0.6, Thres= -1.0, f on , os=40, nv



Zoom= 0.6, Thres= 0.5, f on , os=40, nv

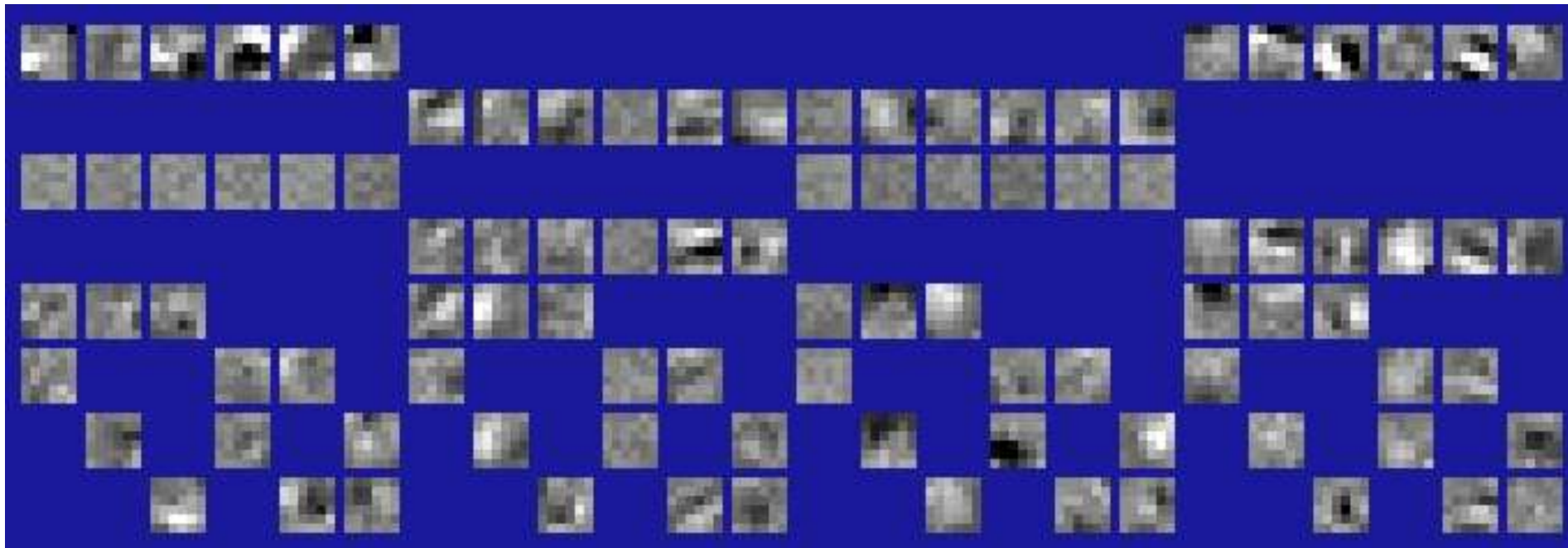




# Learned Features

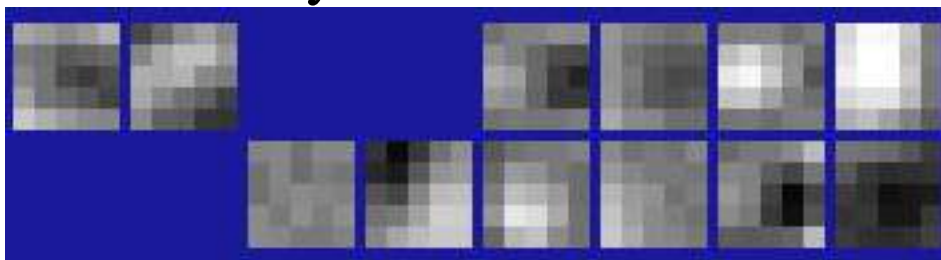
Layer 3

Layer 2

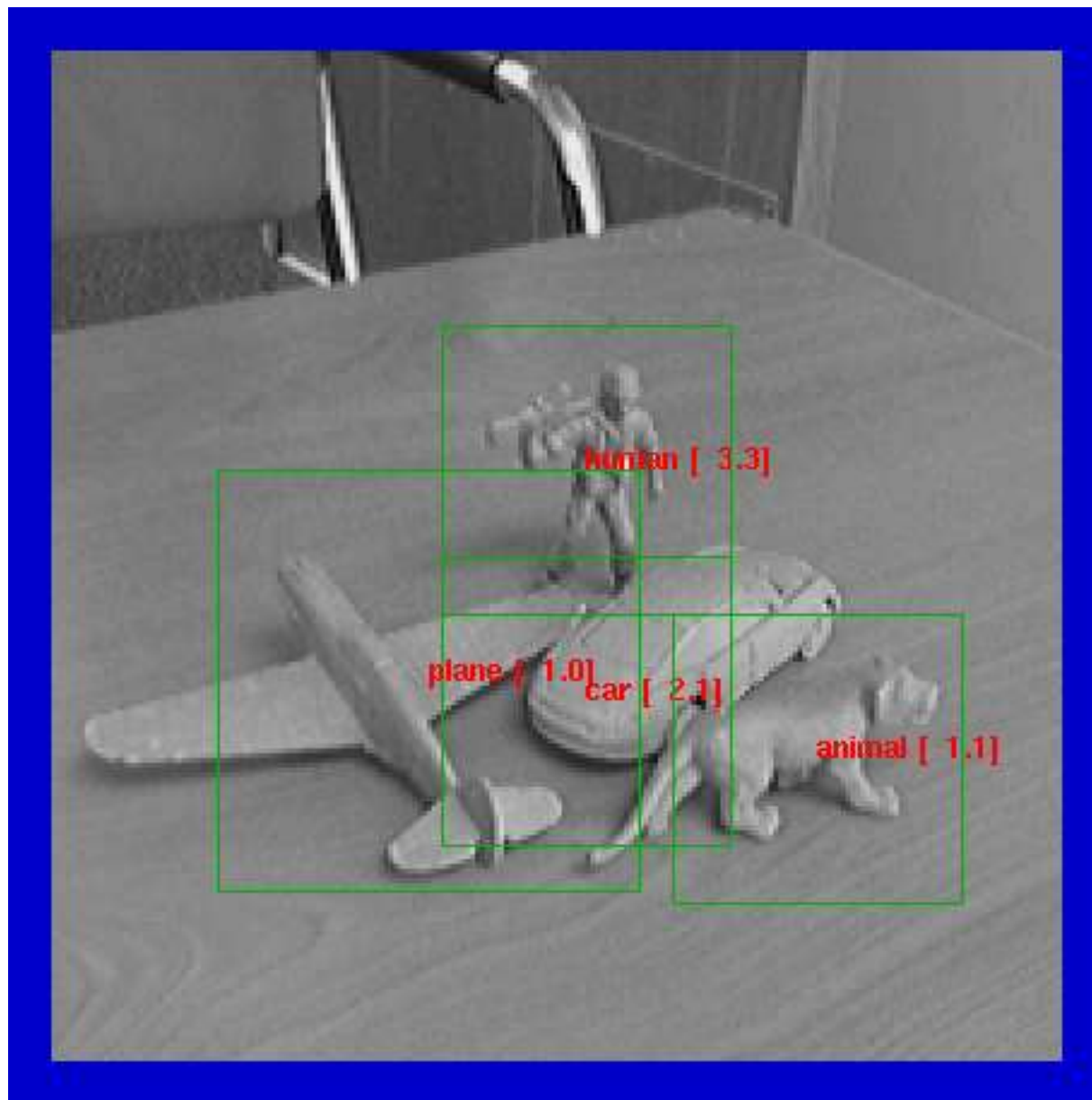


Layer 1

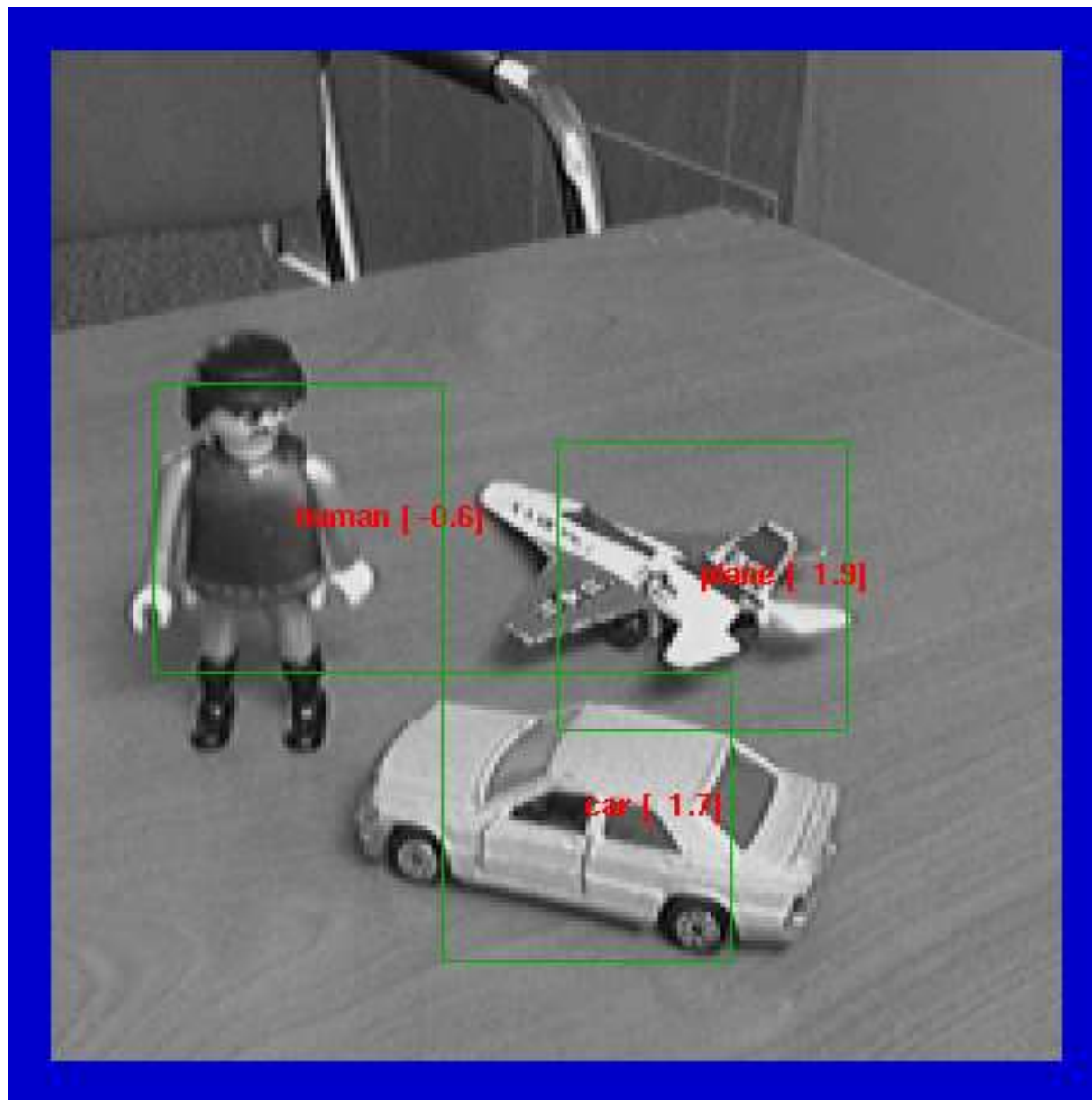
