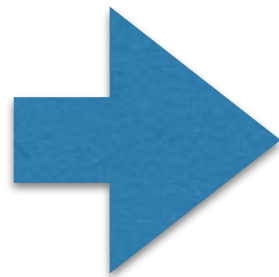# Parts

# Limits of templates
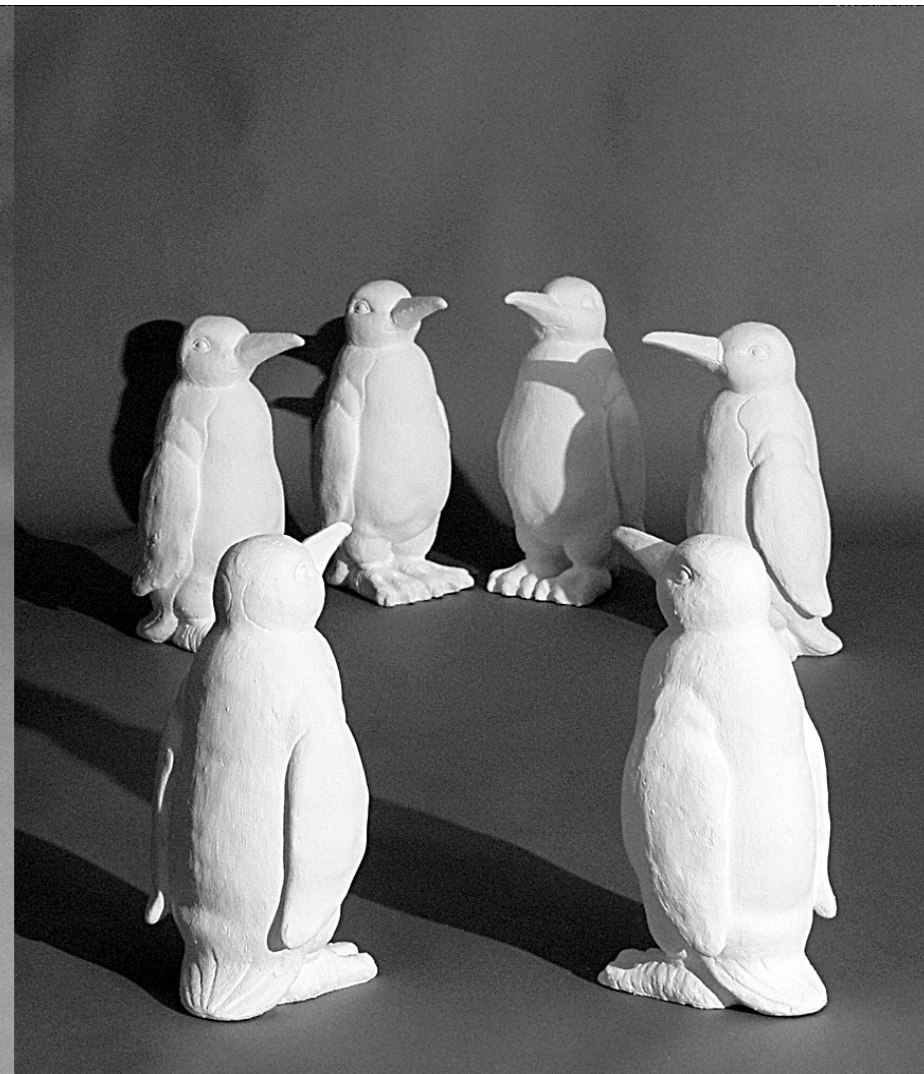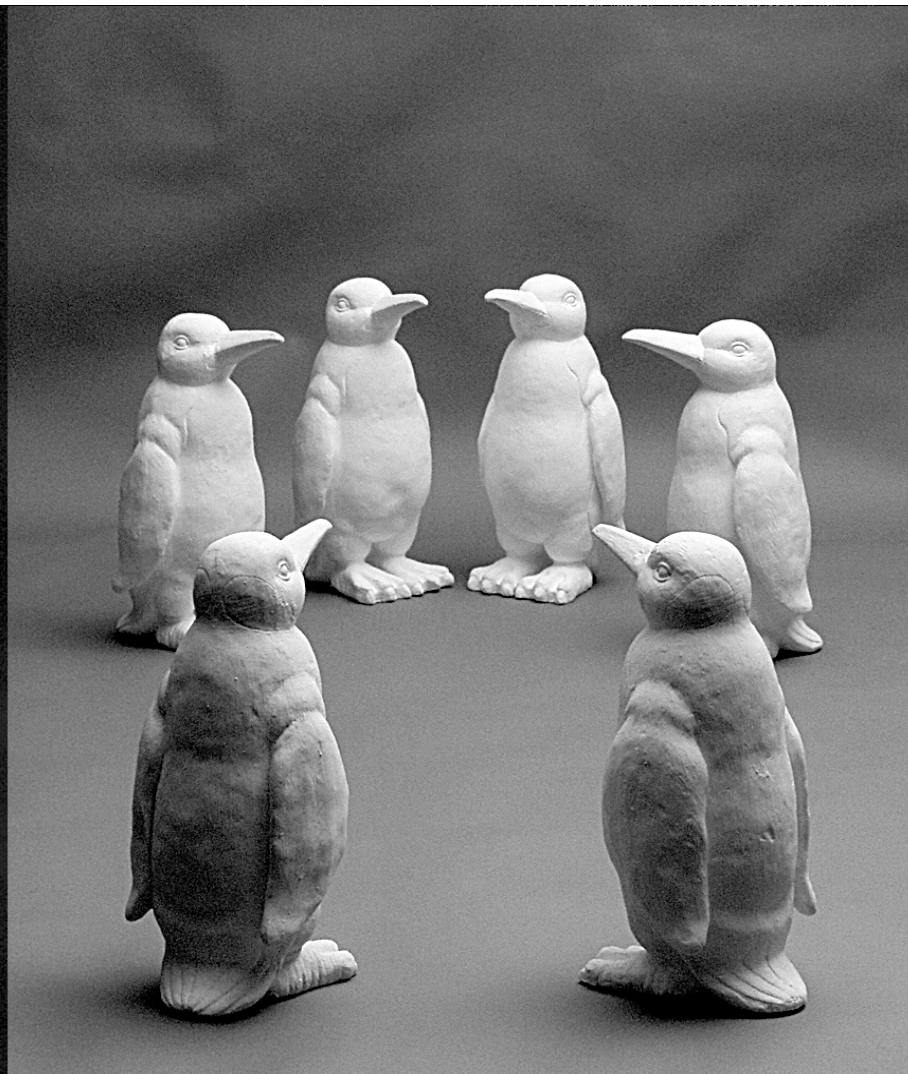
# How to model large variations in appearance?

This is generally regarded as a "central challenge" for recognition

# Challenges: viewpoint variation



Michelangelo 1475-1564

# Challenges: illumination

# Challenges: scale

# Challenges: background clutter



Kilmeny Niland. 1995

# Challenges: intra-class variation

# Why is finding people difficult?


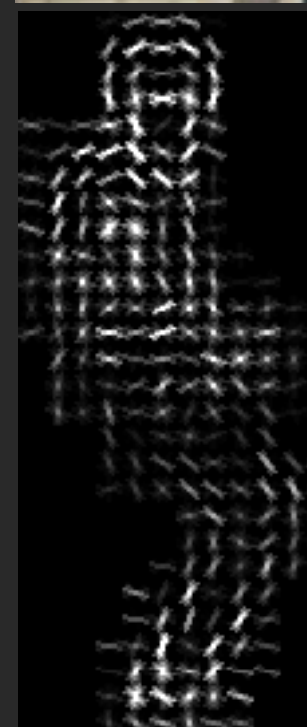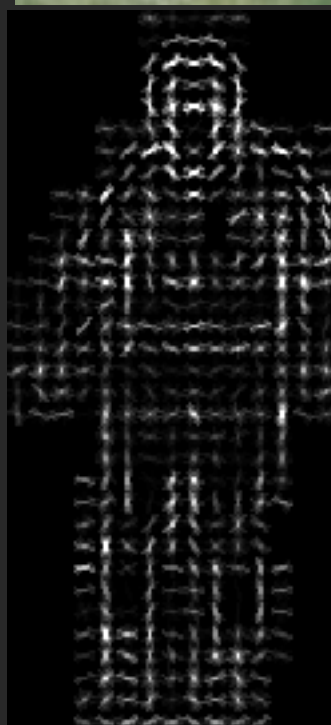
variation in illumination



variation in appearance



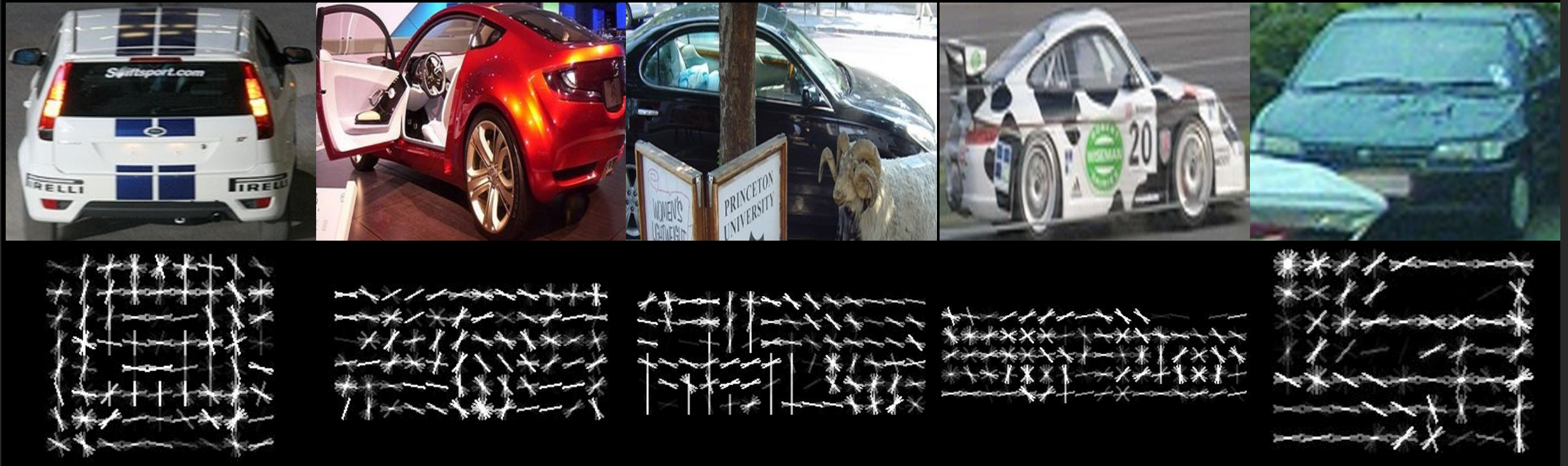variation in pose, viewpoint



occlusion & clutter

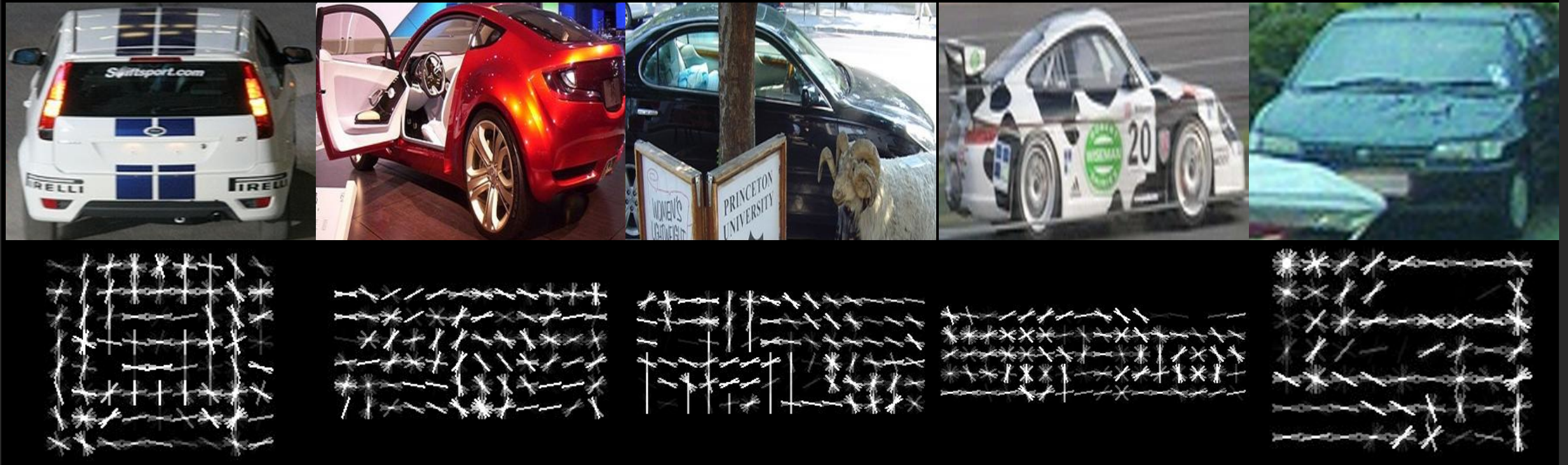Classic "nuisance factors" for general object recognition

# "Sub"categories



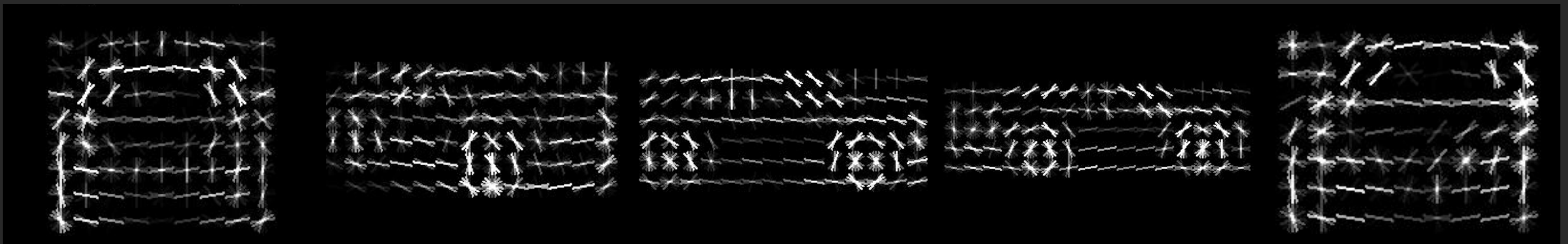Train sub-category templates for each type of pose, body-shape, etc.

# Why not treat each positive example as a unique subcategory?



Iter1

Figure 4: We visualize examples training images on the top. We show initial exemplar models trained with them

# Why not treat each positive example as a unique subcategory?



Iter0

Single positive example

Average of 50 closest

Figure 4: We visualize examples training images on the
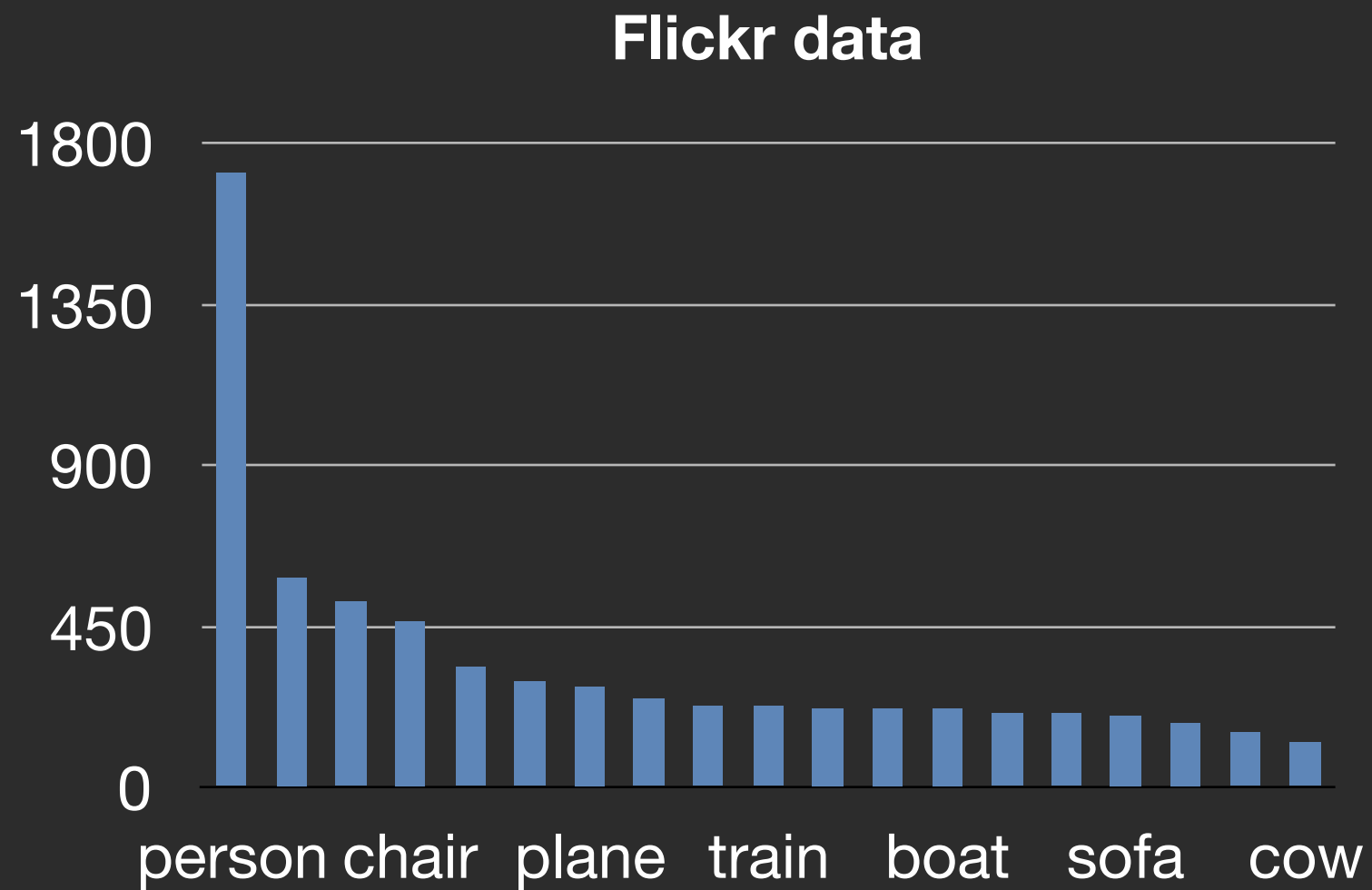
# But how to handle...