Input



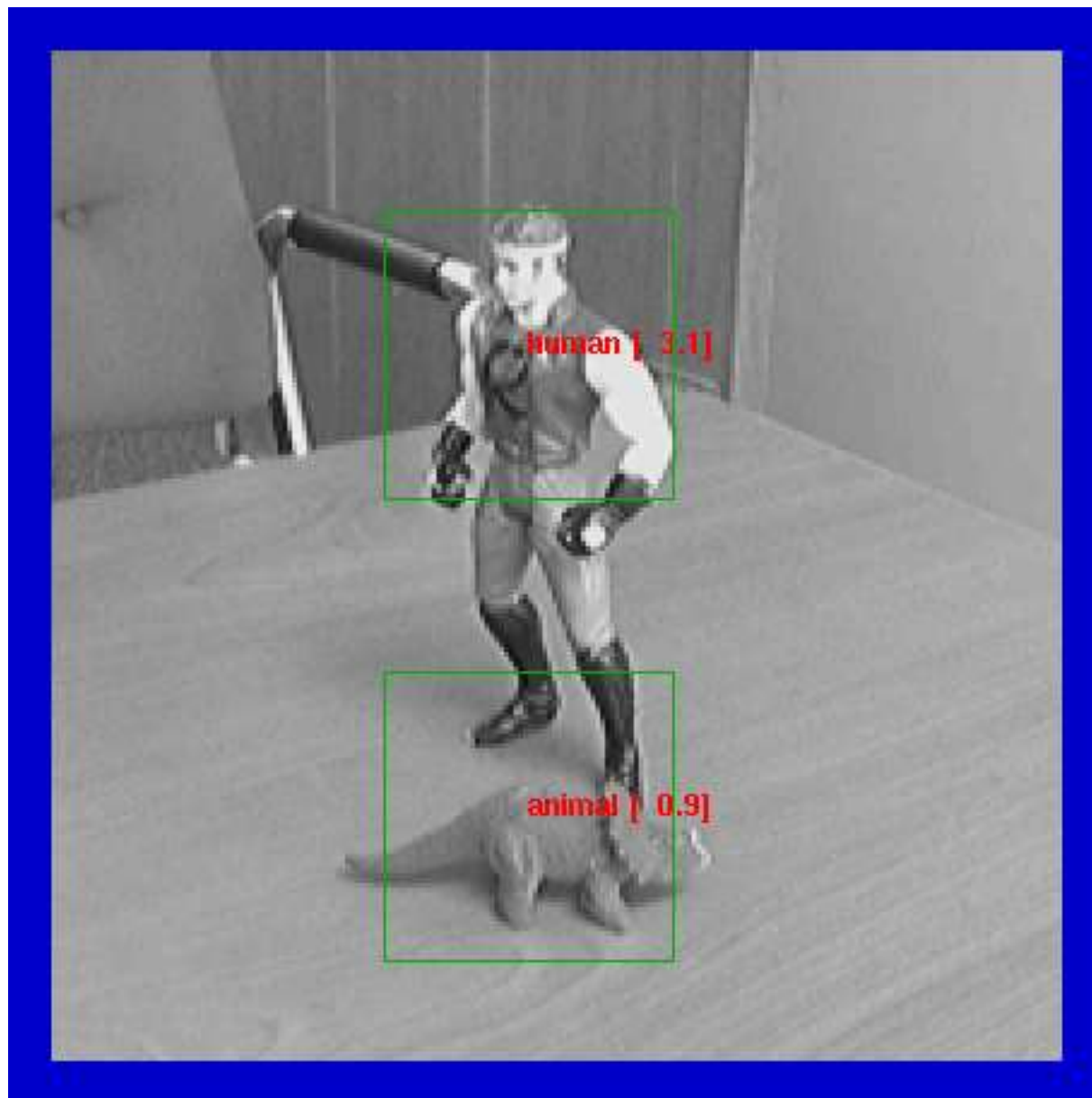
# Examples (Monocular Mode)



# Examples (Monocular Mode)



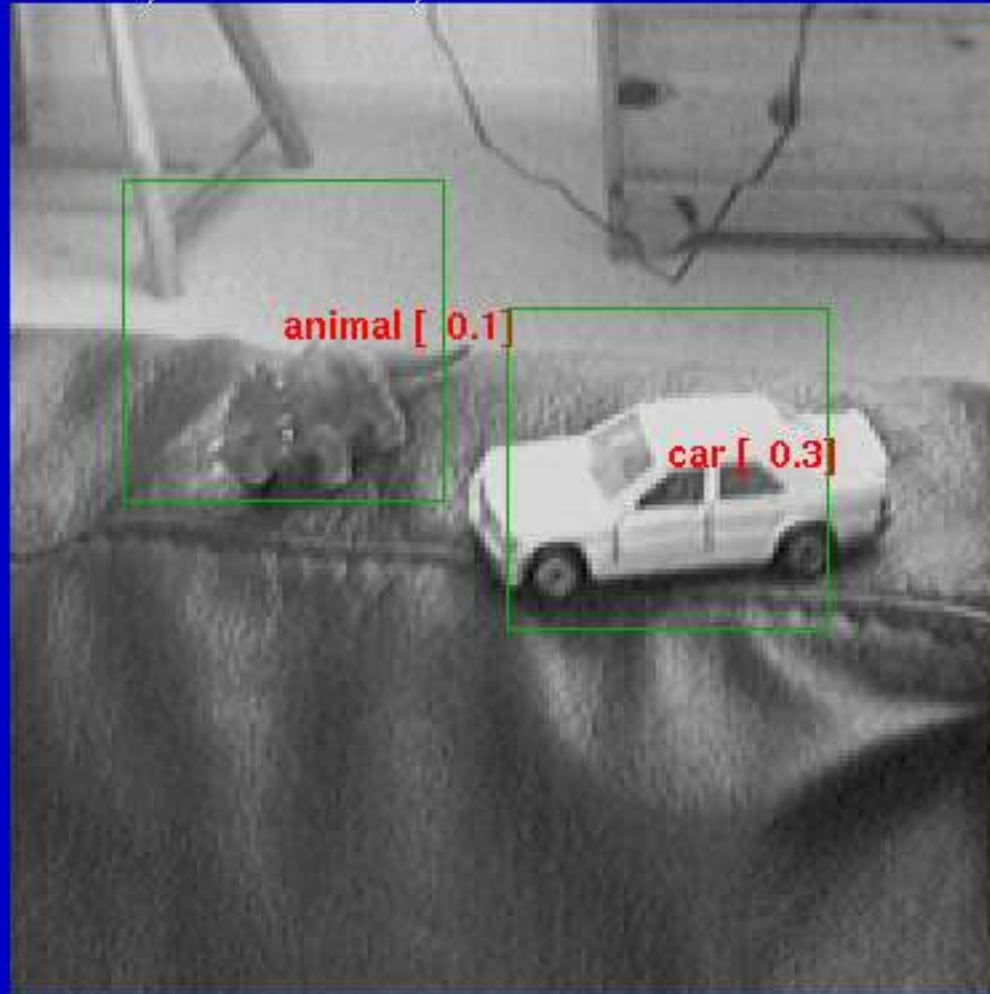
# Examples (Monocular Mode)





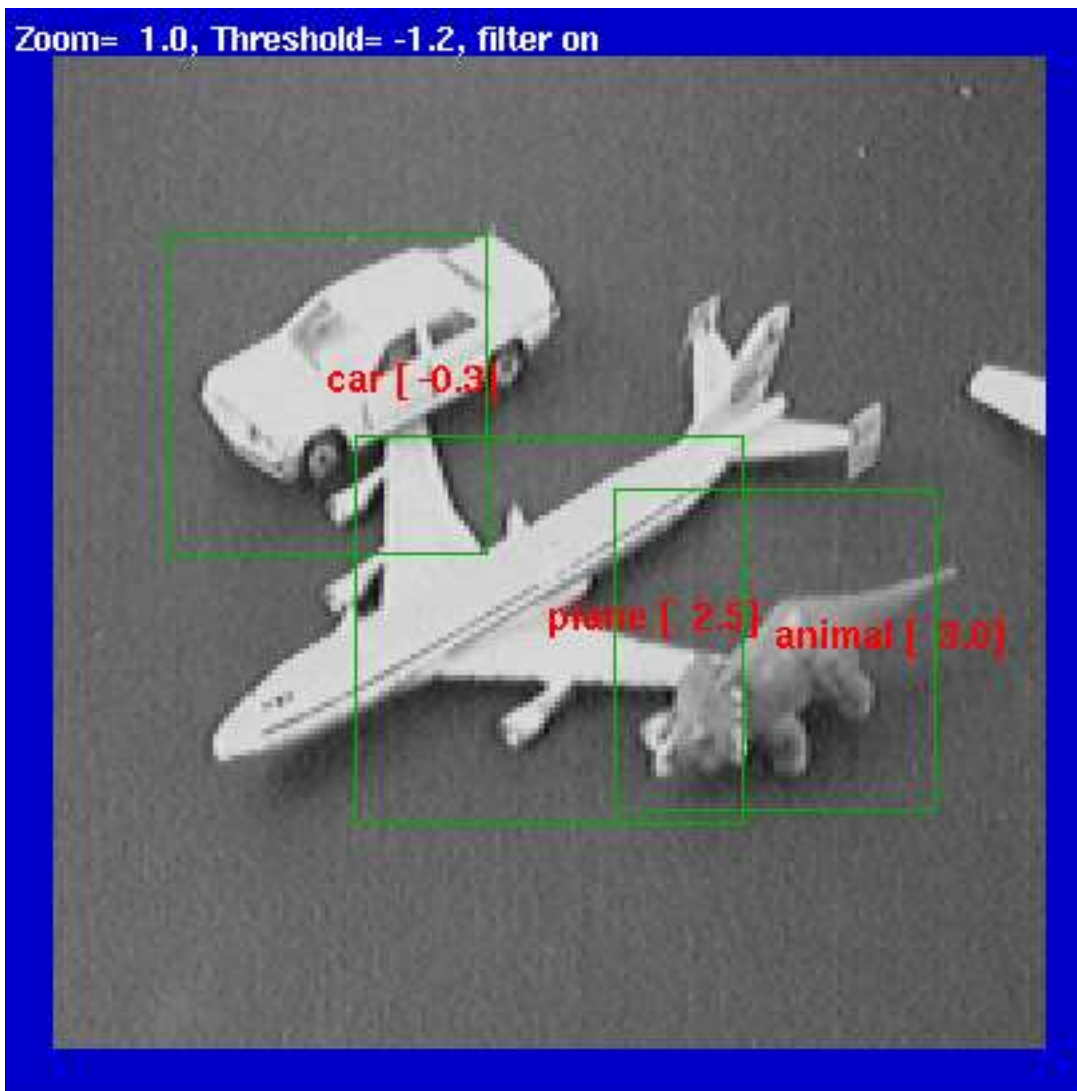
# Examples (Monocular Mode)

Zoom= 1.0, Threshold= -1.0, filter on



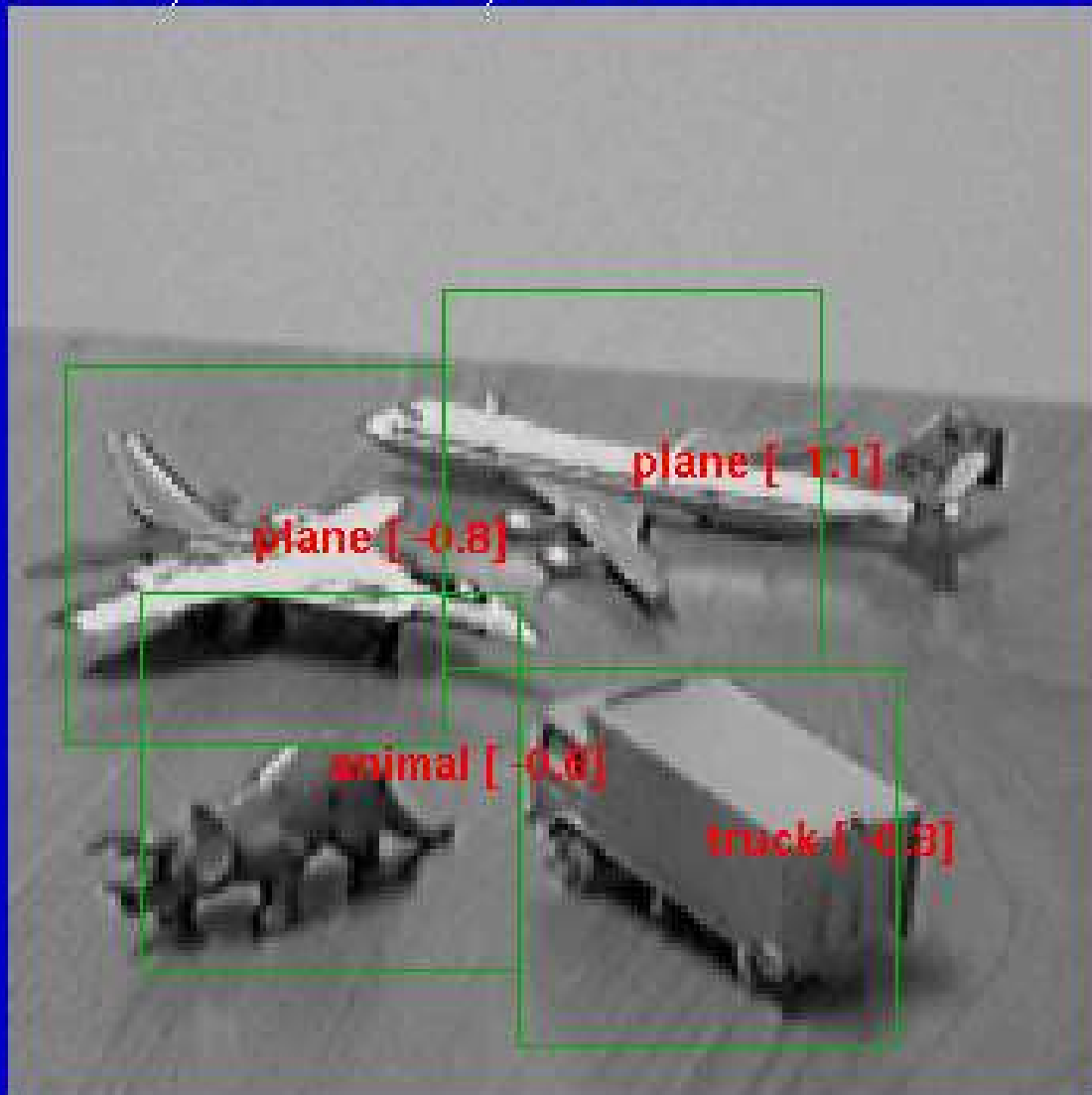
# Examples (Monocular Mode)

Zoom= 1.0, Threshold= -1.2, filter on

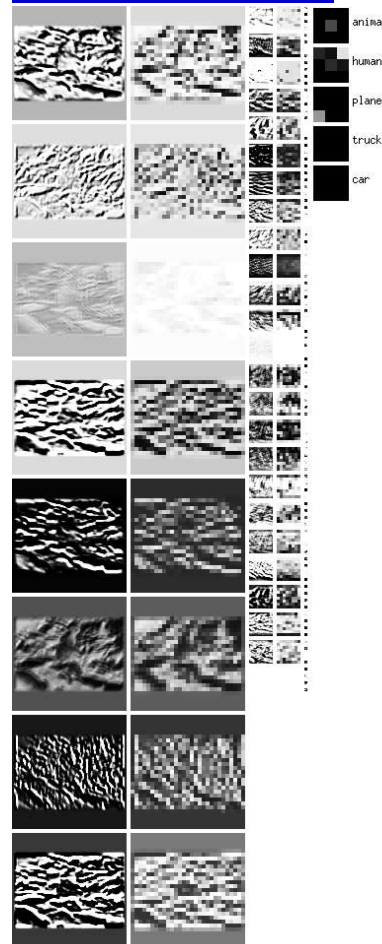
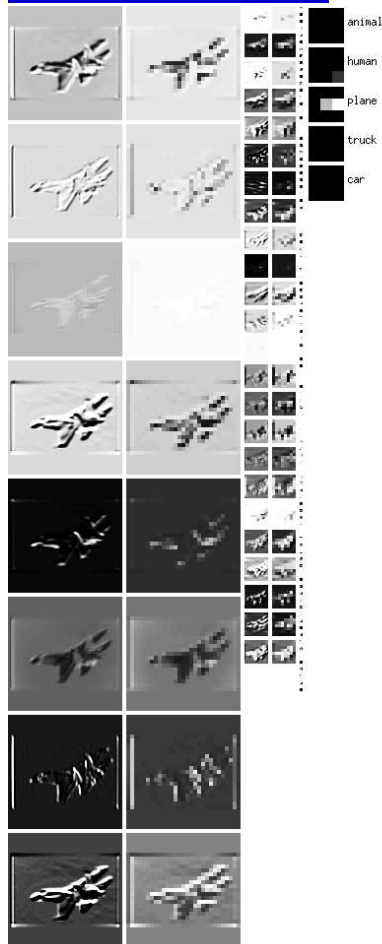
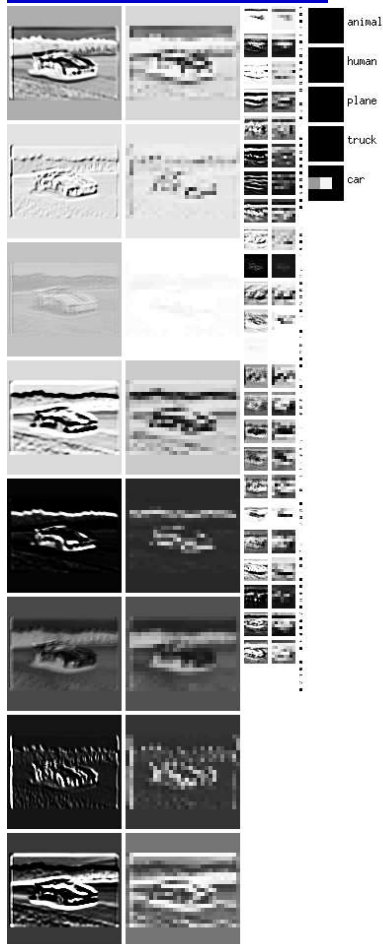


# Examples (Monocular Mode)

Zoom= 0.7, Threshold= -1.8, filter on



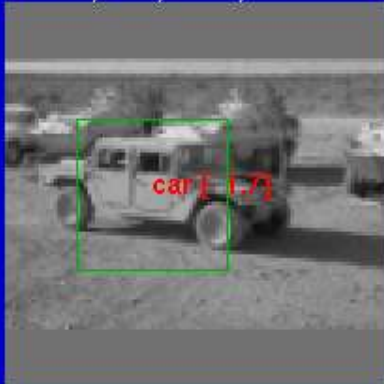
# Natural Images (Monocular Mode)