We need lots of templates, but will likely have little data of 'twisted' poses

# But how to handle...



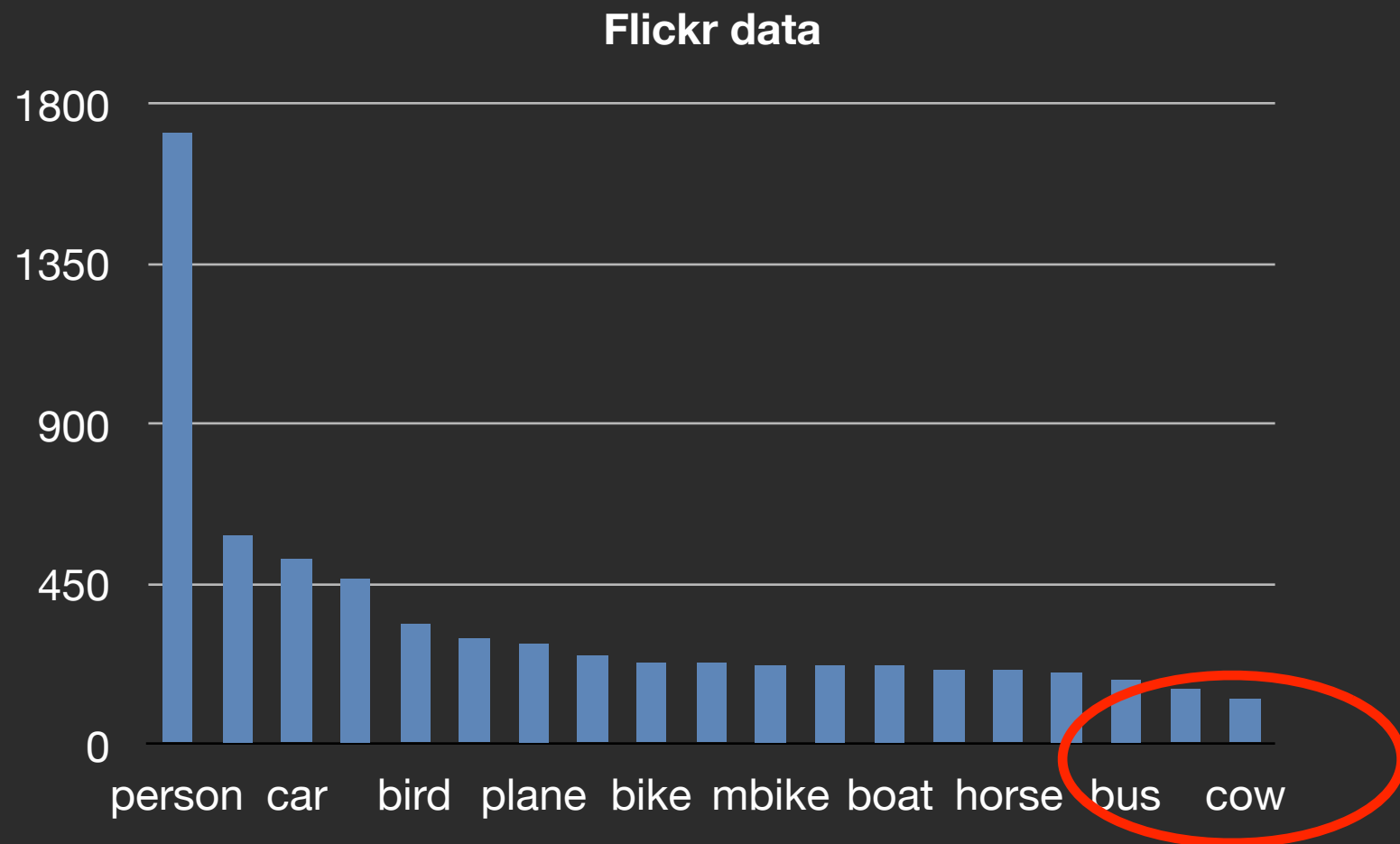We need lots of templates, but will likely have little data of 'rare' car-appearances
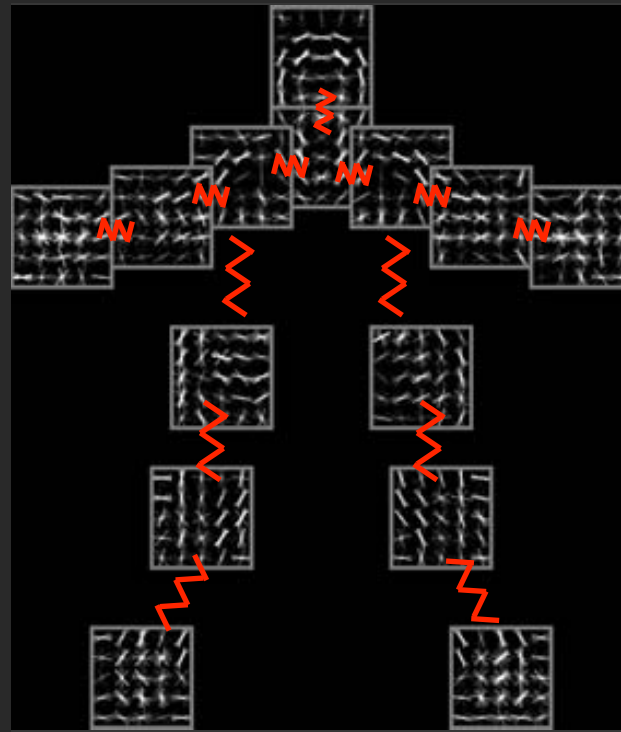
# Difficulties: long tails

**Flickr data**



person chair plane train boat sofa cow

# Difficulties: long tails
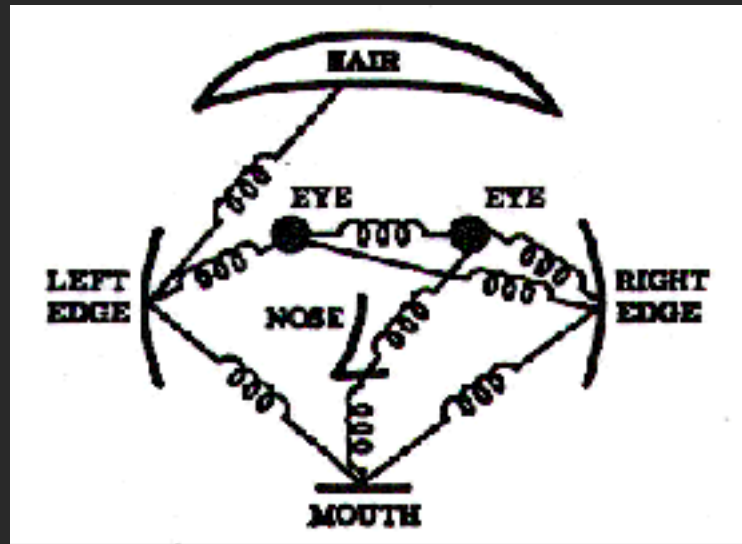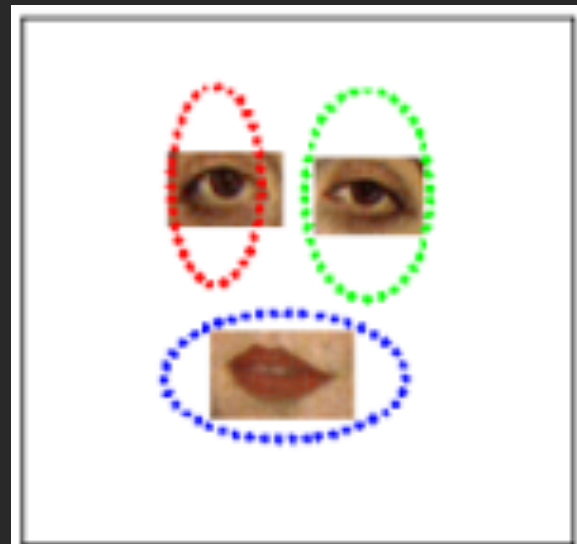
**Flickr data**
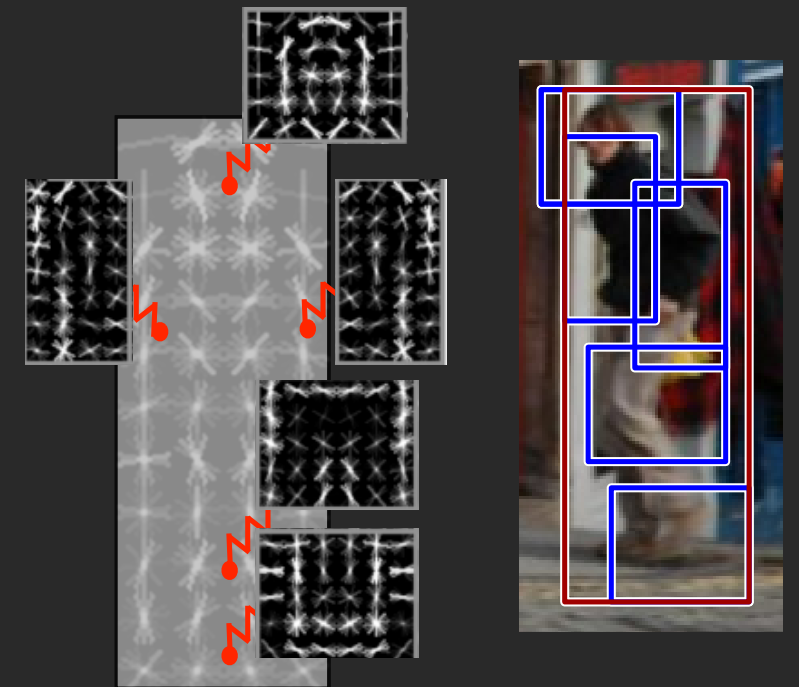


"One-shot learning": sharing

# Parts to the rescue!

# History over 40 years



Pictorial
structures

Constellation
models

Deformable
part models

Model encodes local appearance + pairwise geometry

Pictorial Structures (Fischler & Elschlager 73, Felzenswalb and Huttenlocher 00)
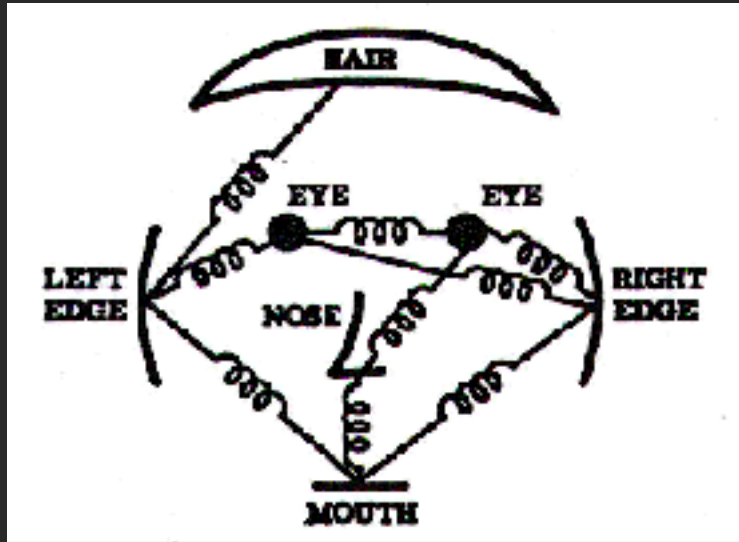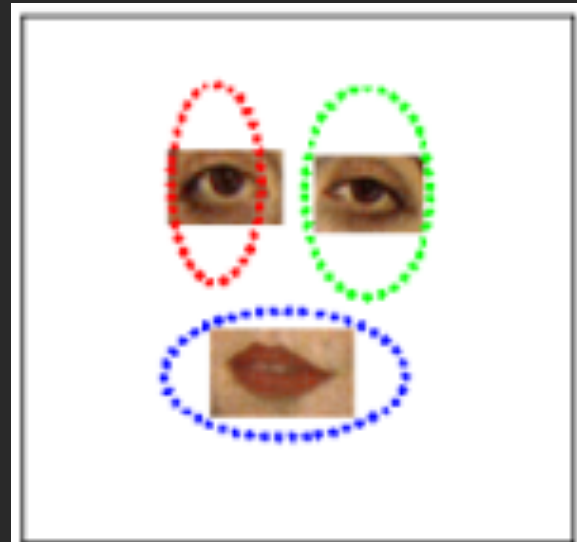Cardboard People (Yu et al 96)
Body Plans (Forsyth & Fleck 97)
Active Appearance Models (Cootes & Taylor 98)
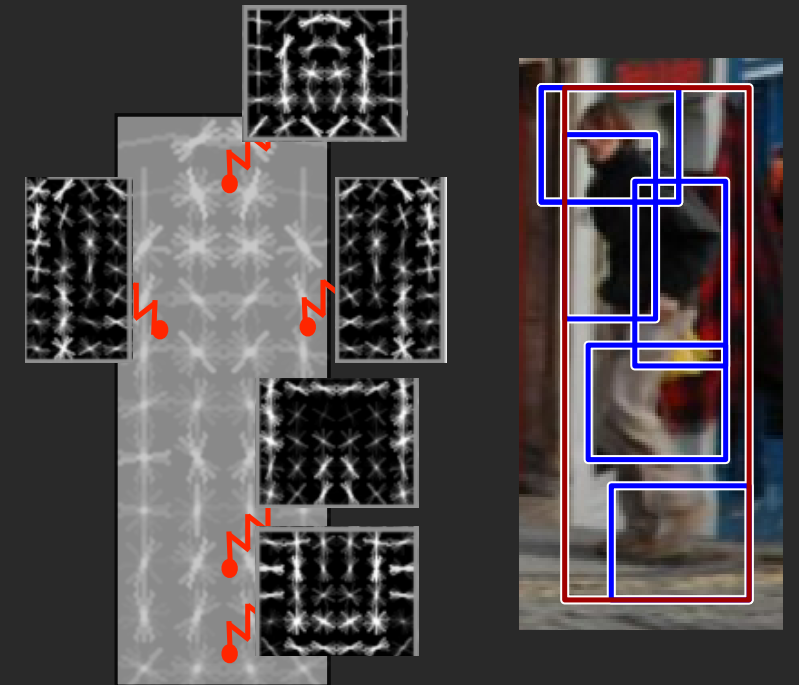Constellation Models (Burl et all 98, Fergus et al 03)

# Part models



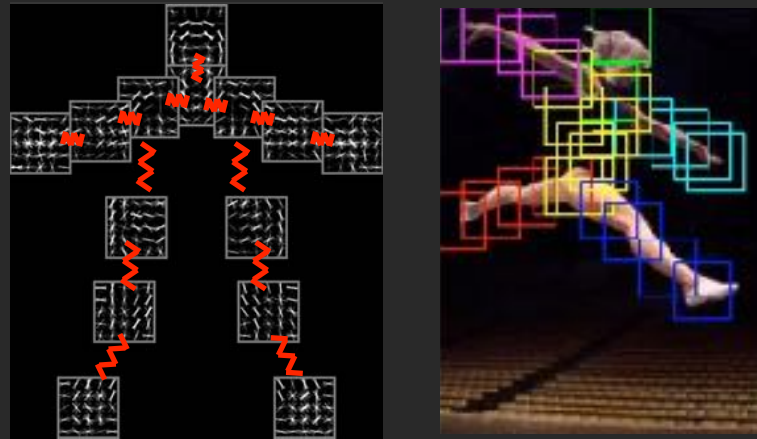Pictorial structures



Constellation models



Deformable part models

I'll talk about DPMs, but give an alternate "long tail" perspective

Felzenszwalb, Girshick, McAllester, Ramanan CACM 2013
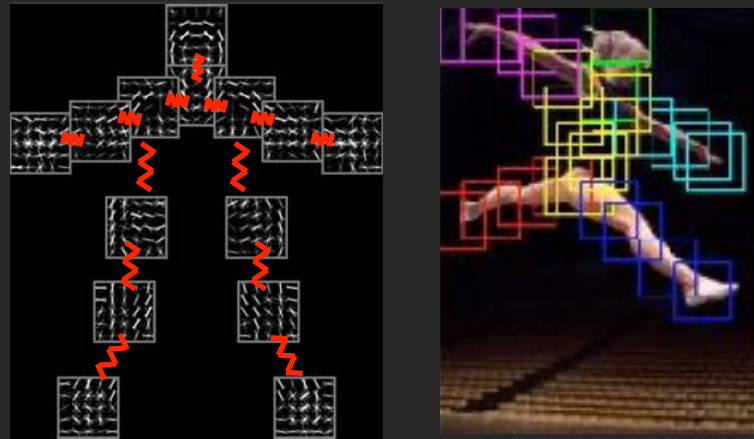
# Scoring function

$$S(x$$

$$x = \text{image}$$
$$z_i = (x_i, y_i)$$
$$z = \{z_1, z_2, \ldots\}$$

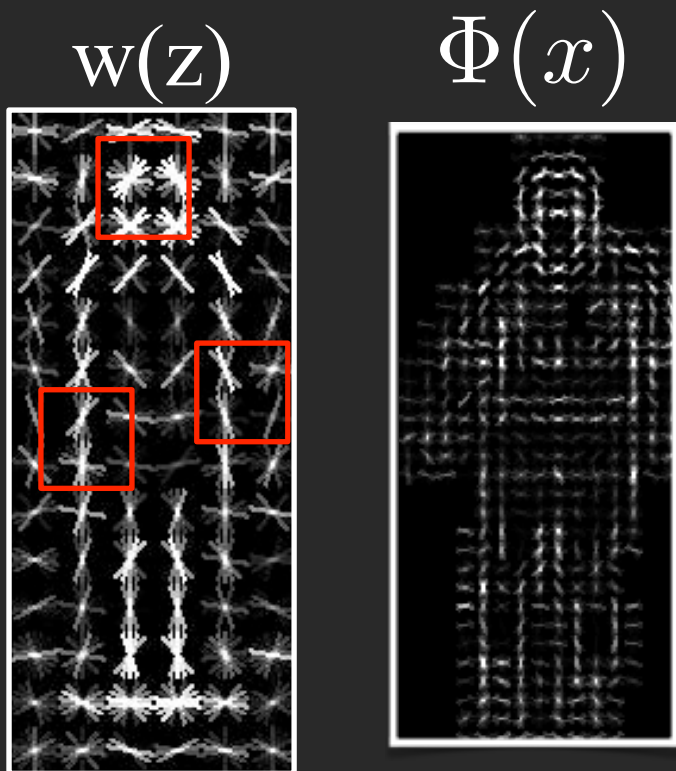# Scoring function

$$S(x, z) = \sum_i$$

$$x = \text{image}$$
$$z_i = (x_i, y_i), \quad i \in \{\text{head,elbow},\dots\}$$
$$z = \{z_1, z_2, \dots\}$$