# Natural Images (Monocular Mode)

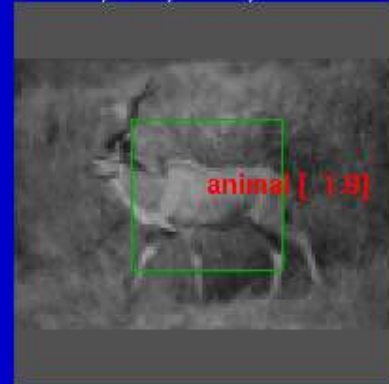
Thrs= 0.5, f on , os=40, nwin=23616



Thrs= 0.5, f on , os=40, nwin=23616



Thrs= 0.5, f on , os=40, nwin=23616



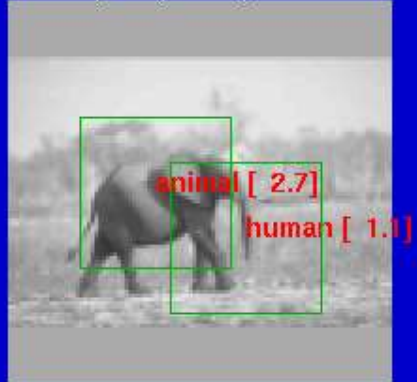
Thrs= 0.5, f on , os=40, nwin=23616



Thrs= 0.5, f on , os=40, nwin=23616



Thrs= 0.5, f on , os=40, nwin=23616



Thrs= 0.5, f on , os=40, nwin=23616



Thrs= 0.5, f on , os=40, nwin=23616



# Natural Images (Monocular Mode)

Thrs= 0.5, f on , os=40, nwin=23616



Thrs= 0.5, f on , os=40, nwin=23616



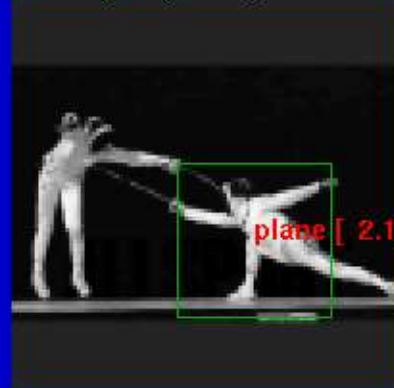
Thrs= 0.5, f on , os=40, nwin=23616



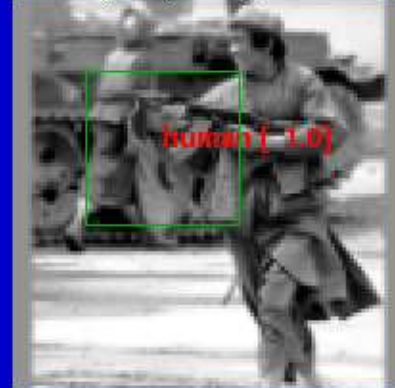
Thrs= 0.5, f on , os=40, nwin=23616



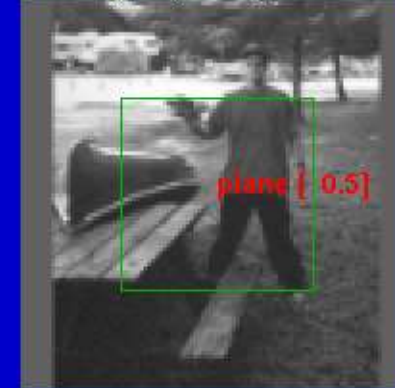
Thrs= 0.5, f on , os=40, nwin=23616



Thrs= 0.5, f on , os=40, nwin=23616

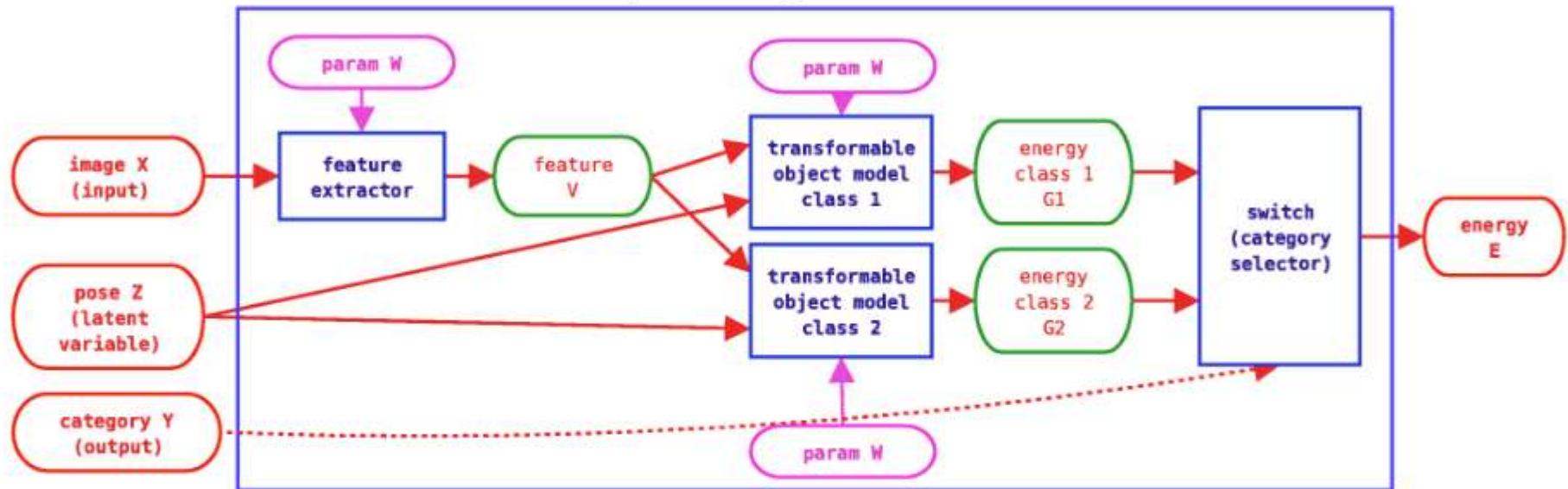


Thrs= 0.5, f on , os=40, nwin=23616



# EBM with Latent Variable for Pose Invariance

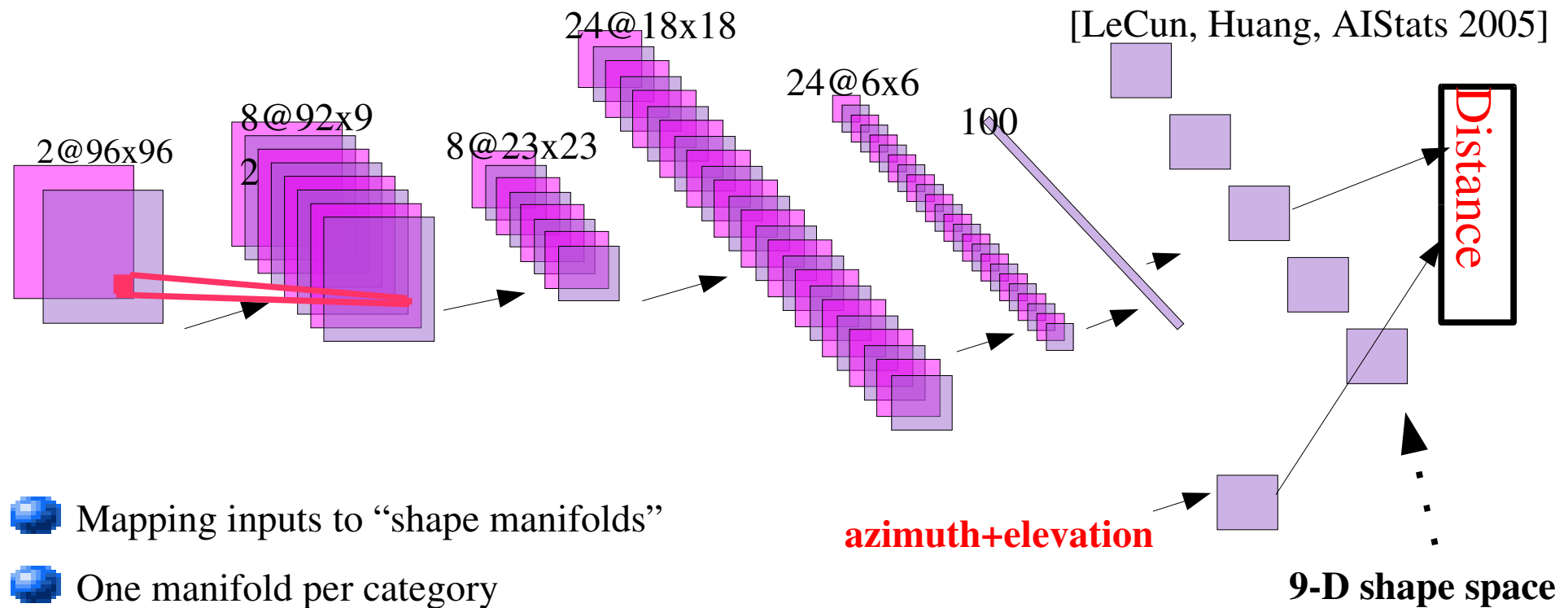
EBM Architecture for invariant object recognition



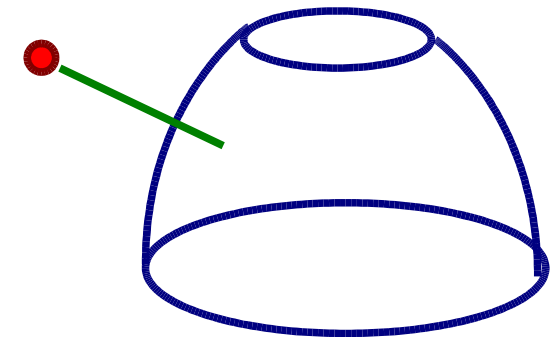
Each object model matches the output of the feature extractor to a reference representation that is transformed by the pose parameters.

**Inference** finds the category and the pose that minimize the energy.

# EBM with a latent pose variable

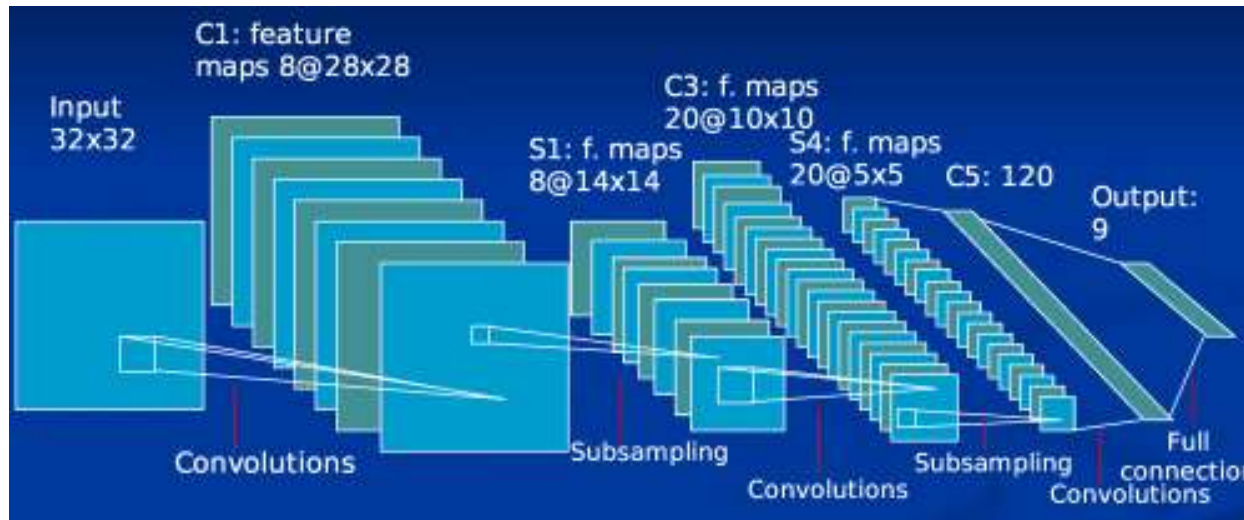


- Mapping inputs to “shape manifolds”
- One manifold per category
- Each manifold is a 2-D half sphere embedded in a 9-D space
- 2 latent variables parameterize position on the manifold (azimuth and elevation).
- Loss function:** pulls the network output toward manifold of desired class, and repel from manifolds of non-desired classes.
- Result on uniform set: 5.1% error (vs 6.6%)**



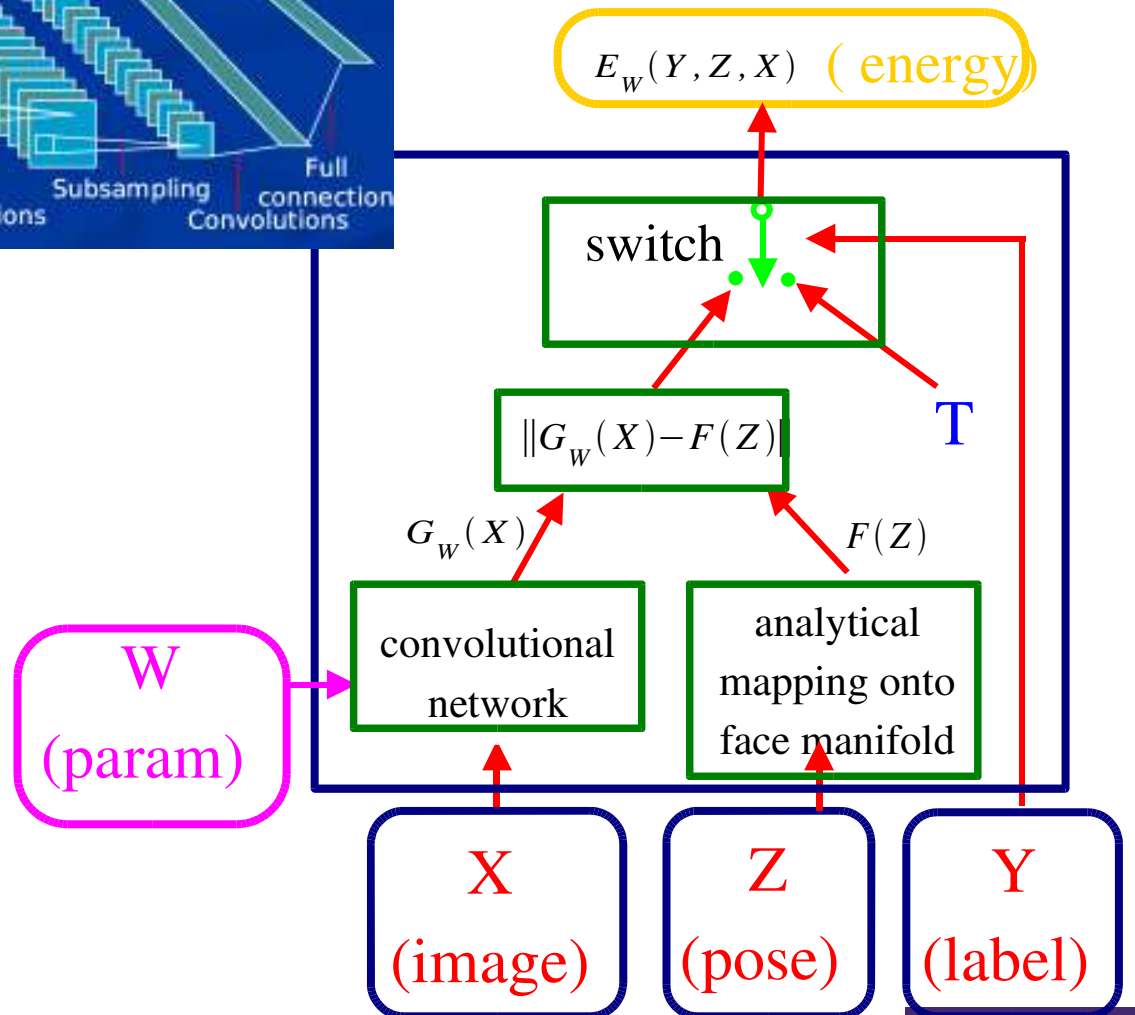


# Face Detection and Pose Estimation with a Convolutional EBM



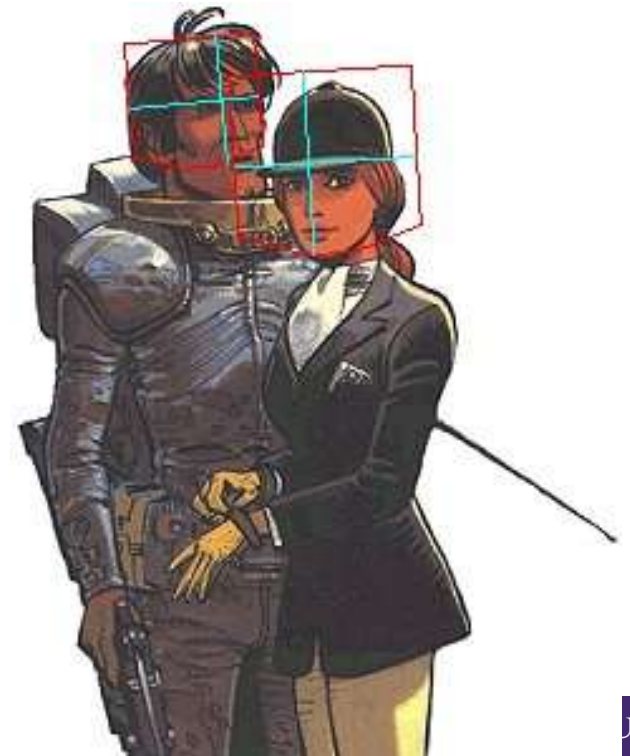
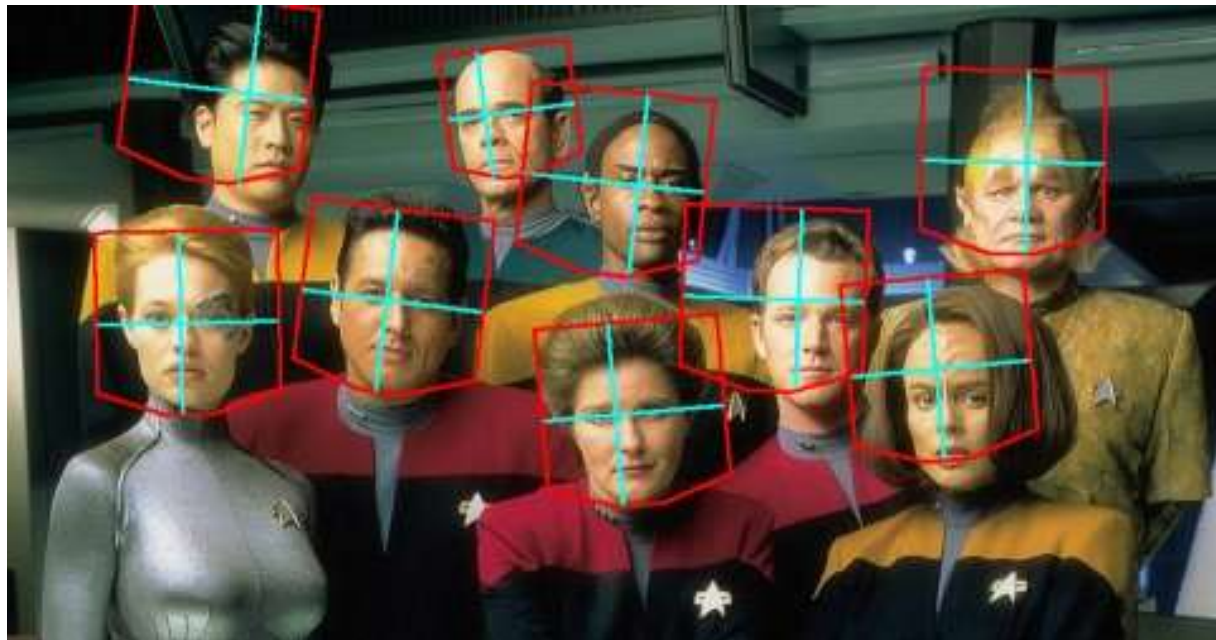
[Osadchy, Miller, LeCun, NIPS 2004]

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



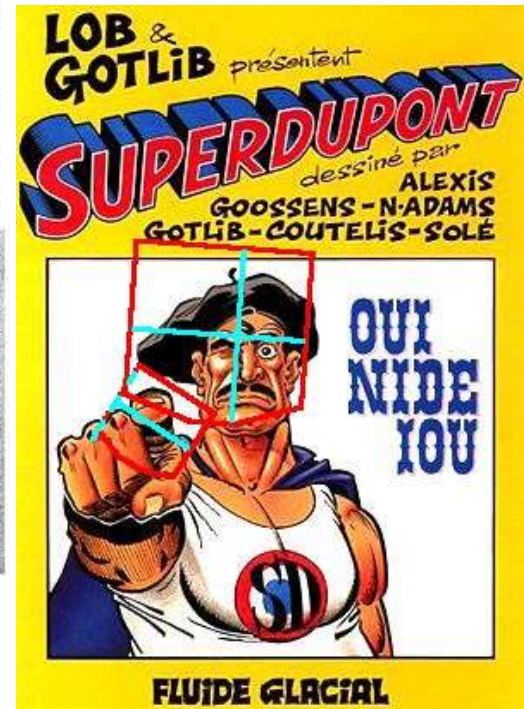
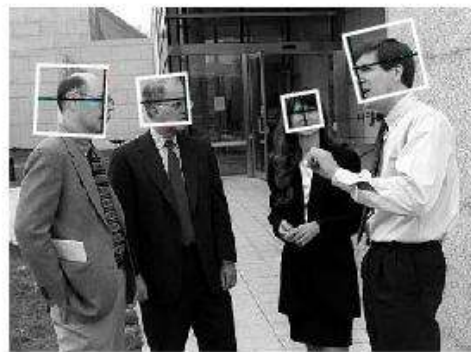
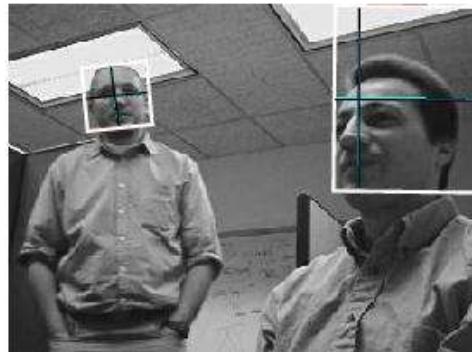
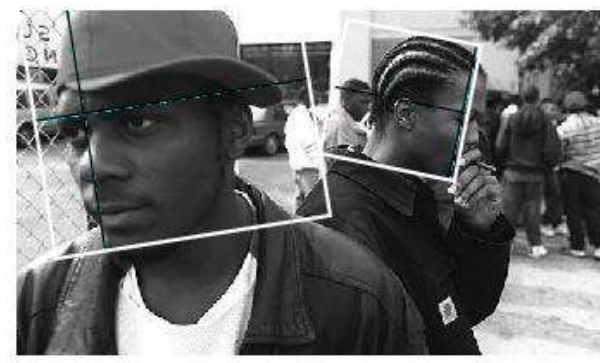
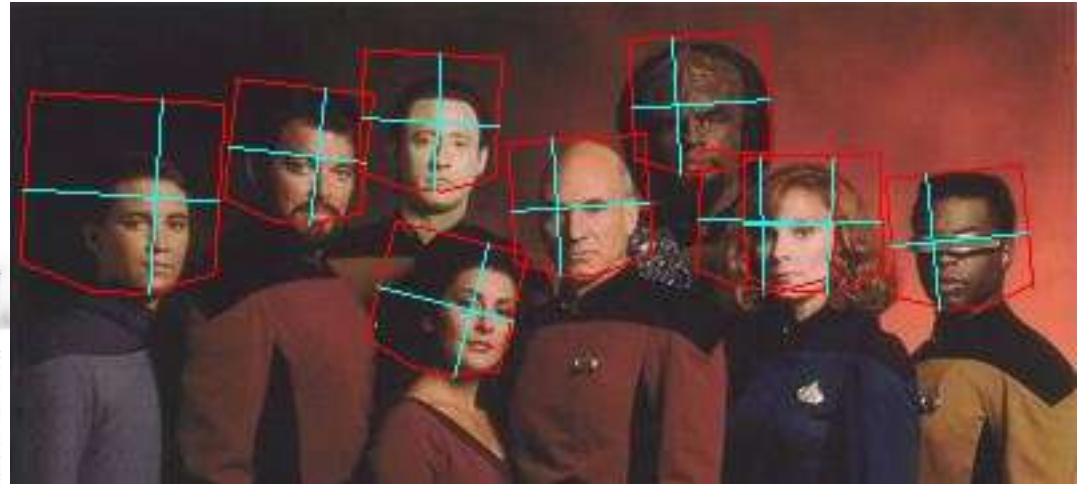
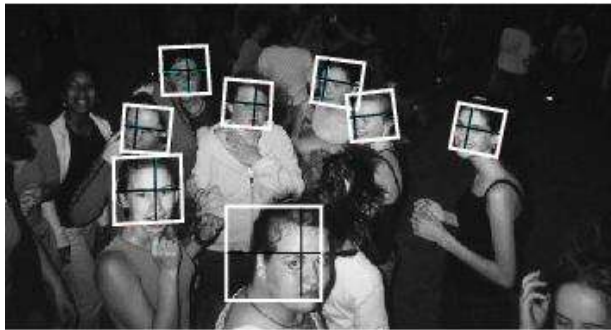
# Face Detection: Results