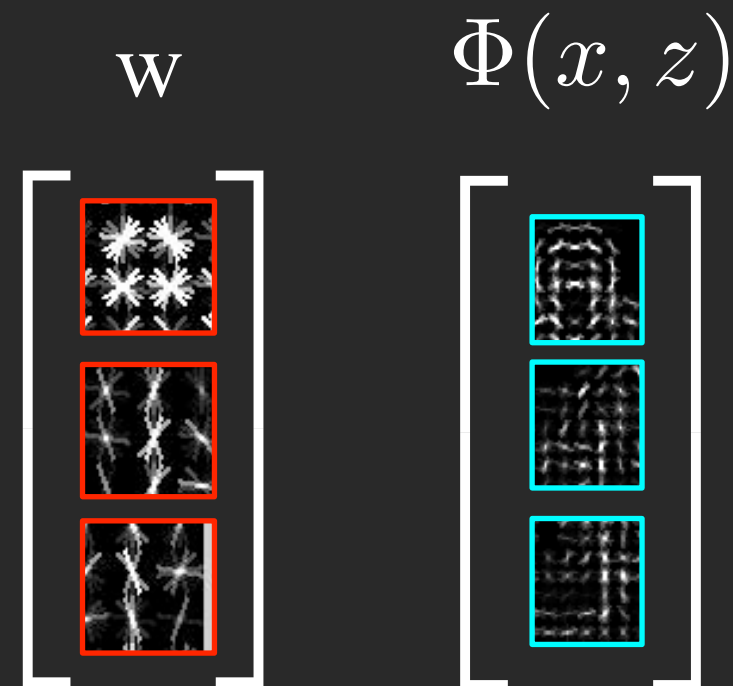(often the scoring function includes an additional "spring term"; let's ignore for now)

# Alternative formulations

$$S(x, z) = \sum_i \boxed{w_i} \cdot \boxed{\phi(x, z_i)}$$

w(z)   $\Phi(x)$



w   $\Phi(x, z)$



$$S(x, z) = w(z) \cdot \Phi(x)$$

$$S(x, z) = w \cdot \Phi(x, z)$$

[Useful for visualizing model]   [Useful for learning model parameters]

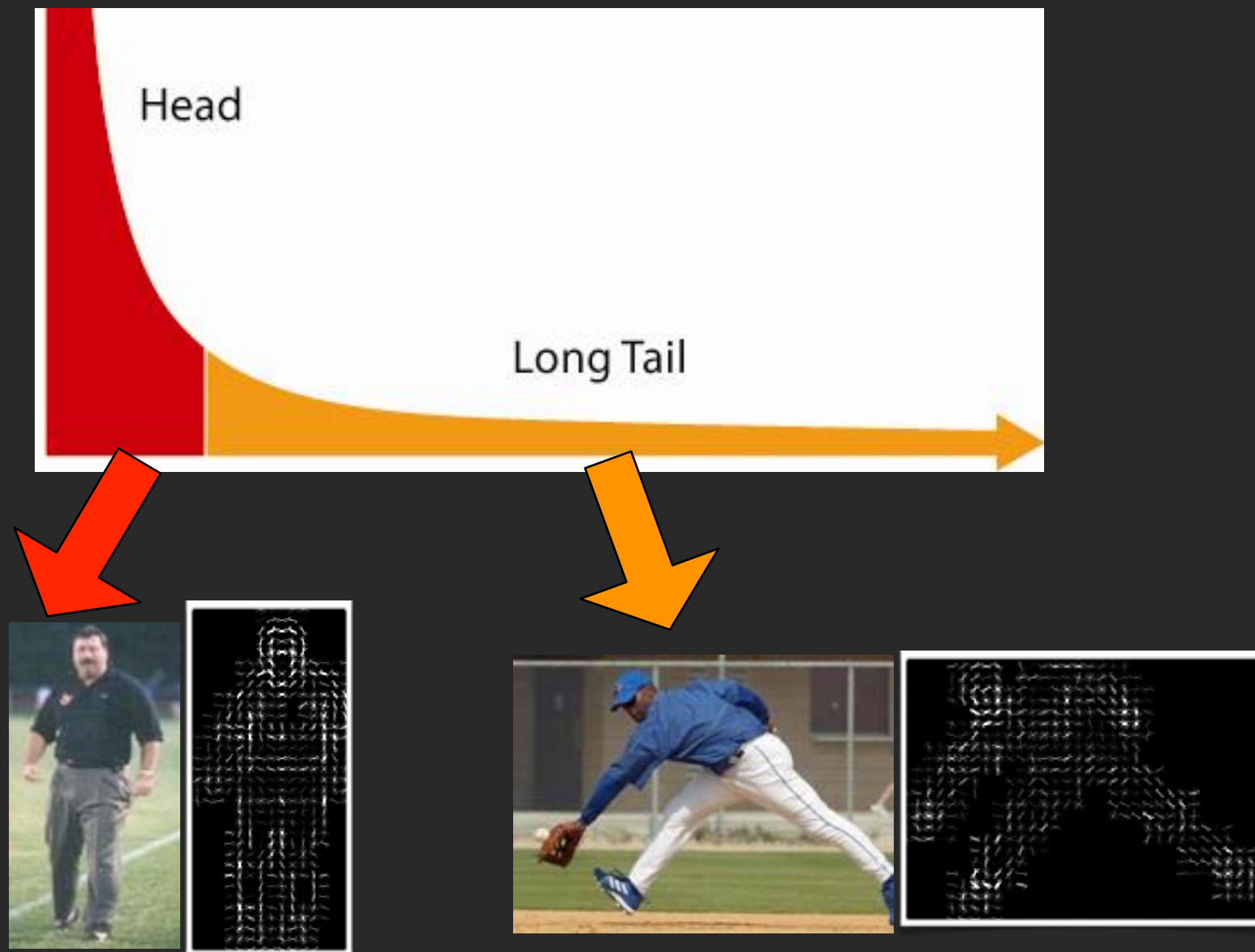# Visualizing family of classifiers

$$S(x, z) = w(z) \cdot \Phi(x)$$



How do we define set of valid $z \in \Omega$ ?

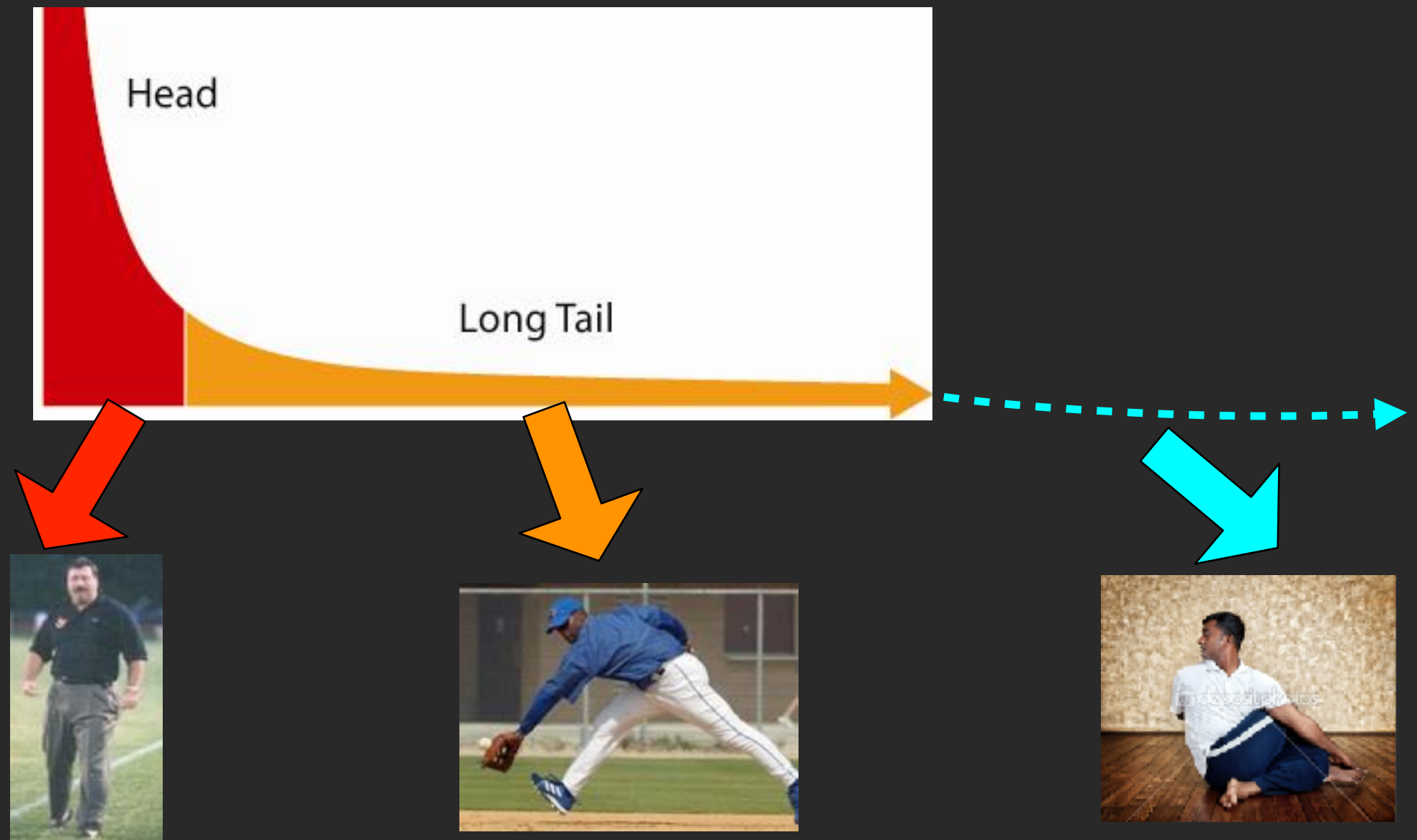One option: just use set of poses observed in training set

# Sharing

Helps address "one-shot" learning (subcategory seen at least once)


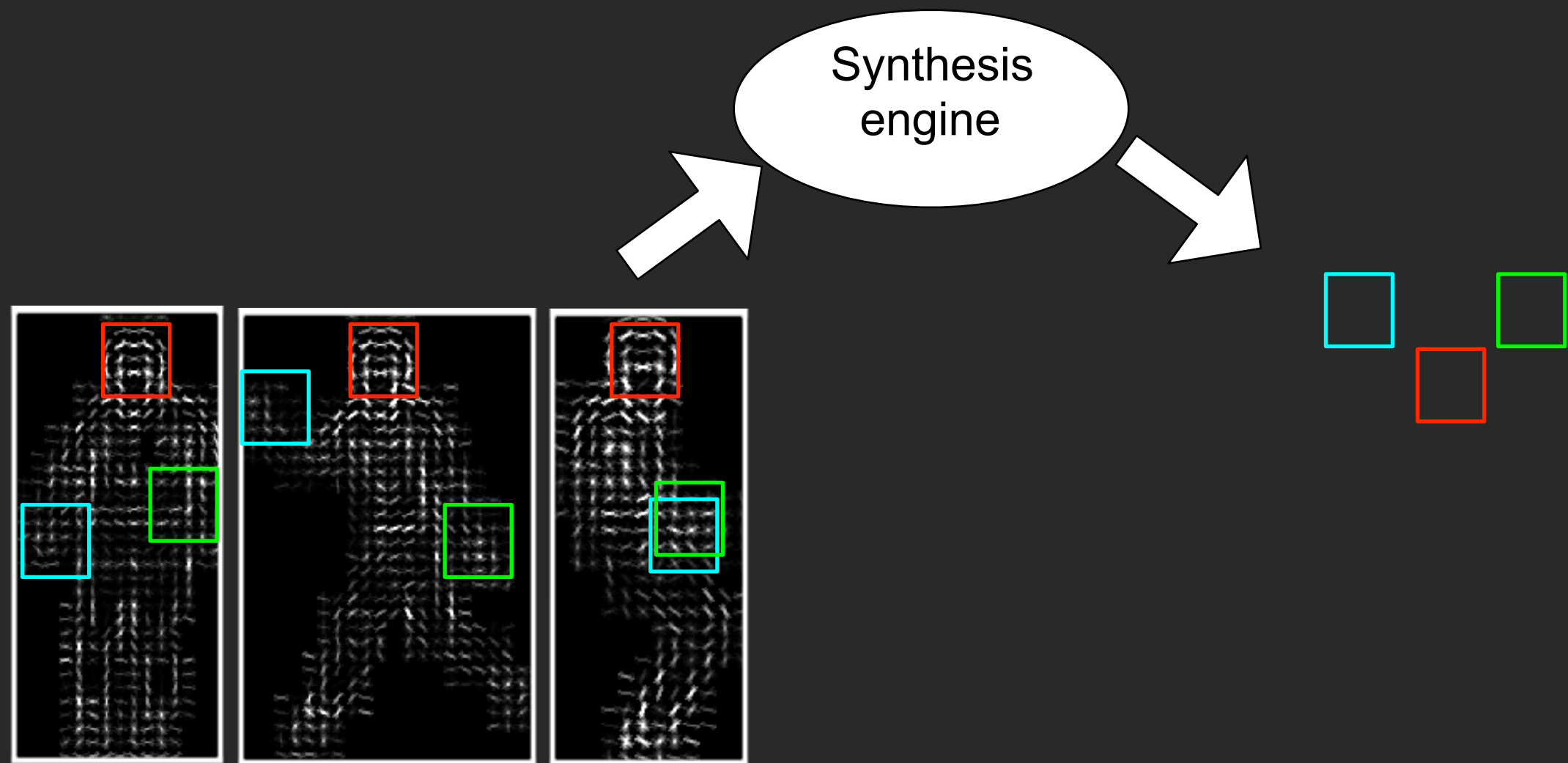
Use parts from common poses to help model rare poses

# Sharing

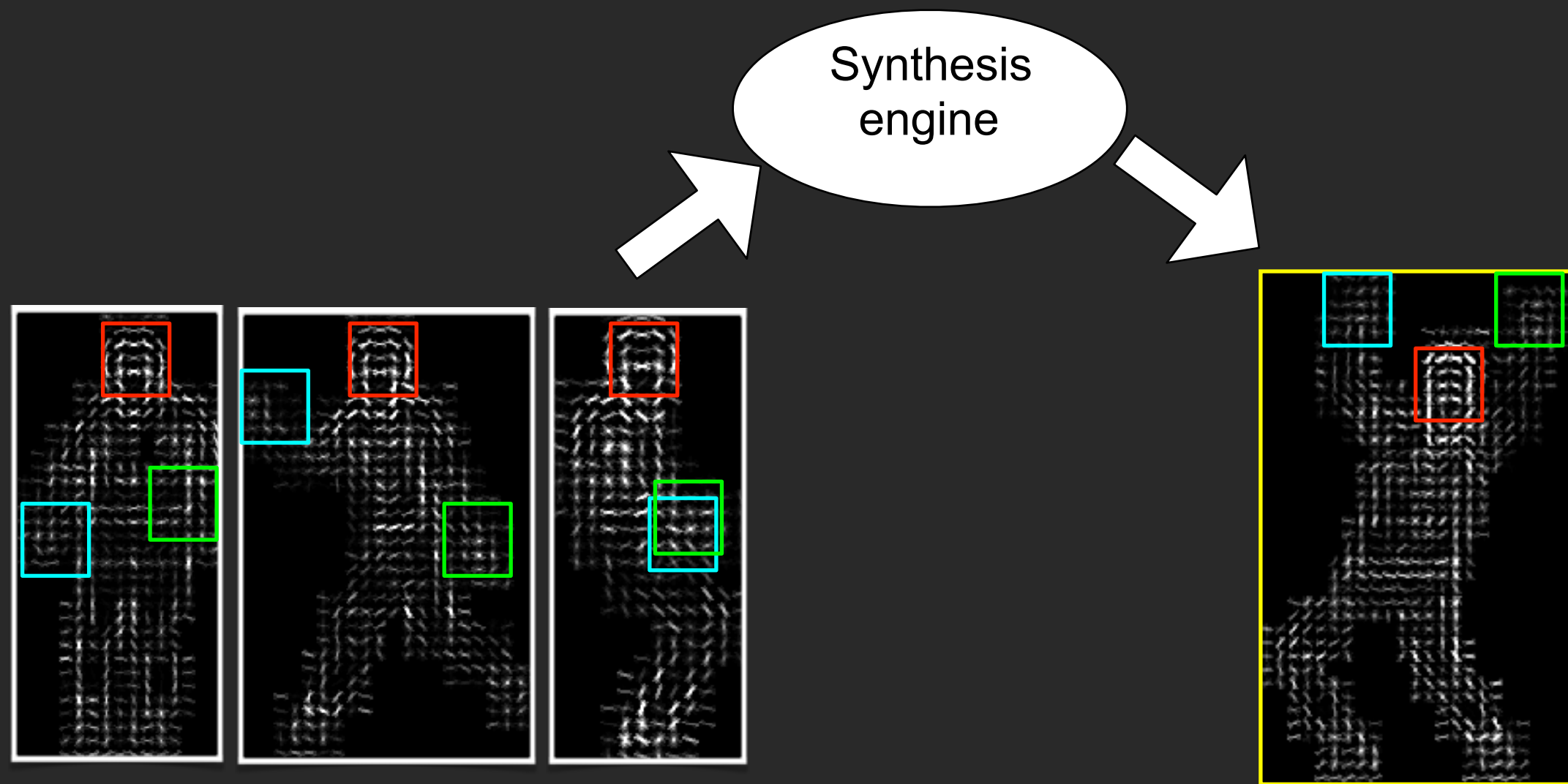Helps address "one-shot" learning (subcategory seen at least once)



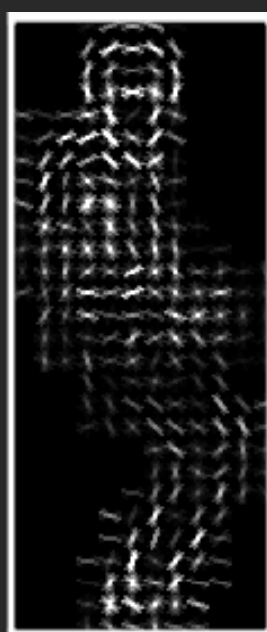What about poses that are never seen ("zero-shot" learning)?
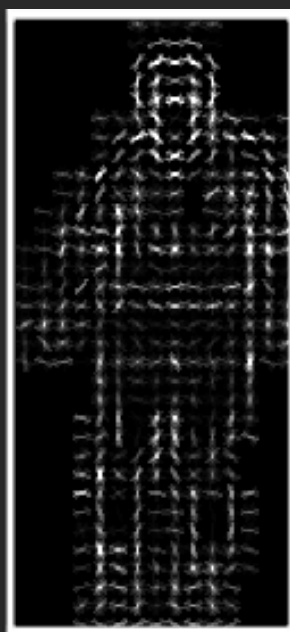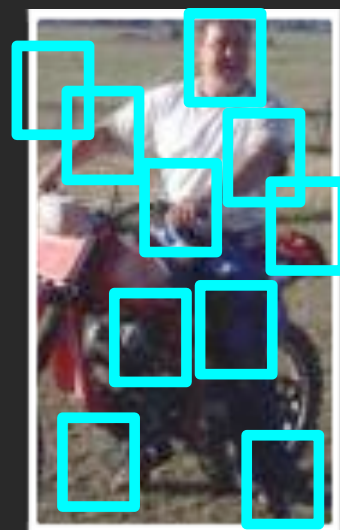
# Shape synthesis

# Shape synthesis

# Algorithmic synthesis

# Learned mo

$\Omega$ : set of observed + synthesized part locations
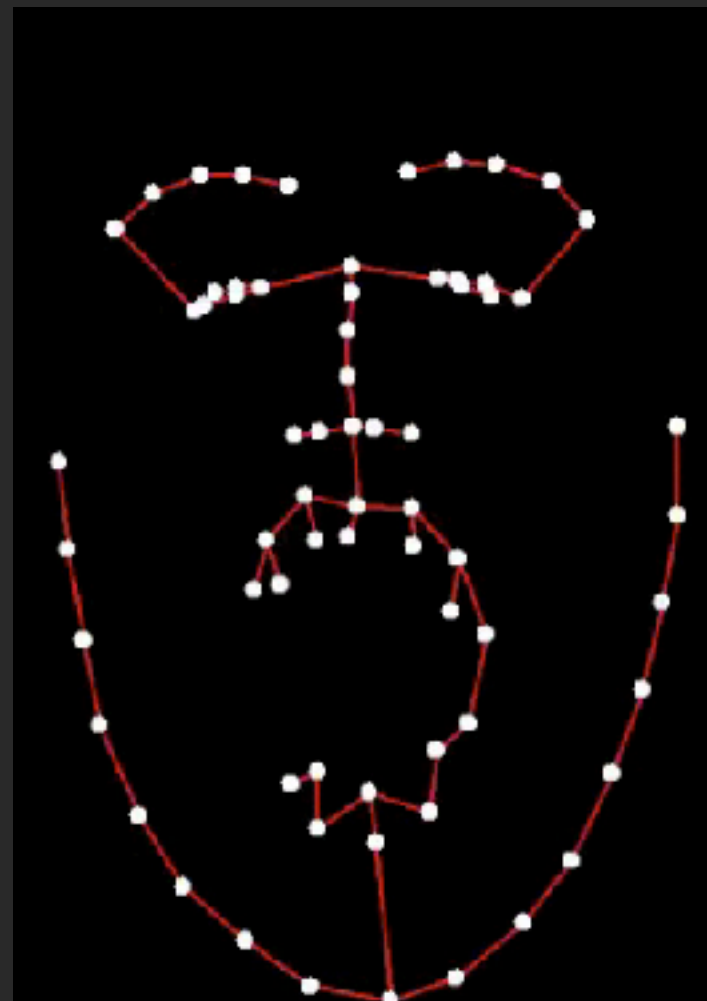
$$f_w(x) = w \cdot \Phi(x)$$

...

# Shape synthesis

$$S(z) = (z - \mu)\Sigma^{-1}(z - \mu)$$



Graphics engine

# Parametric family of classifiers

# Recognition



$$f_w(x) = w \cdot \Phi(x)$$

$$f(x) = w \cdot x$$



$\Phi(x, z)$

$$f(x) > 0$$

$$f_w(x) = \max_z w \cdot \Phi(x, z)$$

# Recognition as reconstruction



$$f(x) > 0$$

$$f(x) = w \cdot x$$

$$f(x) = \max_{z \in \Omega} w(z) \cdot x$$

$$\Phi(x, z)$$

$\Omega$   : set of observed + synthesized part locations

Argmax (z*) reveals pose

$$f(x) = w \cdot x$$

$$f(x) = \max_{z \in \Omega} w(z) \cdot x$$

Score is linear in x

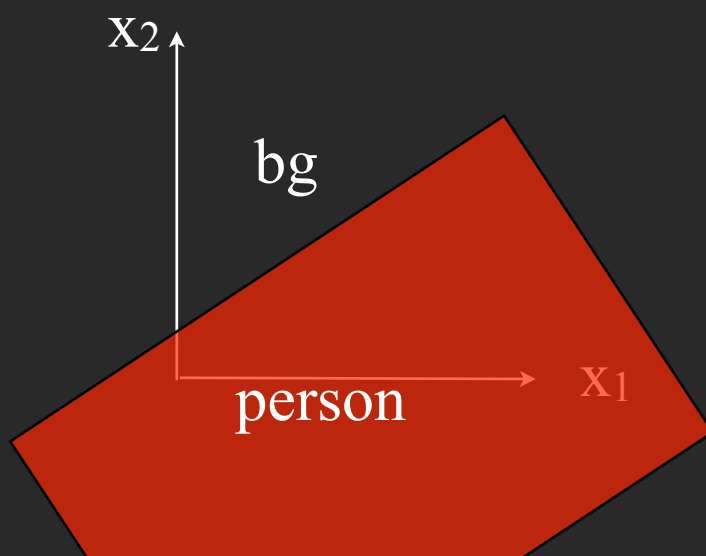Positive set $\{x : f_w(x) > 0\}$ is half-space

Score is ?

Positive set is ?

# Revisit latent (vs linear) classification



$$f(x) = w \cdot x$$

Score $f_w(x)$ is linear in x

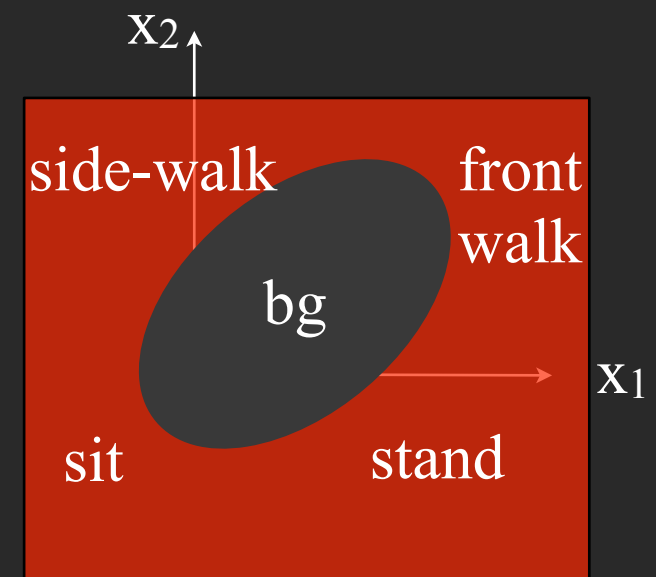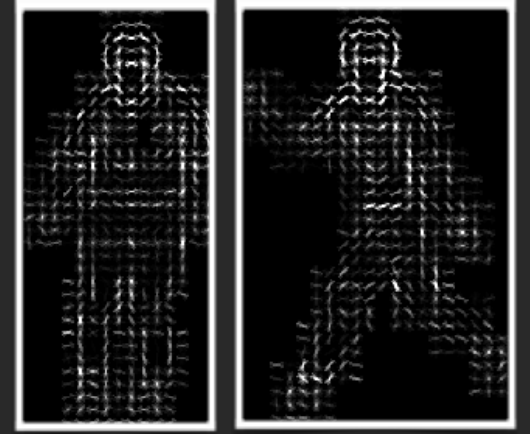Positive set $\{x : f_w(x) > 0\}$ is half-space

$$\Phi(x, z)$$

$$f(x) = \max_{z \in \Omega} w(z) \cdot x$$

Score $f_w(x)$ is convex in x

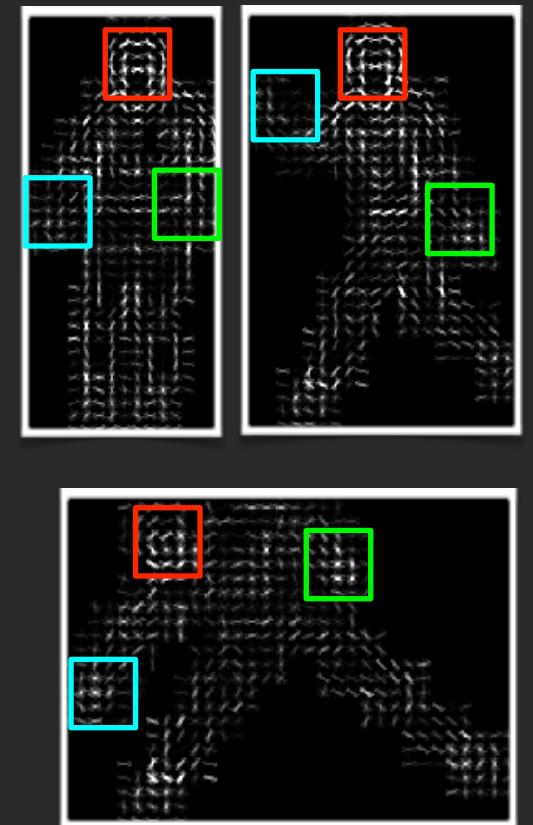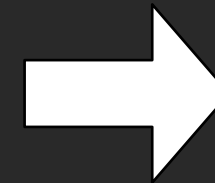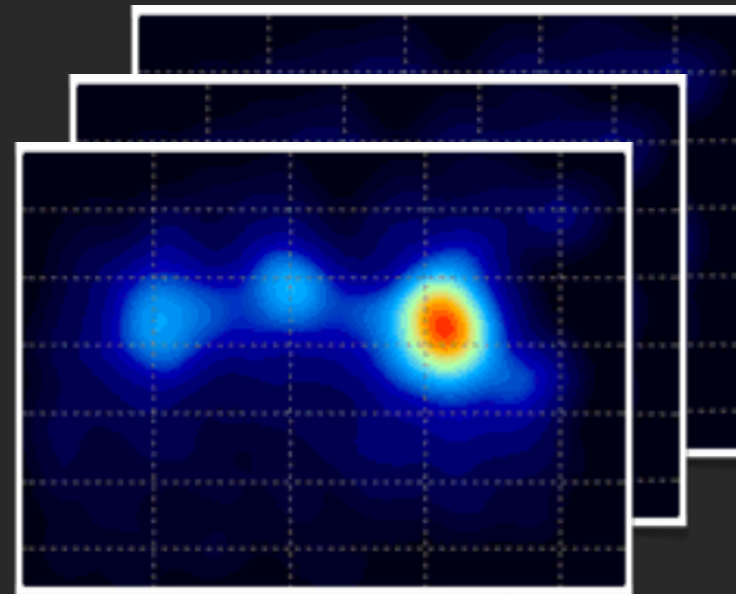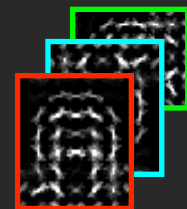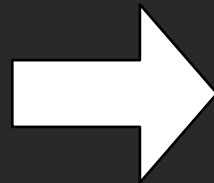Negative set $\{x : f_w(x) <= 0\}$ is convex

$$\Phi(x, z)$$

# Inference

# Inference



(1) Pre-compute tables of part responses

(2) Score each template with lookup table (LUT) queries

Can be implemented as a two-layer convolution
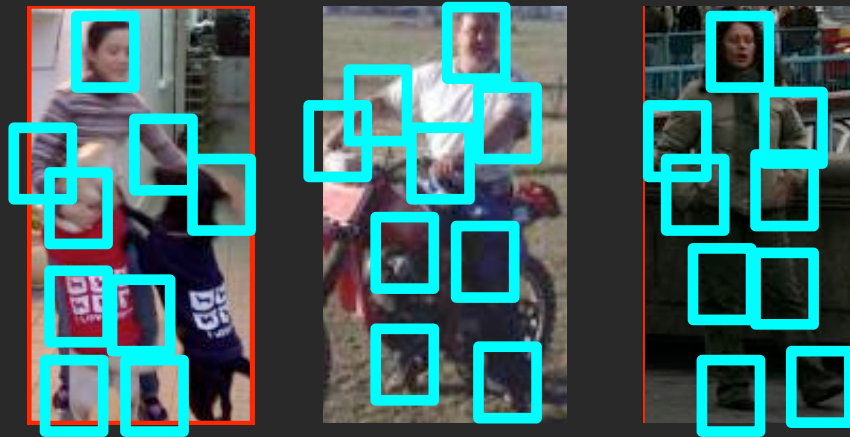
# Learning

$$S(x, z) = \sum_i \boxed{w_i} \cdot \boxed{\phi(x, z_i)}$$

$$w = \begin{bmatrix} \ \\ \ \\ \ \\ \ \end{bmatrix} \qquad \Phi(x, z) = \begin{bmatrix} \ \\ \ \\ \ \\ \ \end{bmatrix}$$