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





# Face Detection: Results



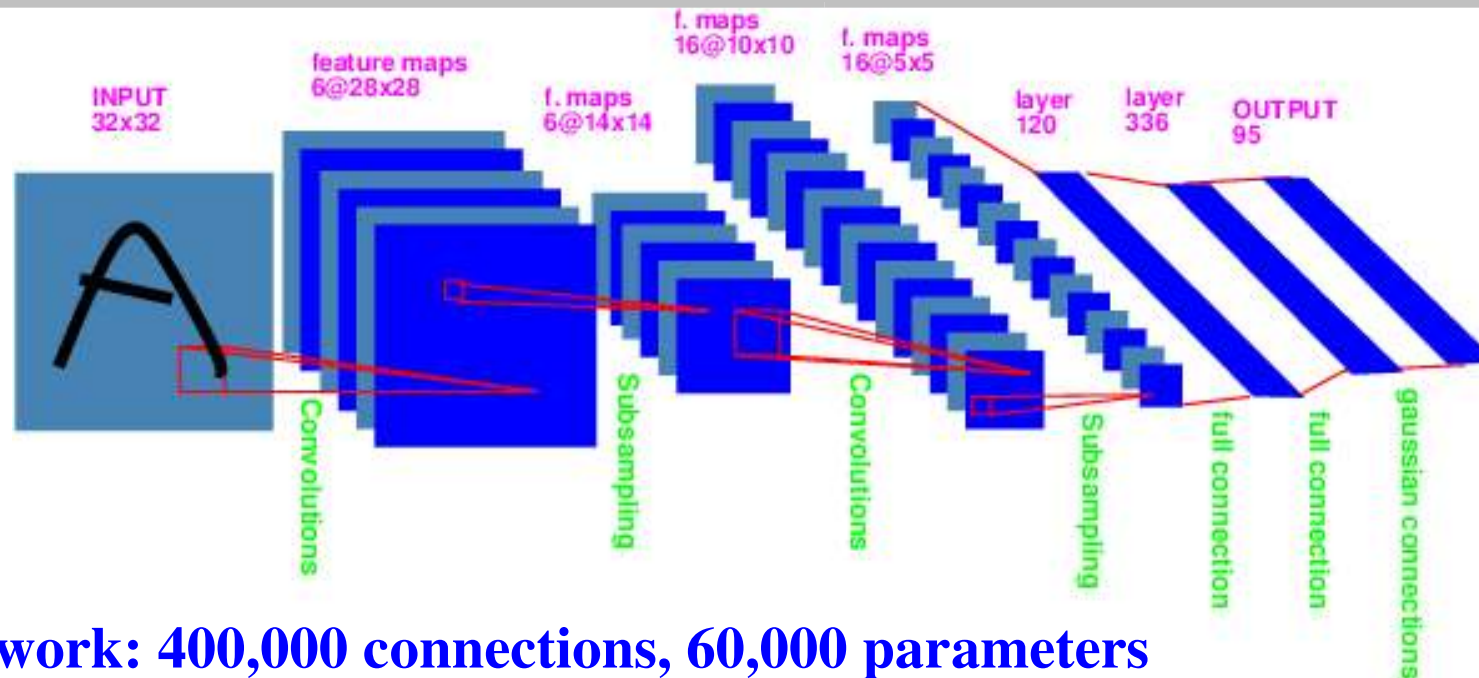


# Face Detection with a Convolutional Net



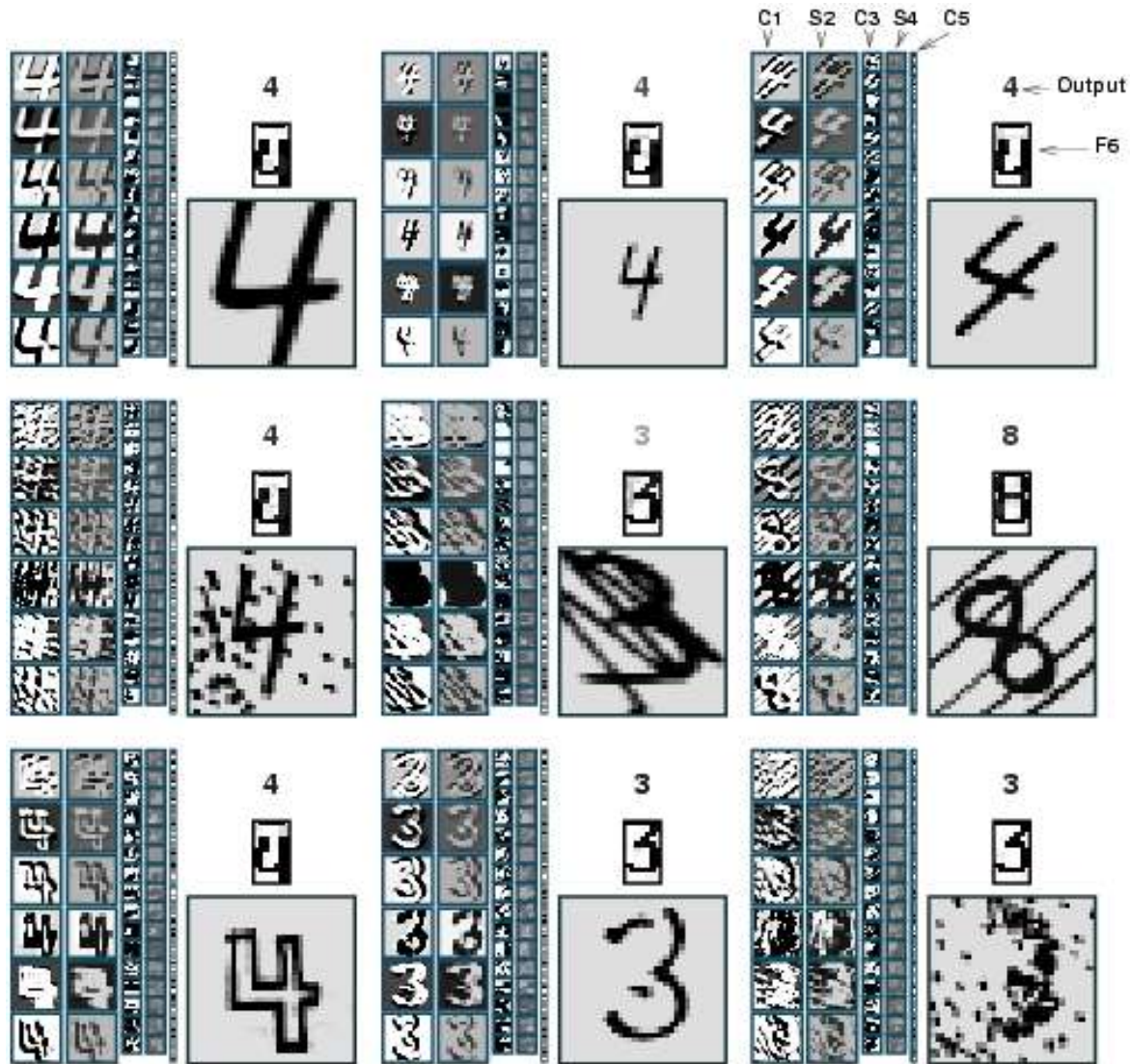


# Handwriting Recognition

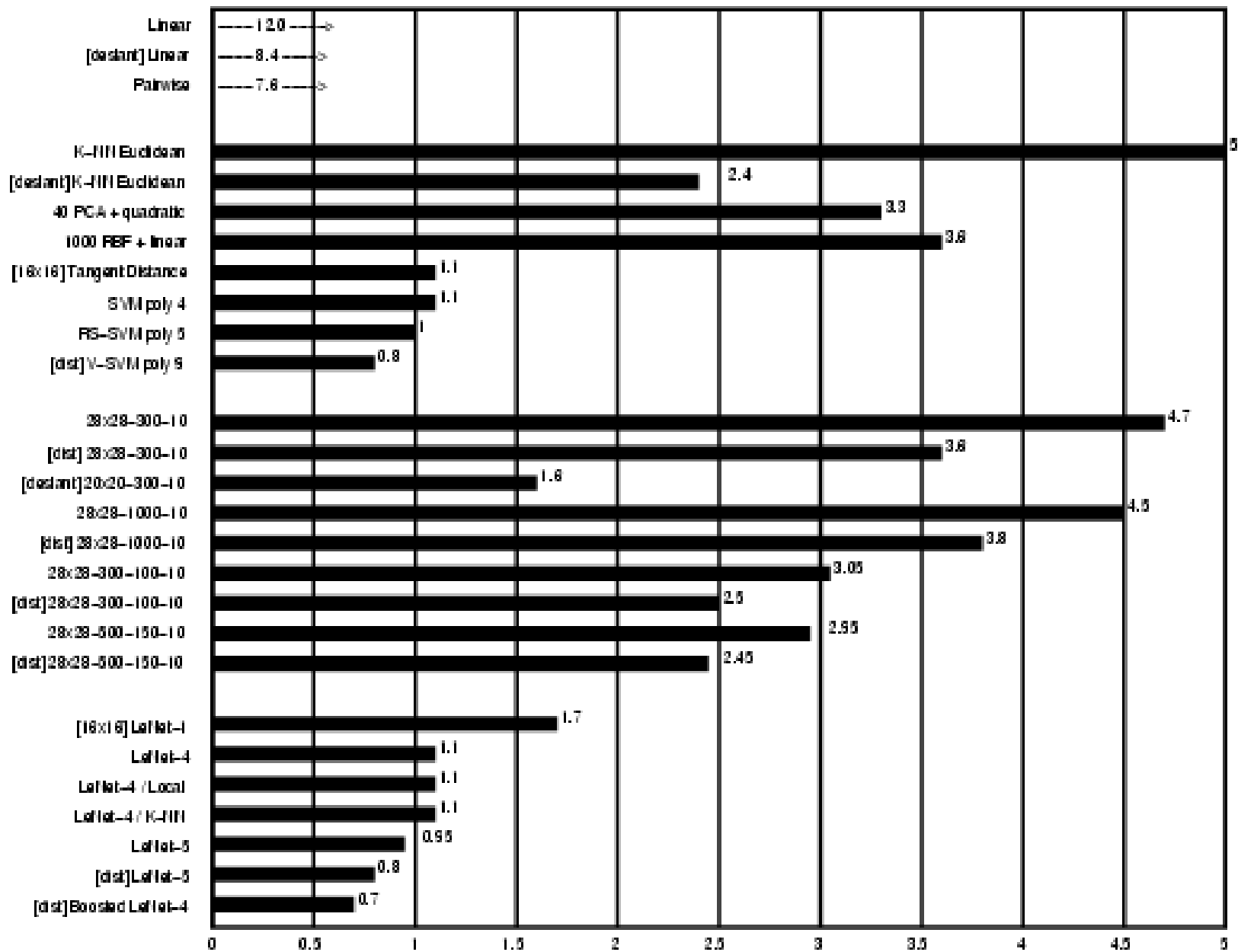


- **Network:** 400,000 connections, 60,000 parameters
- **Input:**  $32 \times 32$  pixels
- **Dataset:** MNIST: 60,000 handwritten digits for training, 10,000 for testing.
- **Results:** 0.8% error on test set
- **Simard et al.** recently obtained 0.4% error with a similar architecture

# Handwriting Recognition



# Handwriting Recognition

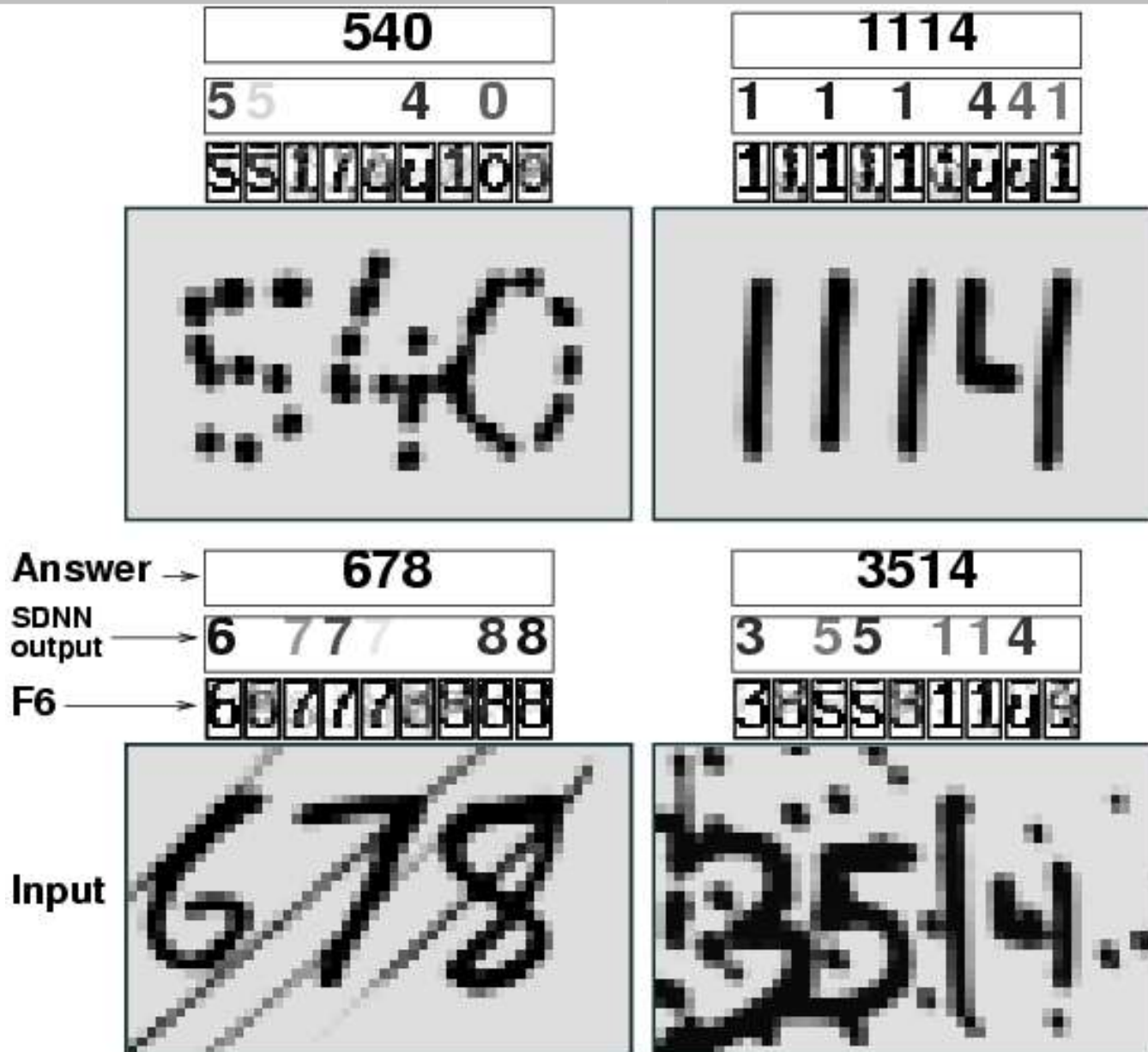


# Handwriting Recognition





# Handwriting Recognition

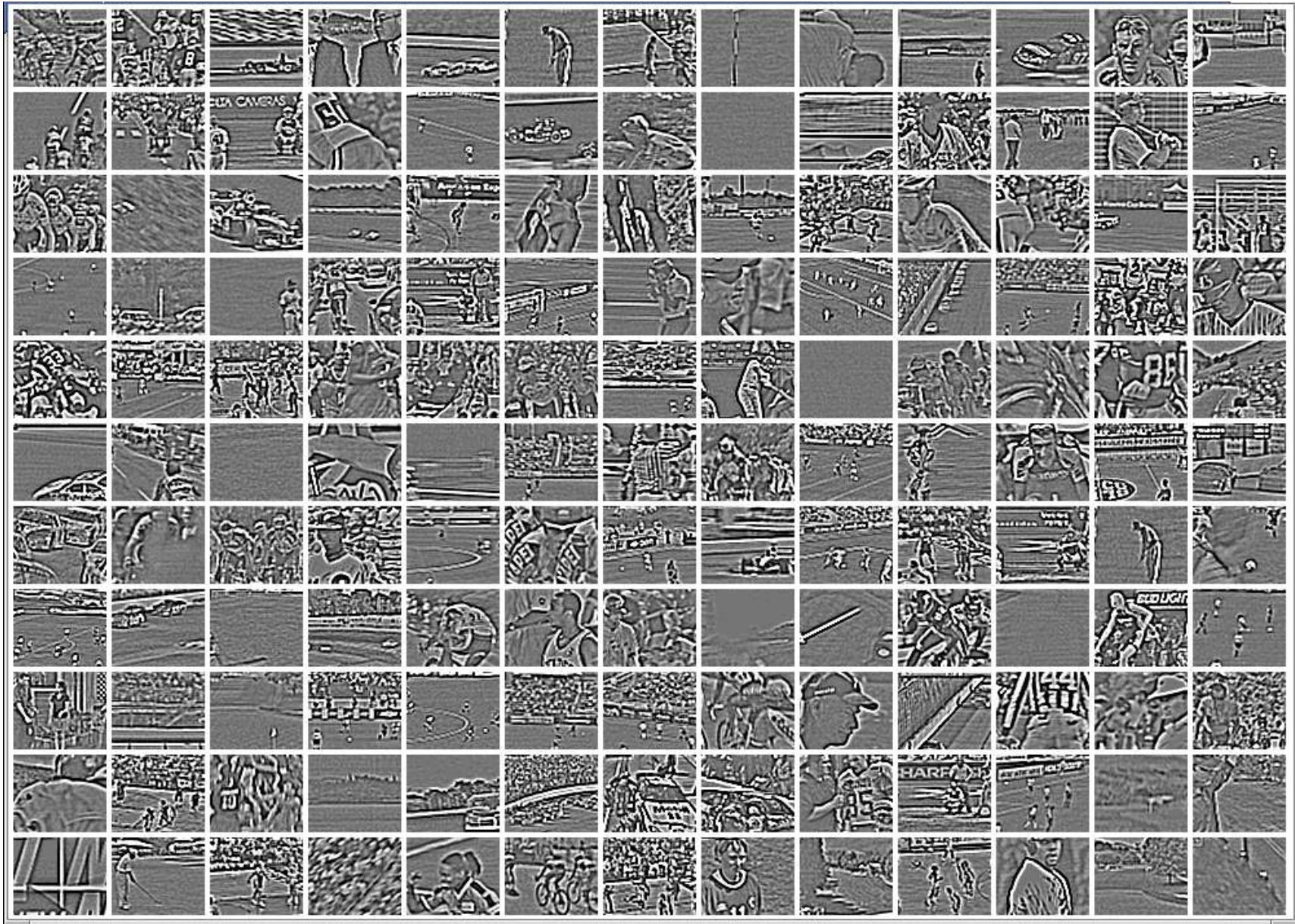


## TV sport categorization (with Alex Niculescu, Cornell)

- **Classifying TV sports snapshots into 7 categories: auto racing, baseball, basketball, bicycle, golf, soccer, football.**
- **123,900 training images (300 sequence with 59 frames for each sport)**
- **82,600 test images (200 sequences with 59 frames for each sport)**
- **Preprocessing: convert to YUV, high-pass filter the Y component, crop, subsample to 72x60 pixels**
- **Results:**
  - ▶ frame-level accuracy: 61% correct
  - ▶ Sequence-level accuracy 68% correct (simple voting scheme).



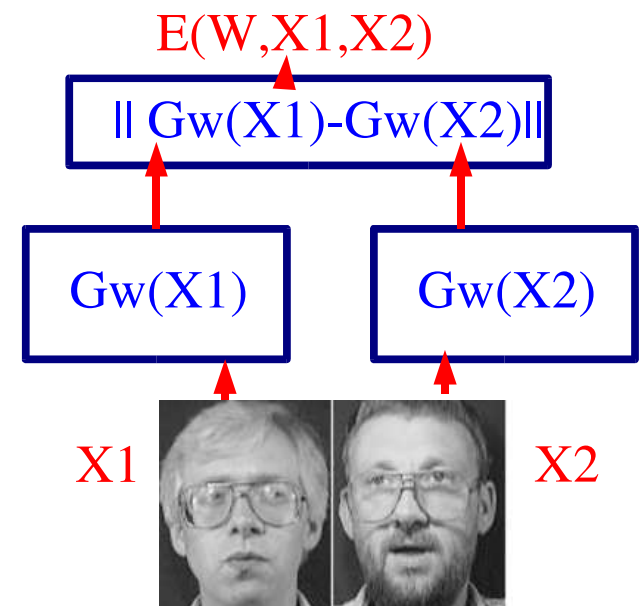
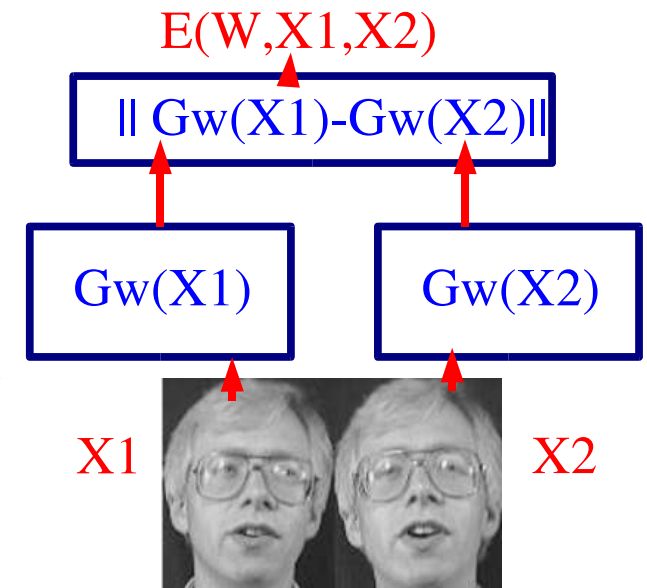
# TV sport categorization (with Alex Niculescu, Cornell)



# Learning an Invariant Dissimilarity Metric with EBMs

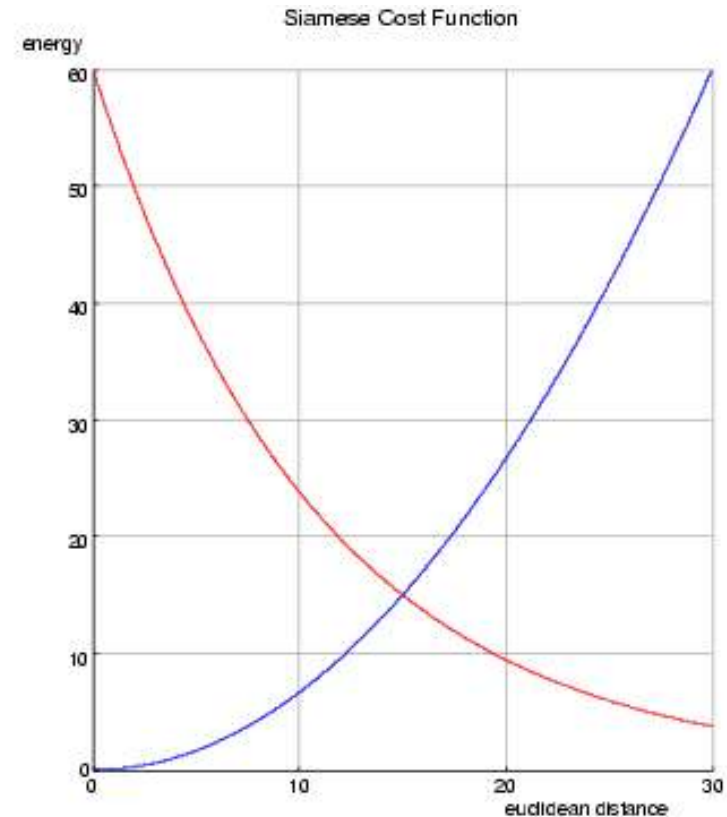
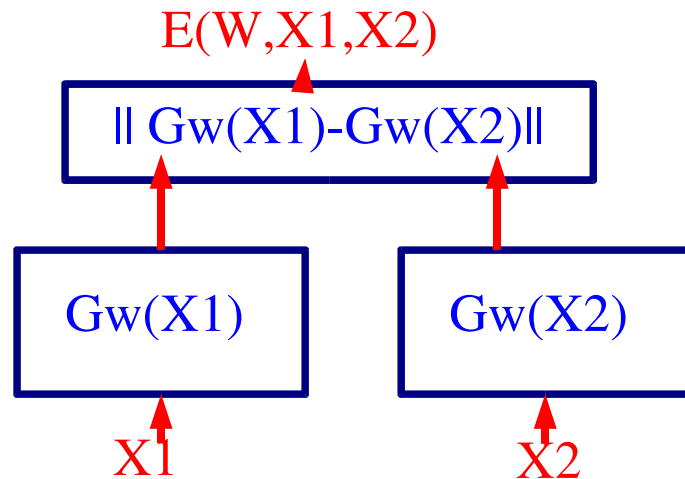
[Chopra, Hadsell, LeCun CVPR 2005]

- Training a **parameterized, invariant dissimilarity metric** may be a solution to the **many-category problem**.
- Find a mapping  $G_w(X)$  such that the Euclidean distance  $\|G_w(X1) - G_w(X2)\|$  reflects the “semantic” distance between  $X1$  and  $X2$ .
- Once trained, a trainable dissimilarity metric can be used to classify **new categories using a very small number of training samples** (used as prototypes).
- This is an example where probabilistic models are too constraining, because we would have to limit ourselves to models that can be normalized over the space of input pairs.
- With EBMs, we can put what we want in the box (e.g. A convolutional net).
- Siamese Architecture**
- Application:** face verification/recognition





# Learning an Invariant Dissimilarity Metric with EBMs



- **Siamese models:** distance between the outputs of two identical copies of a model.
- $E(W, X_1, X_2) = \|G_w(X_1) - G_w(X_2)\|$
- If  $X_1$  and  $X_2$  are from the **same category**, train the two copies of the model to produce **similar outputs**
- If  $X_1$  and  $X_2$  are from **different categories**, train the two copies of the model to produce **different outputs**
- Loss function: square-exponential loss:

$$L(W, Y, X_1, X_2) = (1 - Y) \cdot \frac{2}{R} (\|G_w(X_1) - G_w(X_2)\|)^2 + Y \cdot 2R e^{-\frac{K}{R} \|G_w(X_1) - G_w(X_2)\|}$$