$$S(x, z) = w \cdot \Phi(x, z)$$

# Supervised lear

$$S(x, z) = w \cdot \Phi(x, z),$$

pos

neg

# Learned me
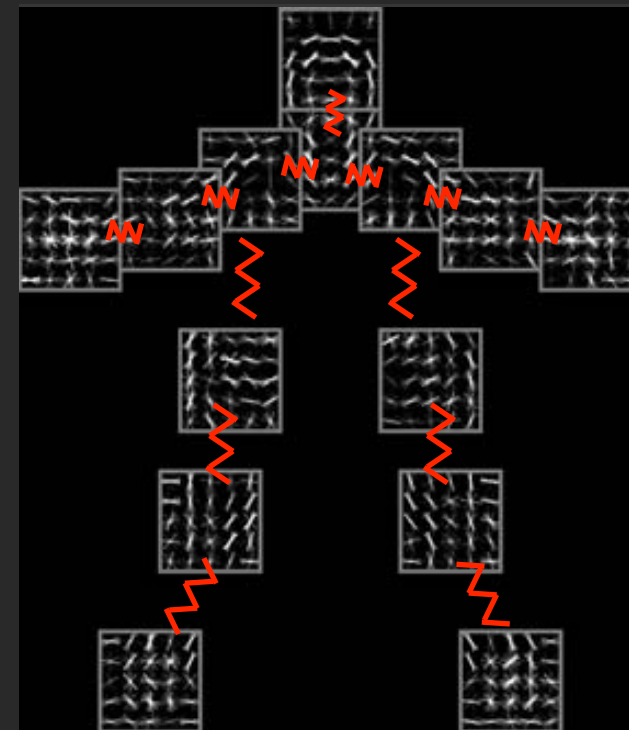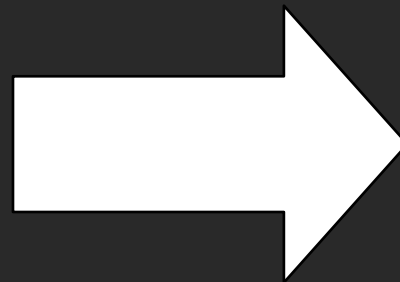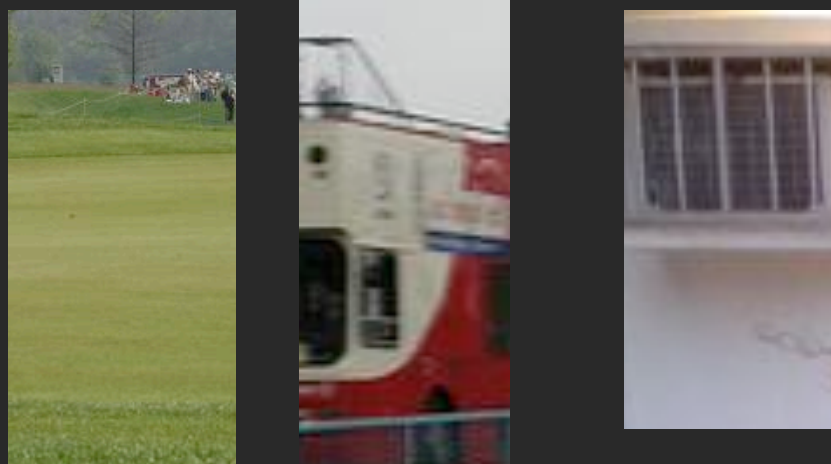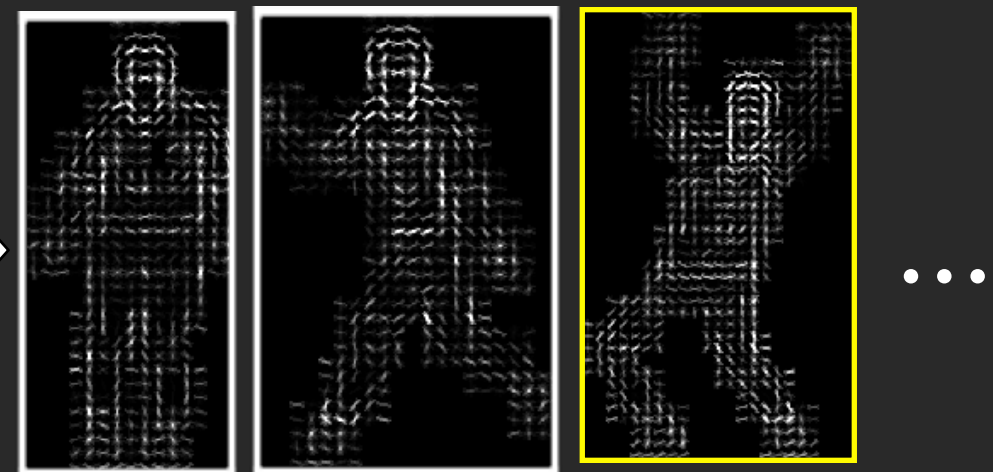
Supervised learning
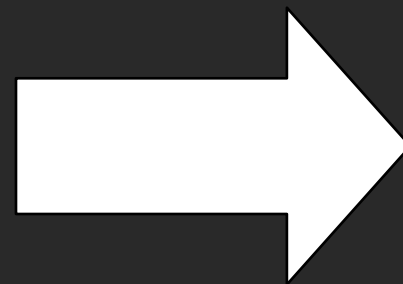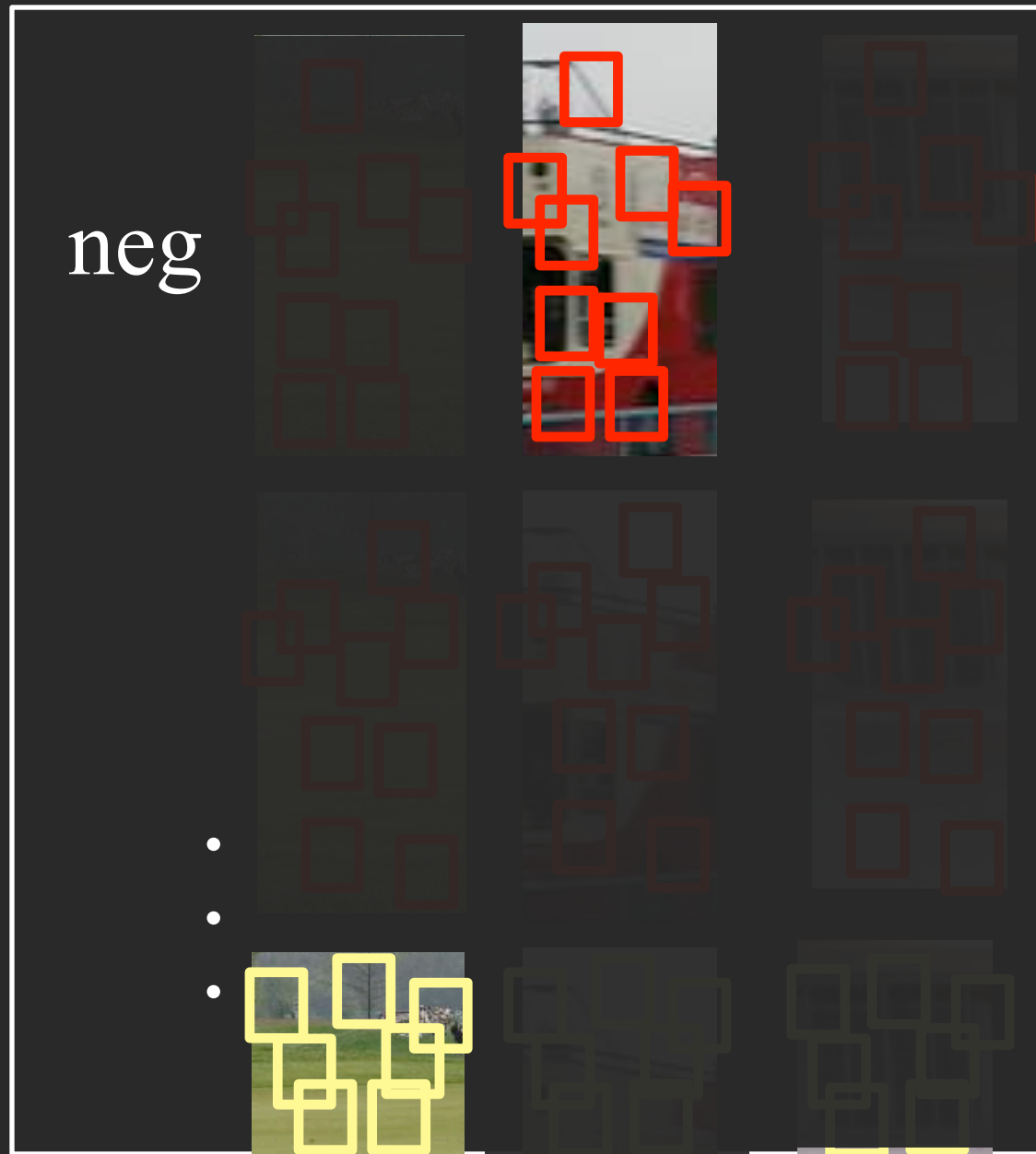$$S(x, z) = w \cdot \Phi(x, z), \quad z \in \Omega$$

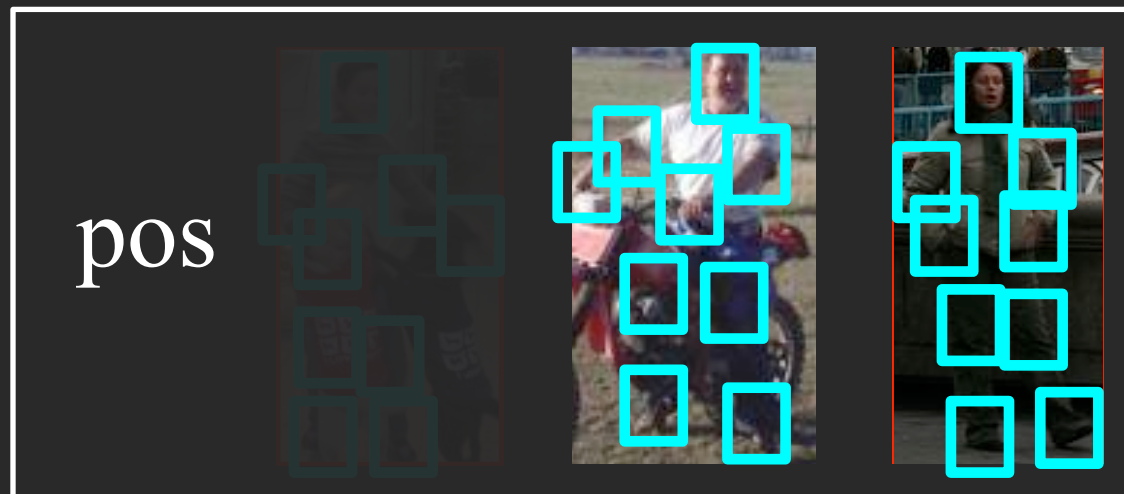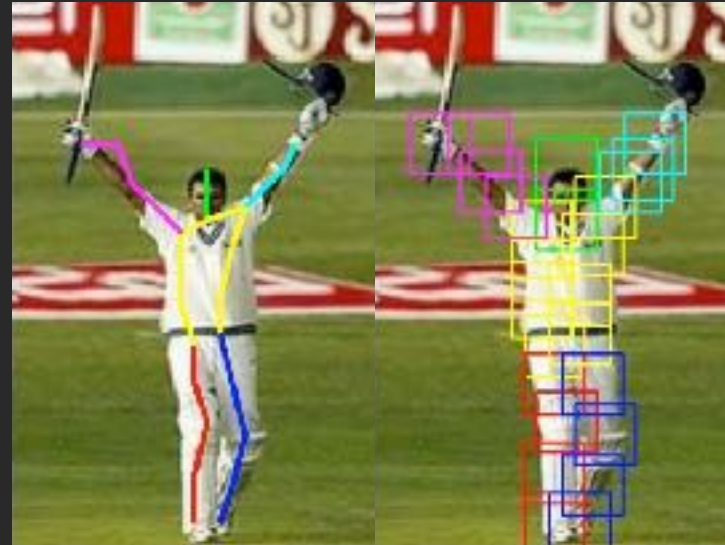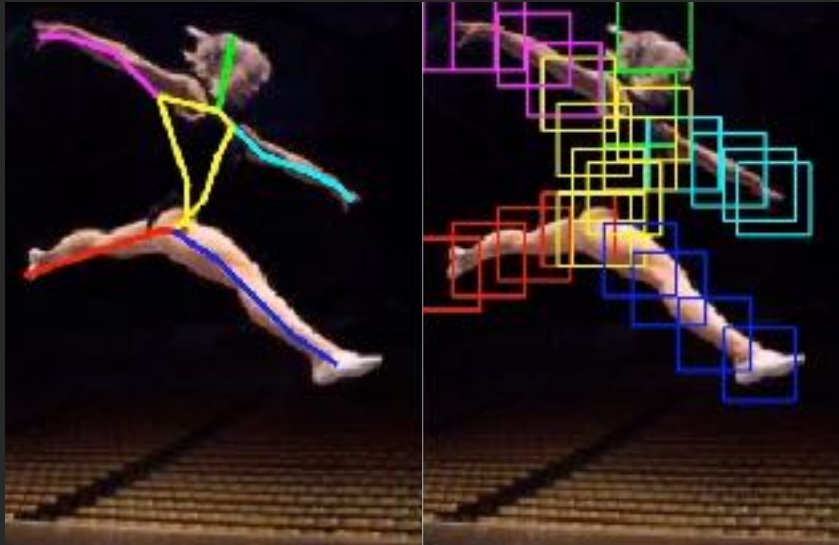$$f_w(x) = w \cdot \Phi(x)$$

pos



neg



$\vdots$

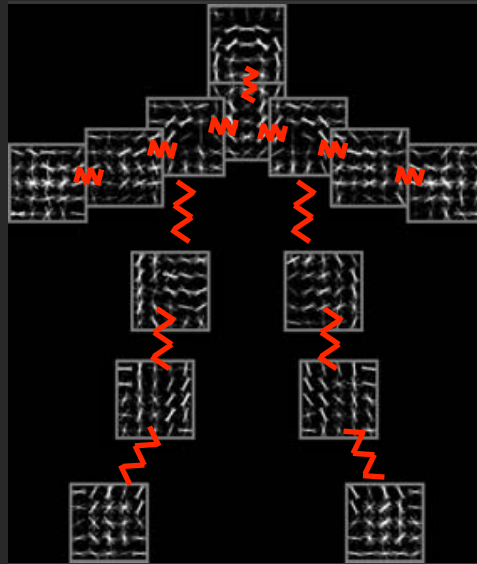Learn classifiers for never-before-seen templates with synthesis

(Apply same sparse learning tricks to deal with large set of negatives!)

# Joint recognition + (2D) reconstruction

# Implicit



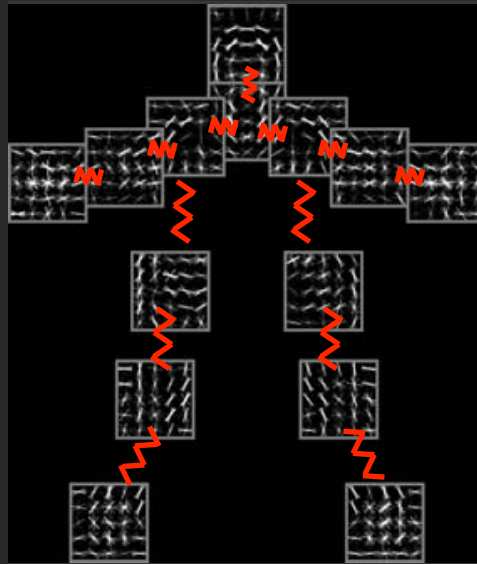Assume K parts and L candidate locations: $|\Omega| = L^K$



$$\max_{z \in \Omega} (x, z_i)$$

Do we really need to search over

# Implicit

Assume K parts and L candidate locations: $|\Omega| = L^K$



$$\max_{z \in \Omega} \sum_i w_i \cdot \phi(x, z_i)$$

Do we really need to search over $\Omega$?
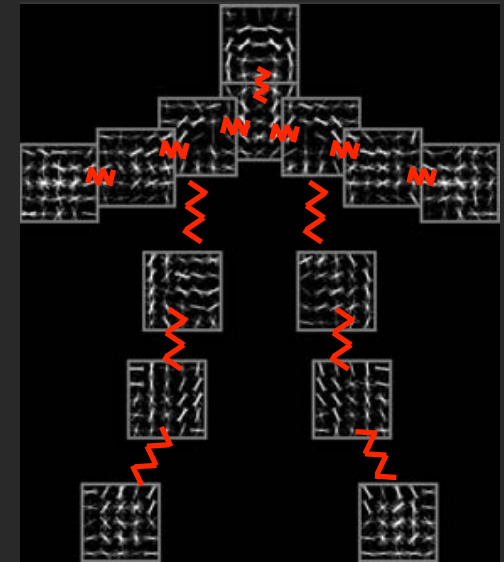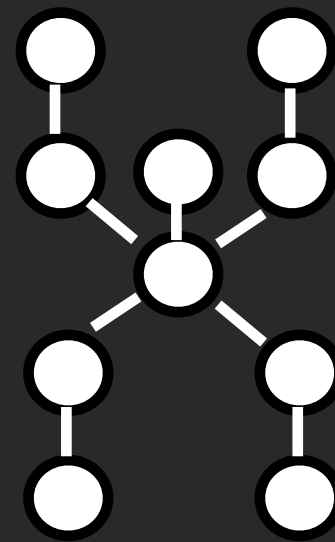
$$= \sum_i \max_{z_i} w_i \cdot \phi(x, z_i)$$

No! Independantly find best location of each part. Allows us to *implicitly* synthesize $\Omega$

# Generalize approach to markov models
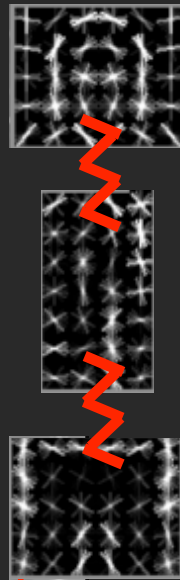
temporal markov model               spatial mark

- For each candidate torso, independently estimate best arm and
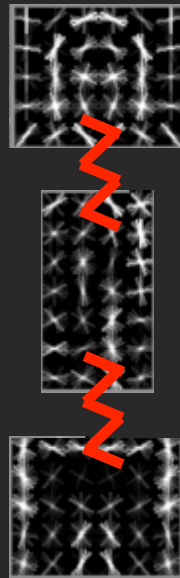- Allows us to model (and learn) priors over exponentially-larg

$$S(x, z) = \sum_{i \in V} w_i \cdot \phi(x, z_i) + \sum_{ij \in E} w_{ij} \cdot \psi(z_i, z_j)$$