# Face Verification datasets: AT&T/ORL

- The AT&T/ORL dataset
- Total subjects: **40**. Images per subject: **10**. Total images: **400**.
- Images had a **moderate** degree of variation in pose, lighting, expression and head position.
- Images from **35** subjects were used for training. Images from **5** remaining subjects for testing.
- Training set was taken from: **3500** genuine and **119000** impostor pairs.
- Test set was taken from: **500** genuine and **2000** impostor pairs.
- <http://www.uk.research.att.com/facedatabase.html>



**AT&T/ORL**  
**Dataset**



# Face Verification datasets: AR/Purdue dataset

- **The AR/Purdue dataset**
- Total subjects: **136**. Images per subject: **26**. Total images: **3536**.
- Each subject has 2 sets of 13 images taken 14 days apart.
- Images had **very high** degree of variation in pose, lighting, expression and position. Within each set of 13, there are 4 images with expression variation, 3 with lighting variation, 3 with dark sun glasses and lighting variation, and 3 with face obscuring scarfs and lighting variation.
- Images from **96** subjects were used for training. The remaining **40** subjects were used for testing.
- **Training set drawn from:** **64896** genuine and **6165120** impostor pairs.
- **Test set drawn from:** **27040** genuine and **1054560** impostor pairs.
- [http://rv11.ecn.purdue.edu/aleix/aleix\\_face\\_DB.html](http://rv11.ecn.purdue.edu/aleix/aleix_face_DB.html)



# Face Verification dataset: AR/Purdue





# Dataset for Verification

# Verification Results

tested on AT&T and AR/Purdue

The AT&T dataset

The AR/Purdue dataset

## AT&T dataset

Number of subjects: 5  
Images/subject: 10  
Images/Model: 5  
Total test size: 5000  
Number of Genuine: 500  
Number of Impostors: 4500

False Accept False Reject

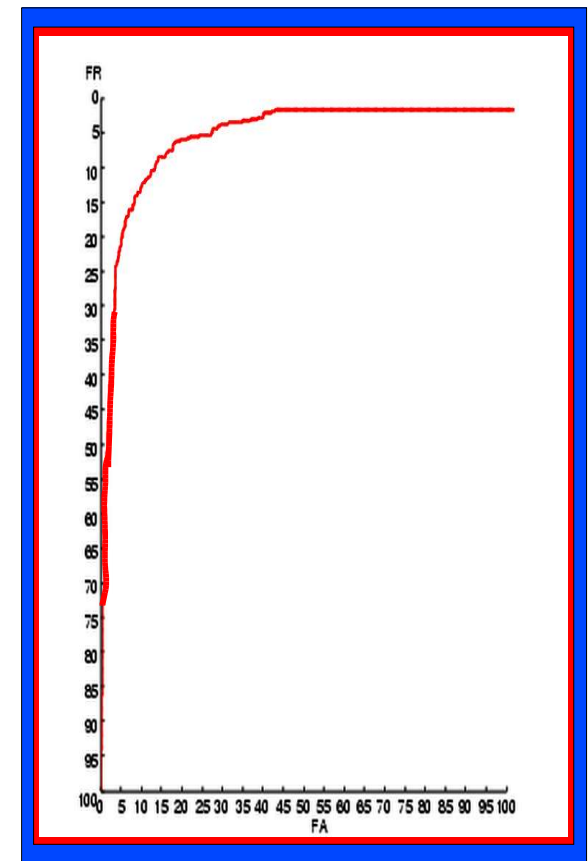
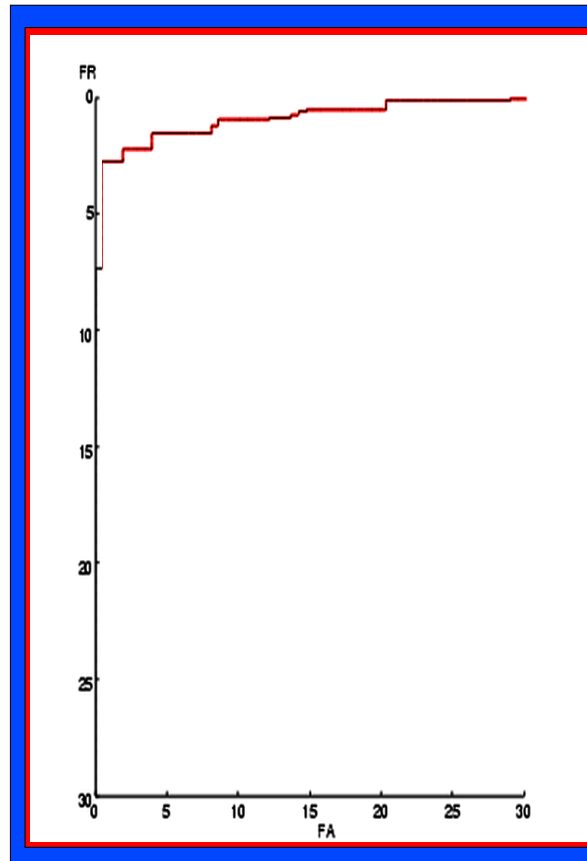
10.00% 0.00%  
7.50% 1.00%  
5.00% 1.00%

False Accept False Reject

10.00% 11.00%  
7.50% 14.60%  
5.00% 19.00%

## Purdue/AR dataset

Number of subjects: 40  
Images/subject: 26  
Images/Model: 13  
Total test size: 5000  
Number of Genuine: 500  
Number of Impostors: 4500

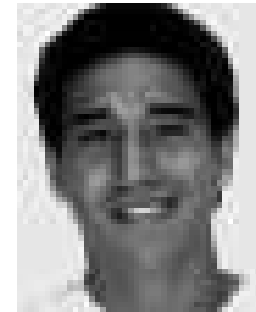


# Internal state for genuine and impostor pairs



# Classification Examples

## Example: Correctly classified genuine pairs

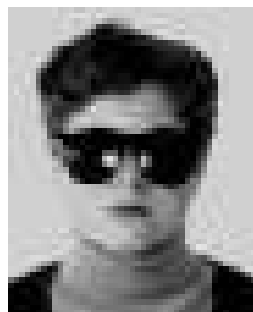


energy: 0.3159

energy: 0.0043

energy: 0.0046

## Example: Correctly classified impostor pairs



energy: 20.1259

energy: 32.7897

energy: 5.7186

## Example: Mis-classified pairs



energy: 10.3209

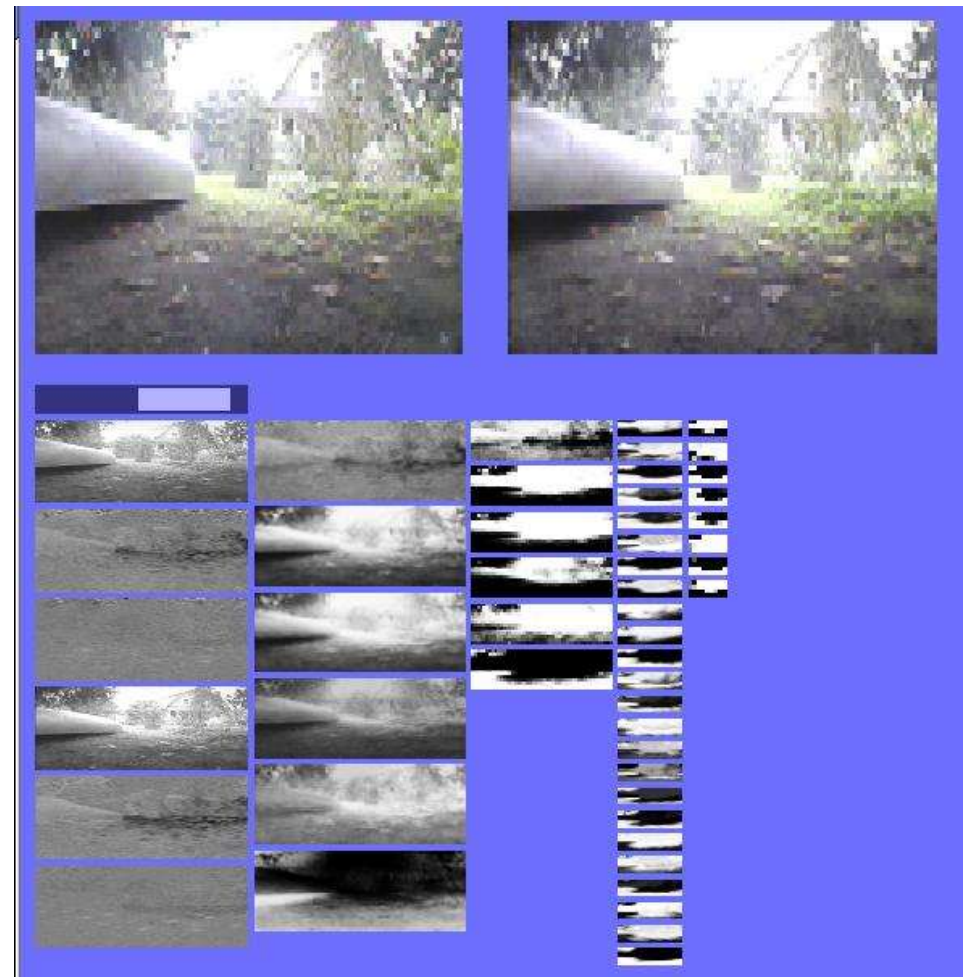
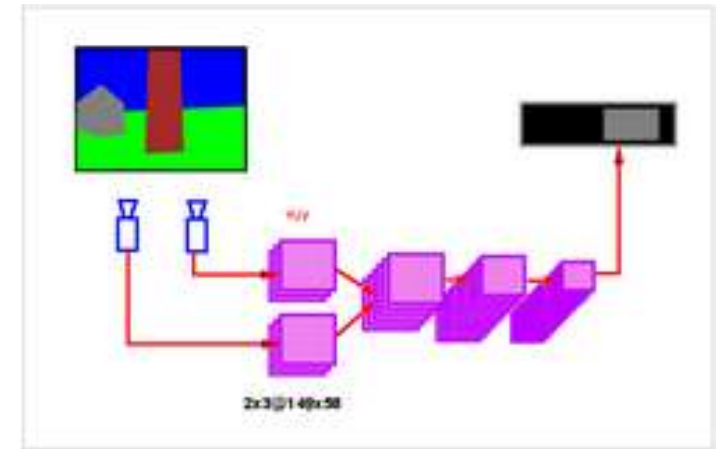


energy: 2.8243



# Visual Navigation for a Mobile Robot

- Mobile robot with two cameras
- The convolutional net is trained to emulate a human driver from recorded sequences of video + human-provided steering angles.
- The network maps stereo images to steering angles for obstacle avoidance





# Invariant Object Recognition

- The old feed-forward architecture can do more than expected.
- **Full invariance to viewpoint and illumination** for detecting and recognizing objects can be learned discriminatively by a **simple feed-forward architecture**.
- With **only 5 training instances** from each category, the model can detect and recognize new instances with high accuracy.
- The model outperforms “traditional” template-based classifiers operating on raw pixels or on PCA features.
- The system takes advantage of the binocular input.
- The convolutional net architecture is generic, and can be applied to a variety of vision tasks with essentially no change.
- Feature tuning produces very parsimonious systems with only a small number of feature detectors at each layer.
- Invariance can be achieved with “deep” architectures, containing multiple, successive layers of feature detection and feature integration/subsampling (Hubel/Wiesel'62, Fukushima'72, LeCun'89, Ullman'02, Riesenhuber/Poggio'02).