# General case



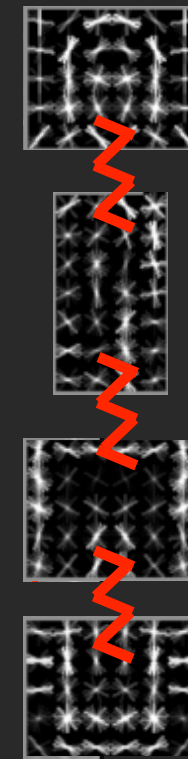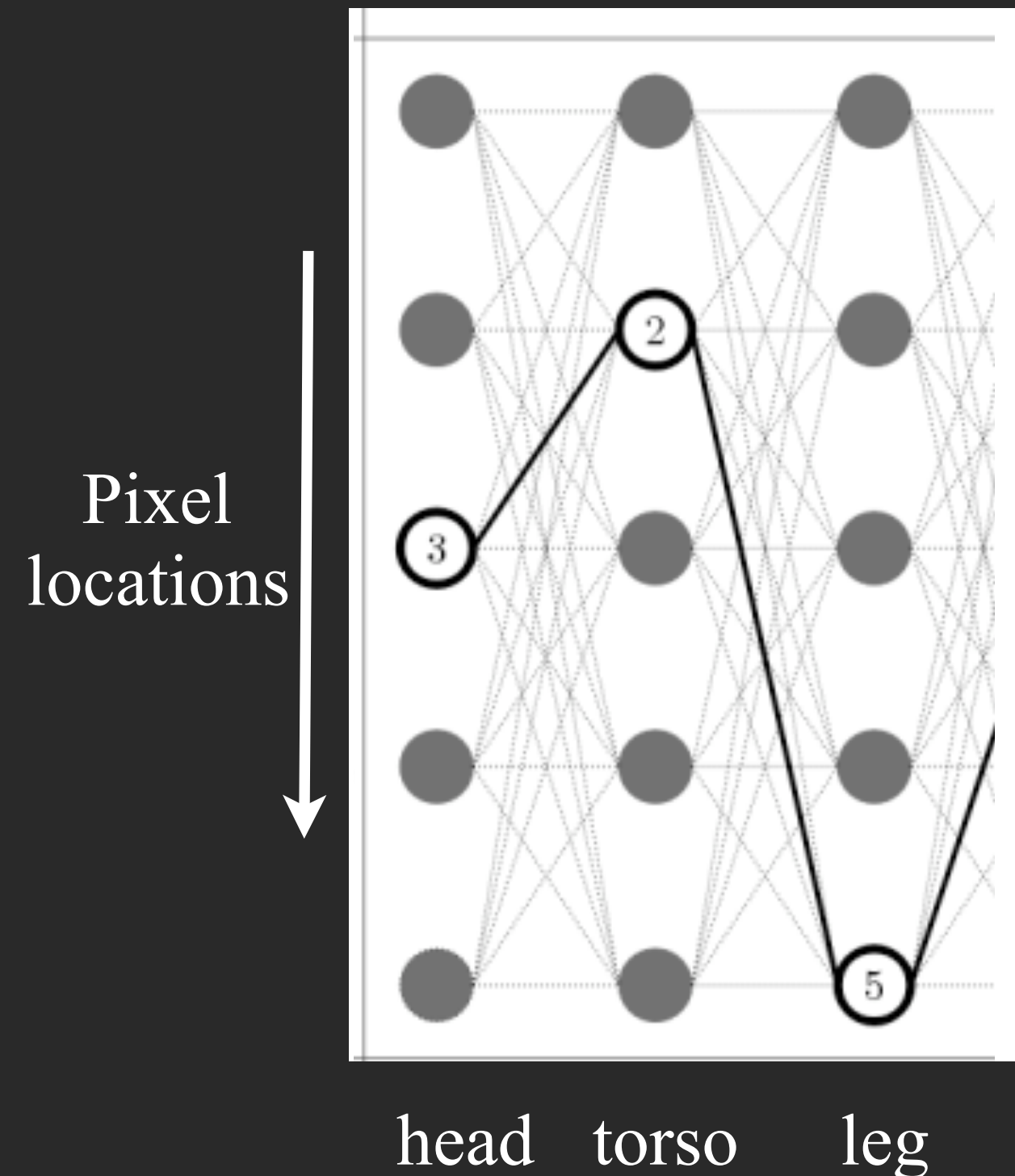$$\max_{z_1, z_2, z_3} \left[ \phi(z_1) + \phi(z_2) + \phi(z_3) + \psi(z_1, z_2) + \psi(z_2, z_3) \right]$$

# General case



$$\max_{z_1, z_2, z_3} \left[ \phi(z_1) + \phi(z_2) + \phi(z_3) + \psi(z_1, z_2) + \psi(z_2, z_3) \right]$$

$$= \max_{z_2, z_3} \left( \psi(z_2) + \psi(z_2, z_3) + \max_{z_1} \left[ \phi(z_1) + \psi(z_1, z_2) \right] \right)$$

$$= \max_{z_2, z_3} \left( \psi(z_2) + \psi(z_2, z_3) + m(z_2) \right)$$

Use simple variatble elimination to reduce inference to $O(KL^2)$

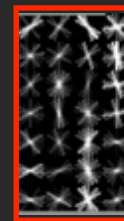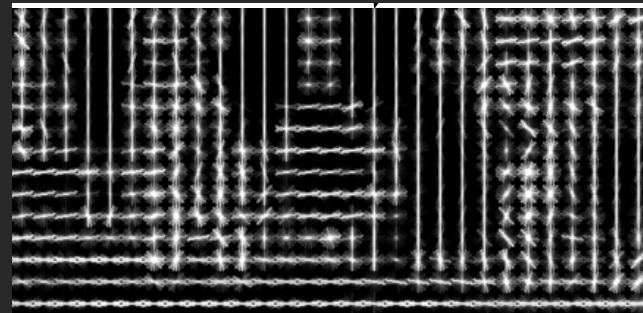# Inference: max$_z$ S(x,z)



Pixel locations

head    torso    leg

1) Initialize nodes with match score

2) Initialize edges with spring score

3) Find best path from left to right

In practice, (1) is bottleneck
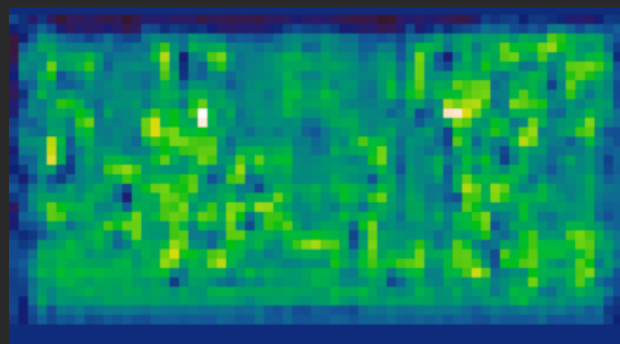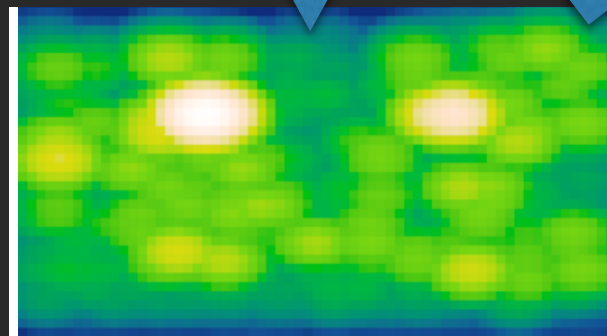
Infe...
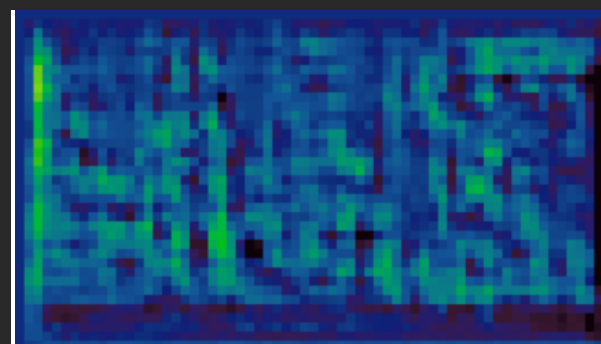
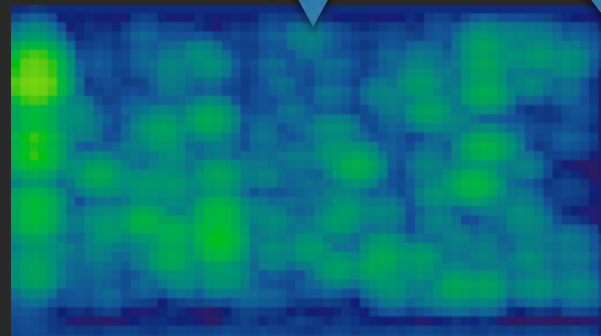implement with convolutions + max pool
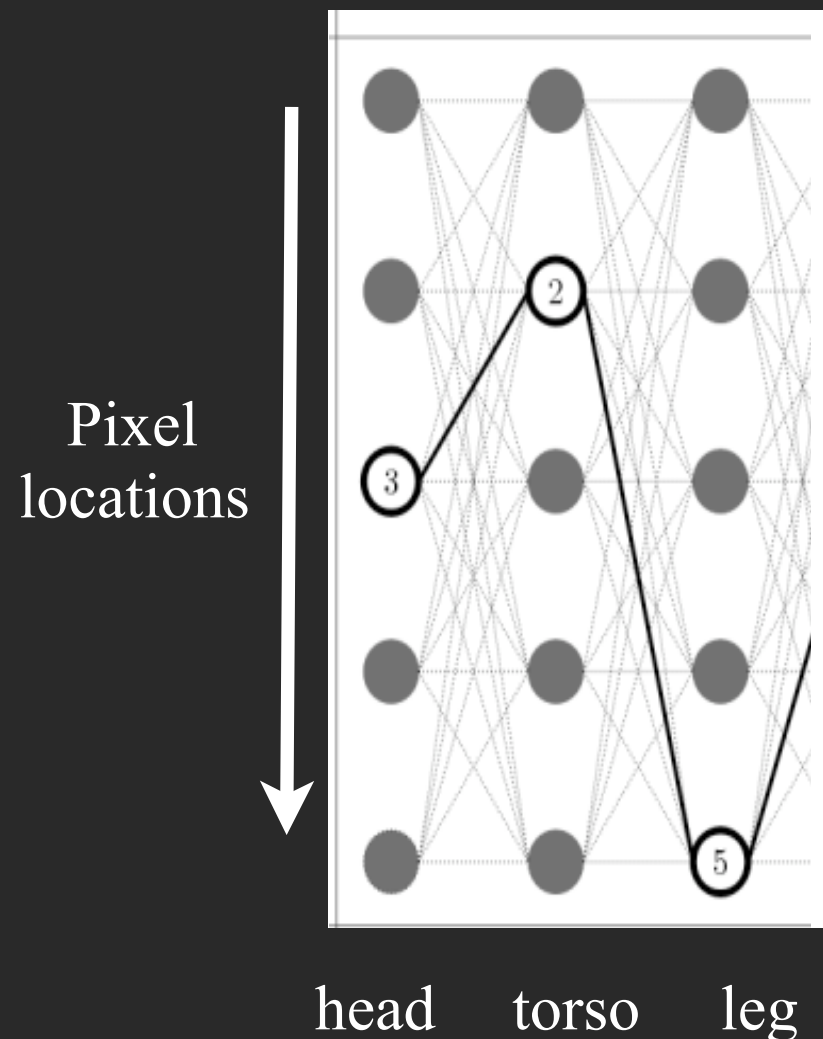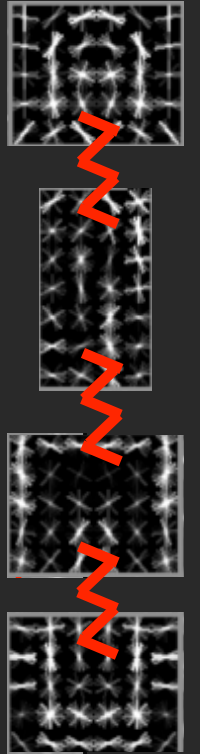
conv

conv

max pool

add

max pool

add

. . .

# General formulation: inference

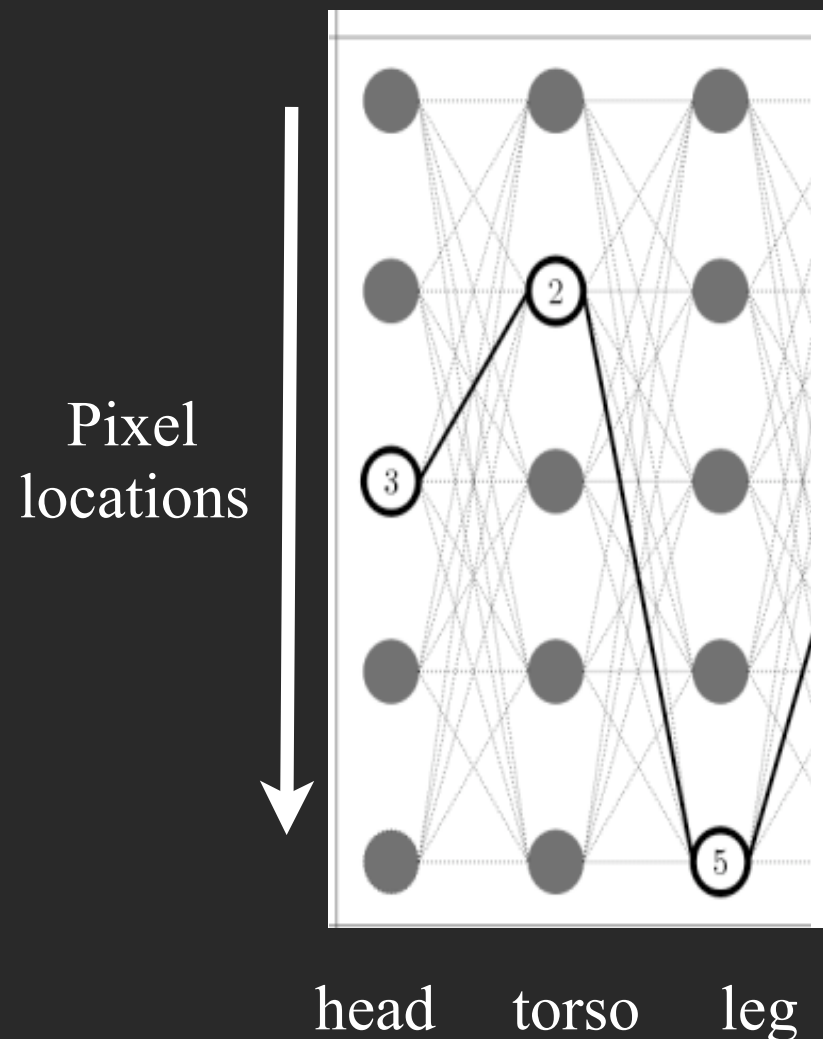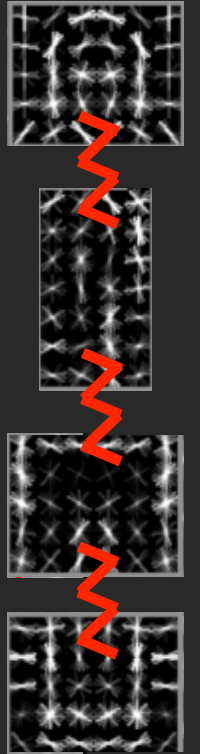$$S(x, z) = \sum_i \phi_i(z_i, x) + \sum_{ij \in E} \psi_{ij}(z_i, z_j, x)$$

Pixel locations

Local and pairwise potentials can be arbitrary nonlinear functions of image and position

(e.g., neural net part model)

(e.g., intervening contour cue on part pairs)

head    torso    leg

# General formulation: inference

$$S(x, z) = \sum_i \phi_i(z_i, x) + \sum_{ij \in E} \psi_{ij}(z_i, z_j, x)$$



Pixel
locations

Local and pairwise potentials can be arbitrary
nonlinear functions of image and position

(e.g., neural net part model)

(e.g., intervening contour cue on part pairs)

head    torso    leg

# General formulation: learning

$$S(x, z) = \sum_{i \in V} w_i \cdot \phi(x, z_i) + \sum_{ij \in E} w_{ij} \cdot \psi(z_i, z_j)$$

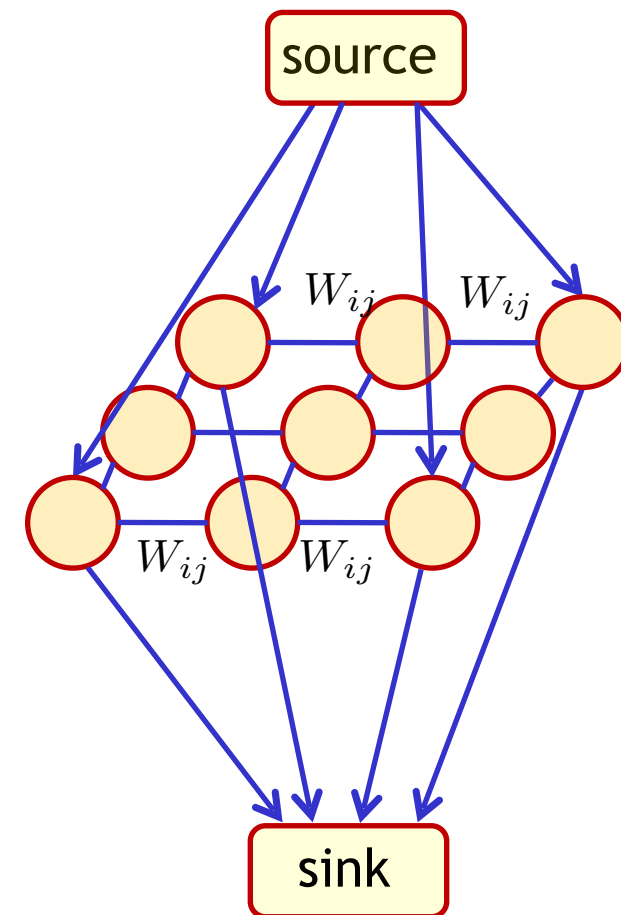$$w = \begin{bmatrix} \ \\ \ \\ \ \\ \ \end{bmatrix} \qquad \Phi(x, z) = \begin{bmatrix} \ \\ \ \\ \ \\ \ \end{bmatrix}$$

$$S(x, z) = w \cdot \Phi(x, z)$$

Jointly learn appearance and geometry with linear classification!

# Aside: structural SVM learning



$$\{(x_n, y_n)\}$$

$$x_n \in R^{H \times W}$$

$$y_n \in \{0, 1, \ldots K\}^{H \times W}$$

$$E(x, y) = w \cdot \Phi(x, y)$$

$$E(x, y) = \sum_{i \in V} w_{local} \cdot \phi(x, y_i) + \sum_{ij \in \mathcal{E}} w_{pair} \cdot \psi(y_i, y_j, x)$$

# Modeling faces

# Face detection

# Face detection



Taigman & Wolf "Leveraging Billions of Faces..." 2011

Picasa

face.com®

part model
(hundreds of faces)

OpenCV

precision

recall
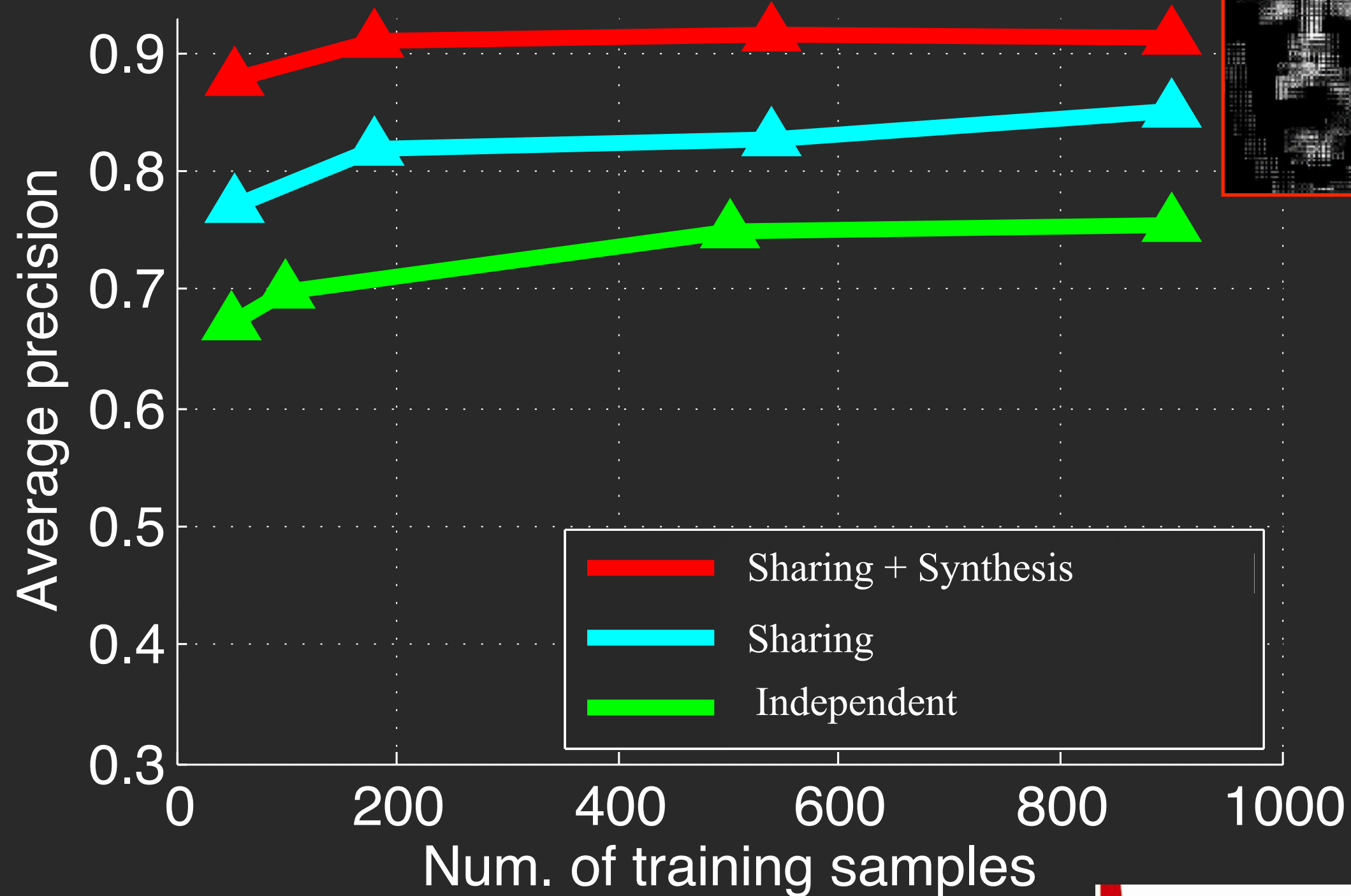
# Could we do better with more data?



Zhu, Vondrick, Ramanan & Fowlkes,
"Do we need more training data or better models?"
IJCV 2015 (accepted)

# Indepedant subcategories vs sharing



Zhu, Vondrick, Ramanan & Fowlkes,
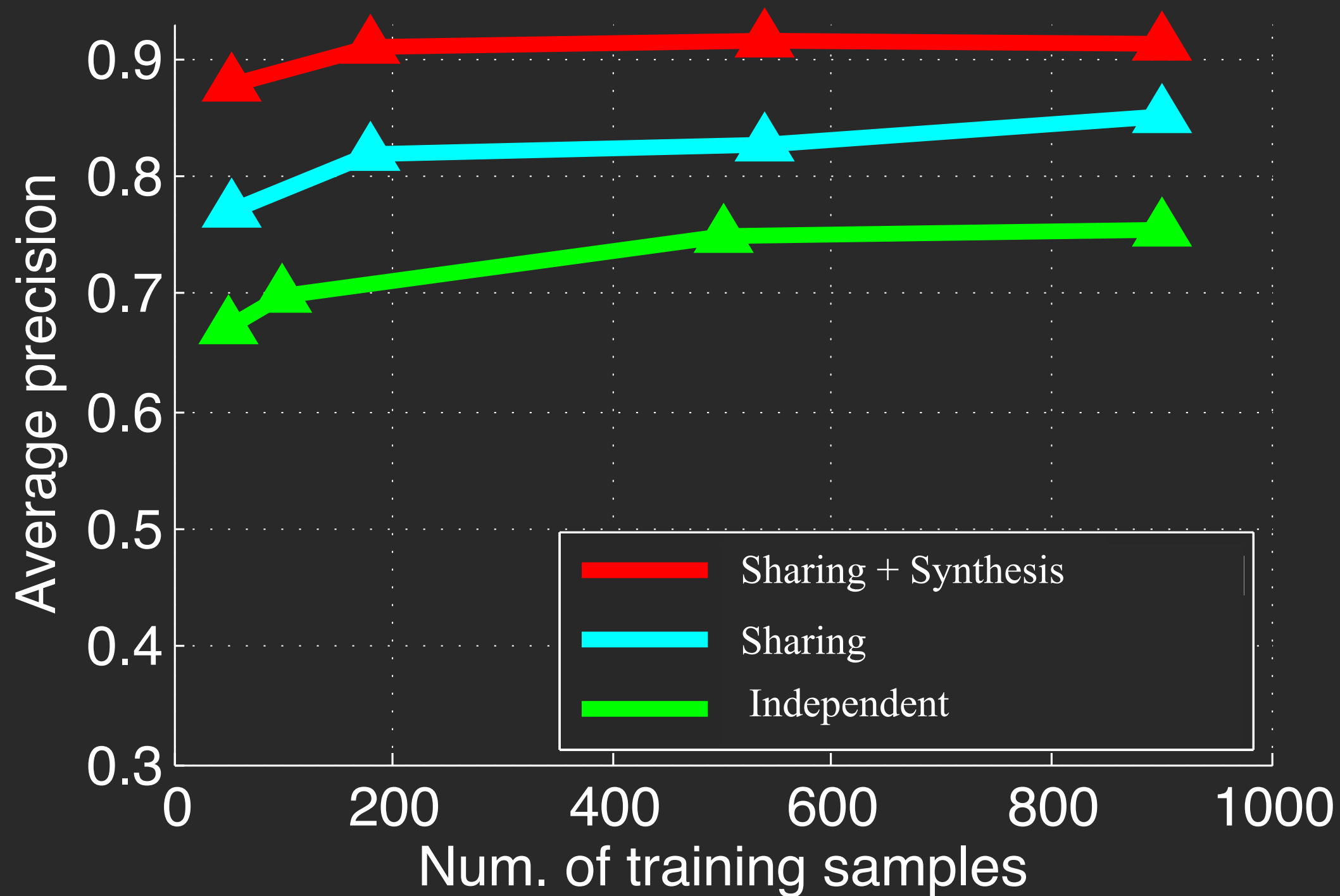"Do we need more training data or better models?"
IJCV 2015 (accepted)

Sharing versus synthesis

Average precision (y-axis): 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9

Num. of training samples (x-axis): 0, 200, 400, 600, 800, 1000

Legend:
- Sharing + Synthesis (red)
- Sharing (cyan)
- Independent (green)

Synthesis even more beneficial than sharing

Head

Long Tail

Sharing versus synthesis

One can train a state-of-art face detector (*c.f.* Google Picassa & Facebook's face.com) with 100 faces!

# What if we want to recognize 3D shapes?



Input:
2D image

Output:
3D shape
camera viewpoint

# Shape synthesis



Problem... local appearances depend on global geometry
(foreshortened or occluded wheels look very different)

# Geometry-conditioned appearance

Enlarge part dictionary to include parts with different local geometries
(learned by clustering)

# Shape and appearance synthesis

Shape Synthesis

"where"

Appearance Synthesis

"what"

# Explicit set of synthesized templates
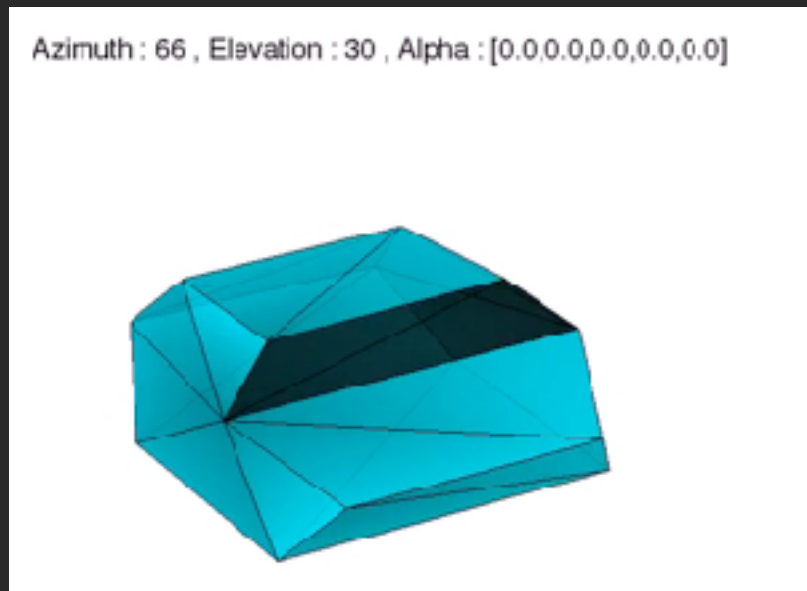


Azimuth : -70 , Elevation : 0 , f : 0.125 Alpha : [0,0,0,0]

(Most viewpoints never seen during training)

# Example detections

# Car detection + reconstruction
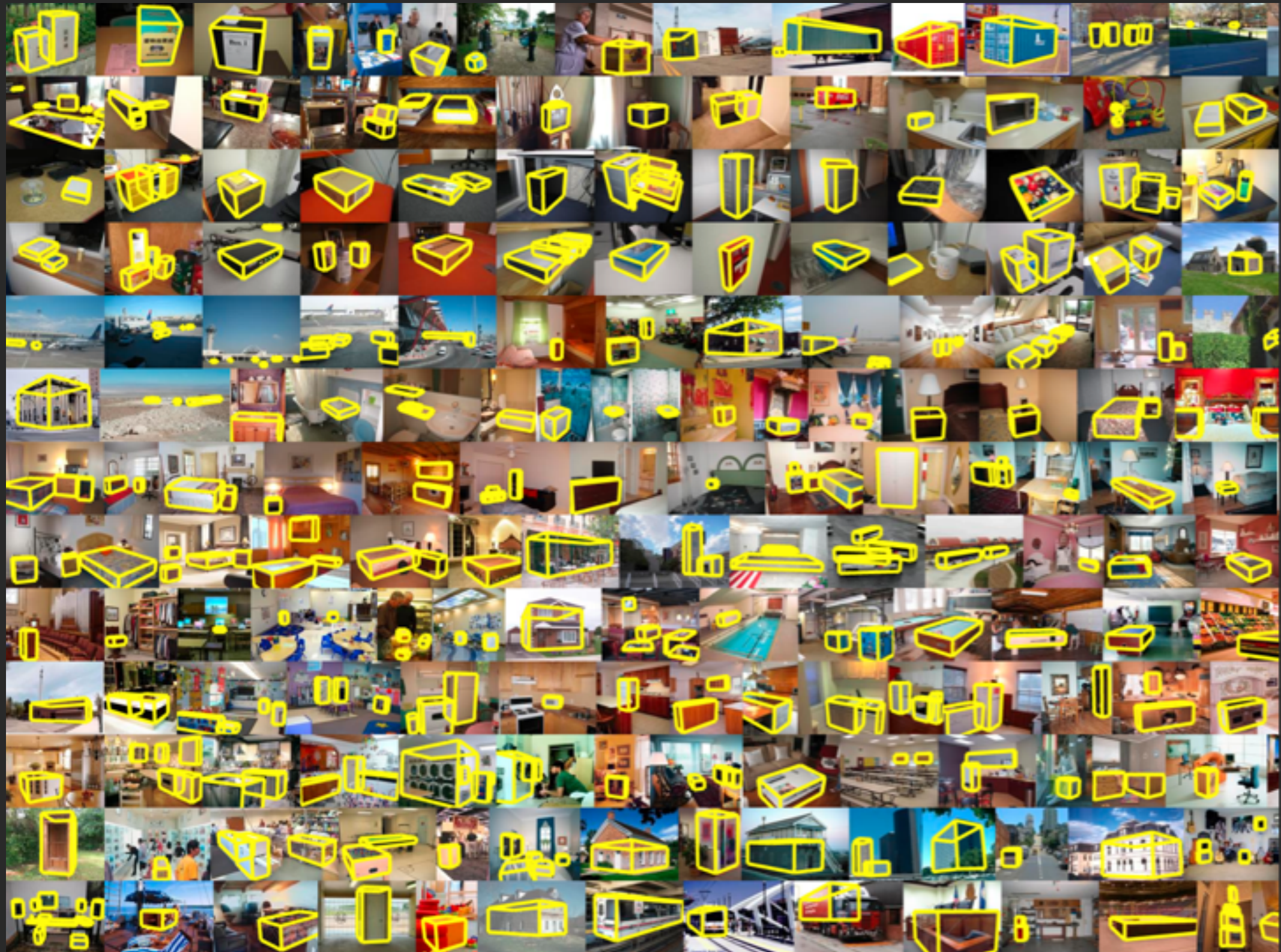


Azimuth : 66 , Elevation : 30 , Alpha : [0.0,0.0,0.0,0.0,0.0,0.0]

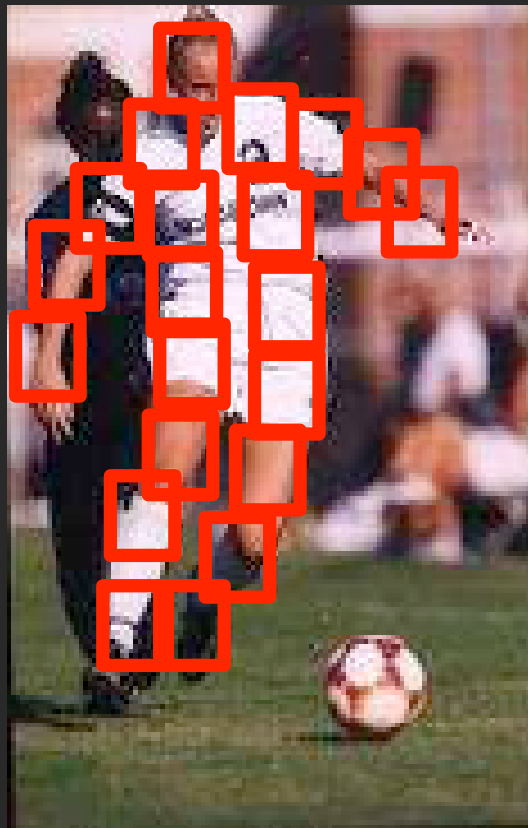parts are beign spatially moved *and* swapped out

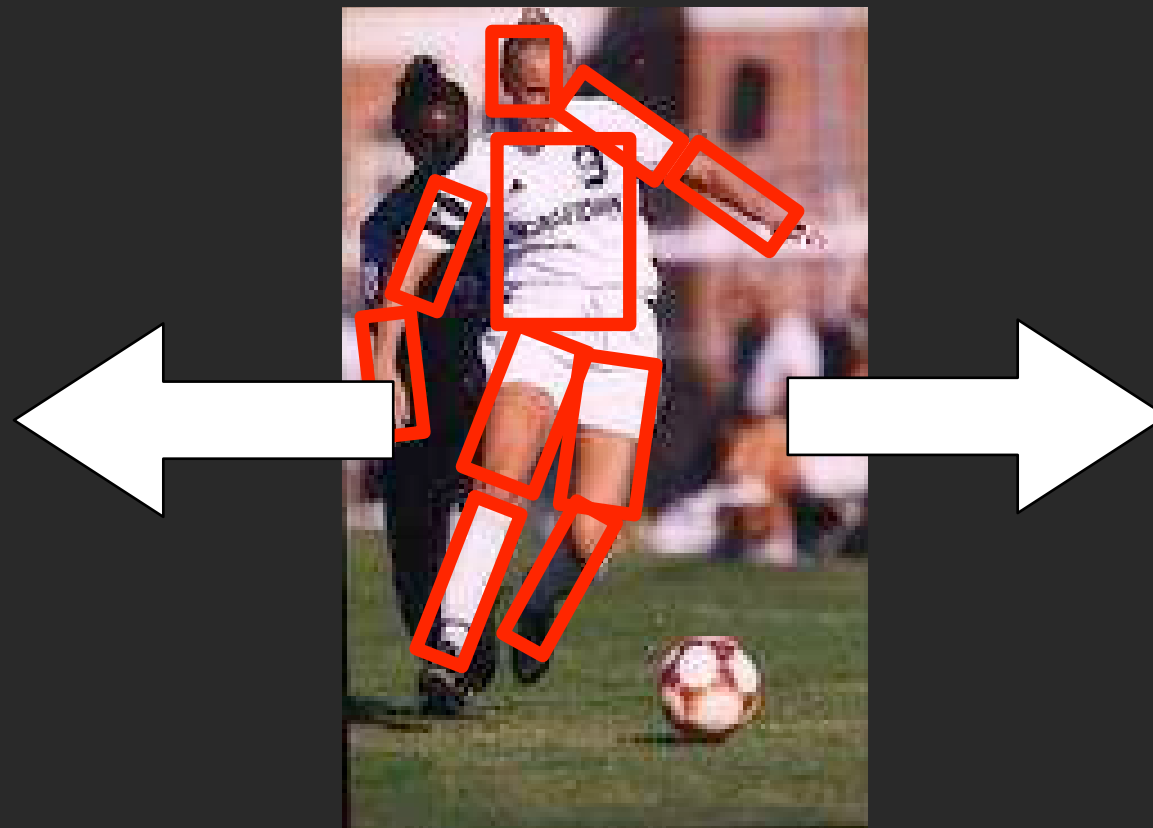# Evaluation: MIT Geometric Primitive Dataset
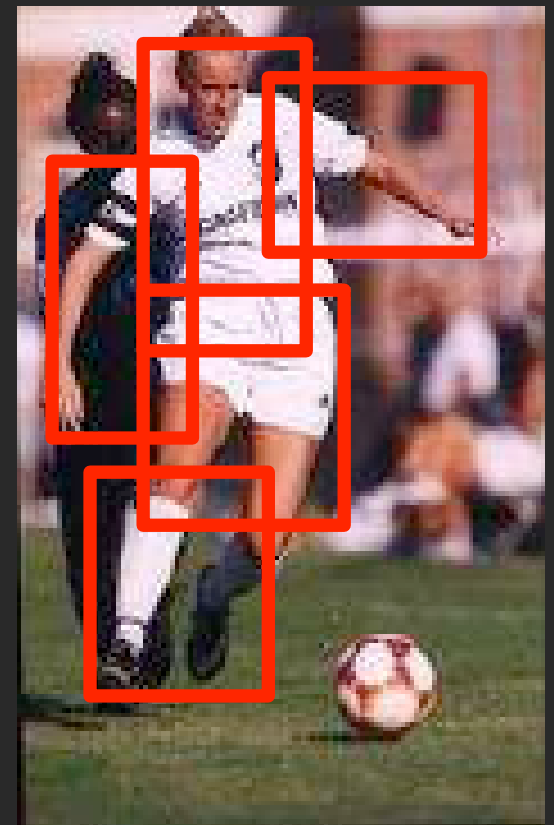
# Semantic vs learned representations

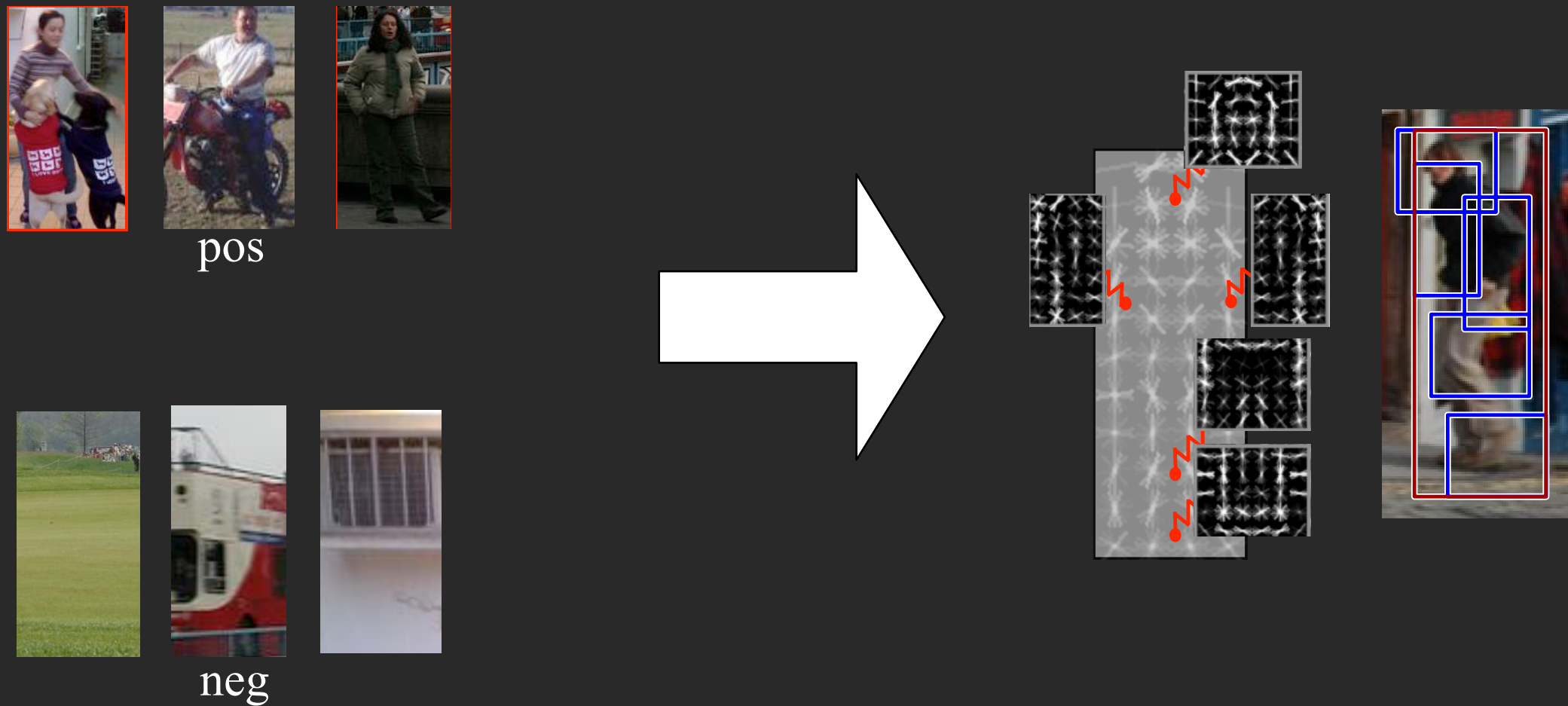What are the right units for sharing and synthesis?
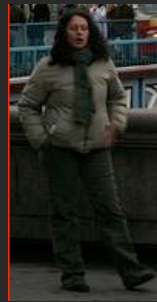


Patches

Skeleton

"Poselets"

# For general objects?

# Learned model

$$f_w(x) = w \cdot \Phi(x)$$

Learn parts that allow for accurate recognition

pos

neg

Felzenszwalb, Girshick, McAllester, and Ramanan *PAMI* 2010

# SVMs



pos

neg
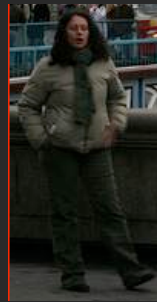
Given positive and negative training windows $\{x_n\}$

$$L(w) = ||w||^2 + \sum_{n \in \text{pos}} \max(0, 1 - f_w(x_n)) + \sum_{n \in \text{neg}} \max(0, 1 + f_w(x_n))$$
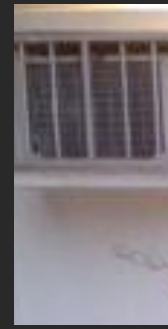
$$f_w(x) = w \cdot \Phi(x)$$

*L(w)* is convex  (Quadratic Program)

# Latent SVMs

Felzenszwalb, McAllester, Ramanan *CVPR* 2008



pos                                    neg

Given positive and negative training windows {x_n}

$$L(w) = ||w||^2 + \sum_{n \in \text{pos}} \max(0, 1 - f_w(x_n)) + \sum_{n \in \text{neg}} \max(0, 1 + f_w(x_n))$$

$$f_w(x) = \max_{z} w \cdot \Phi(x, z)$$

*L(w)* is "almost" convex

# Coordinate descent

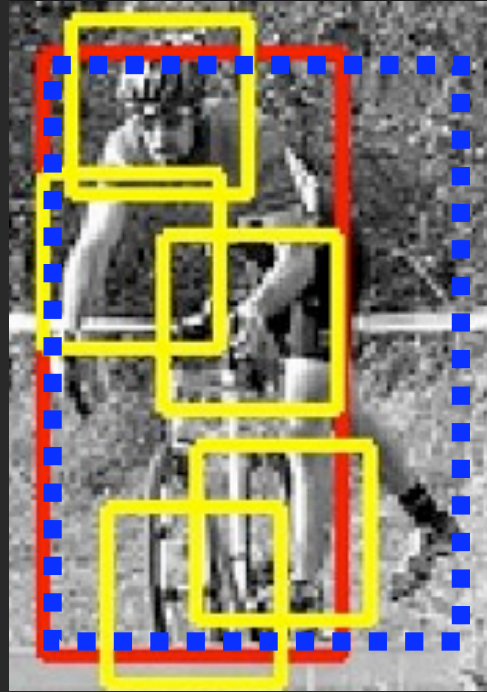1) Given positive part locations, learn w with a convex program

$$w = \underset{w}{\operatorname{argmin}} \, L(w) \quad \text{with fixed} \quad \{z_n : n \in \text{pos}\}$$

2) Given w, estimate part locations on positives

$$z_n = \underset{z}{\operatorname{argmax}} \, w \cdot \Phi(x_n, z) \quad \forall n \in \text{pos}$$
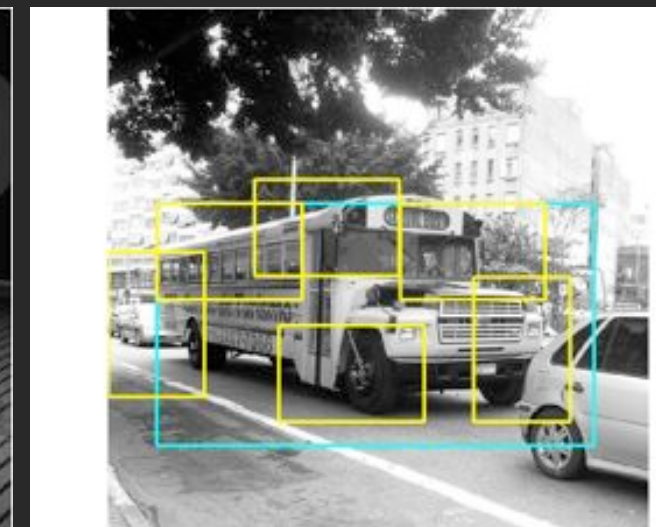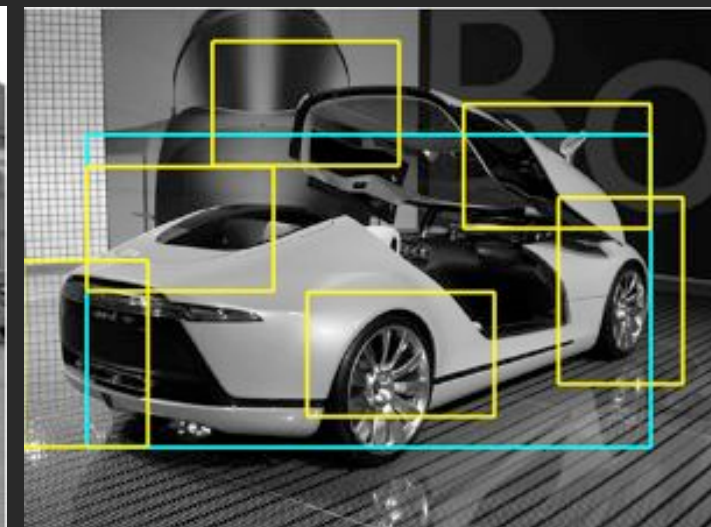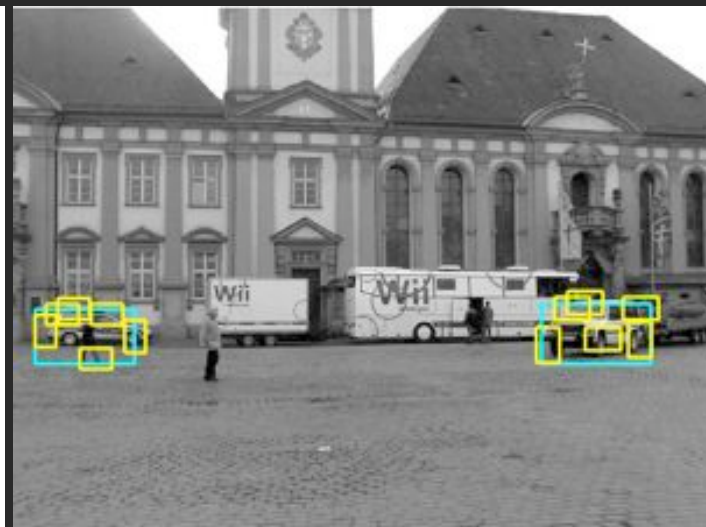
The above steps perform coordinate descent on a joint loss
Can be seen as an instance of the CCCP algorithm (Yuille)

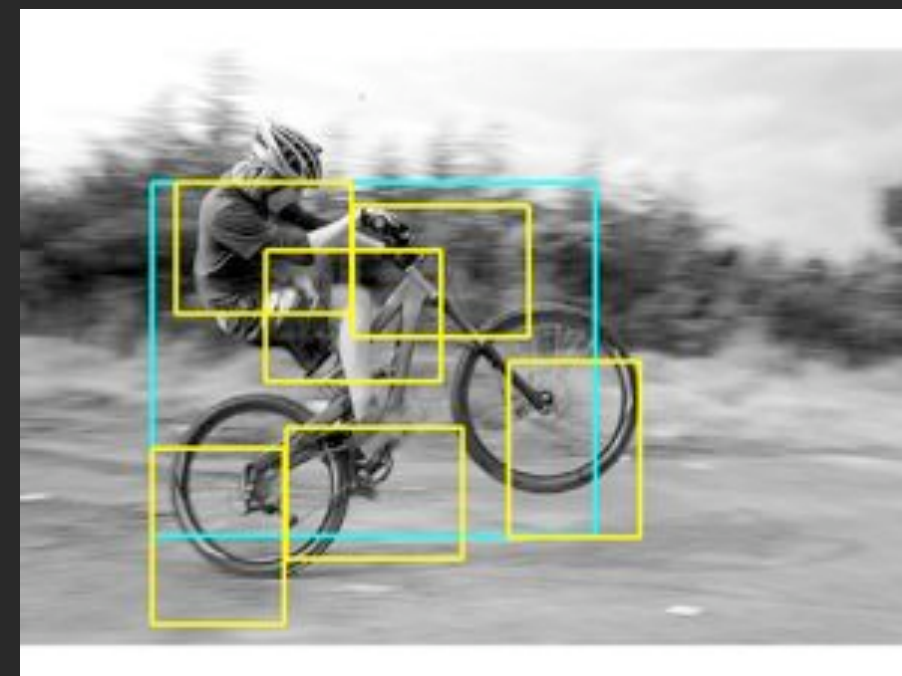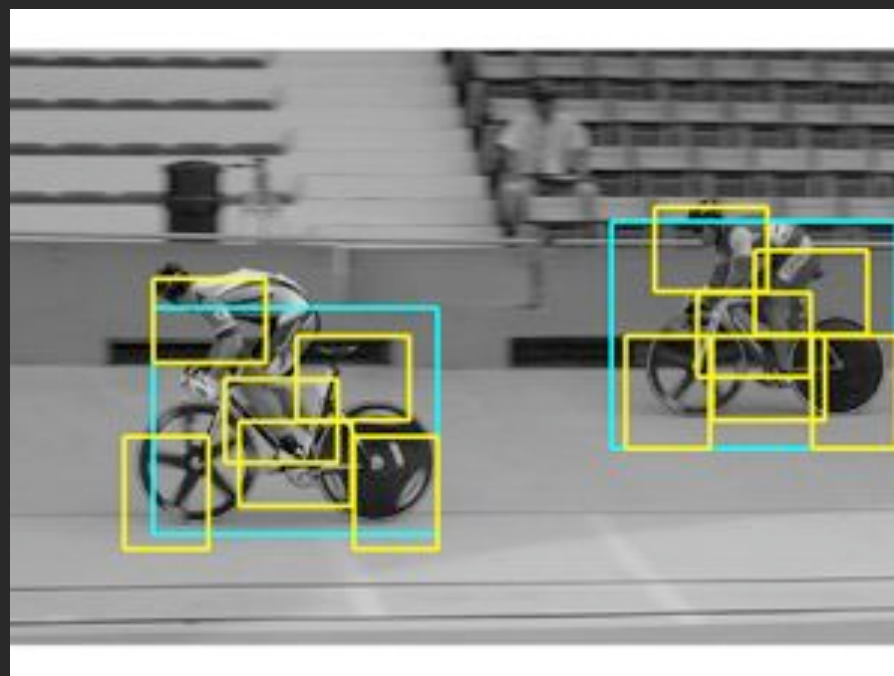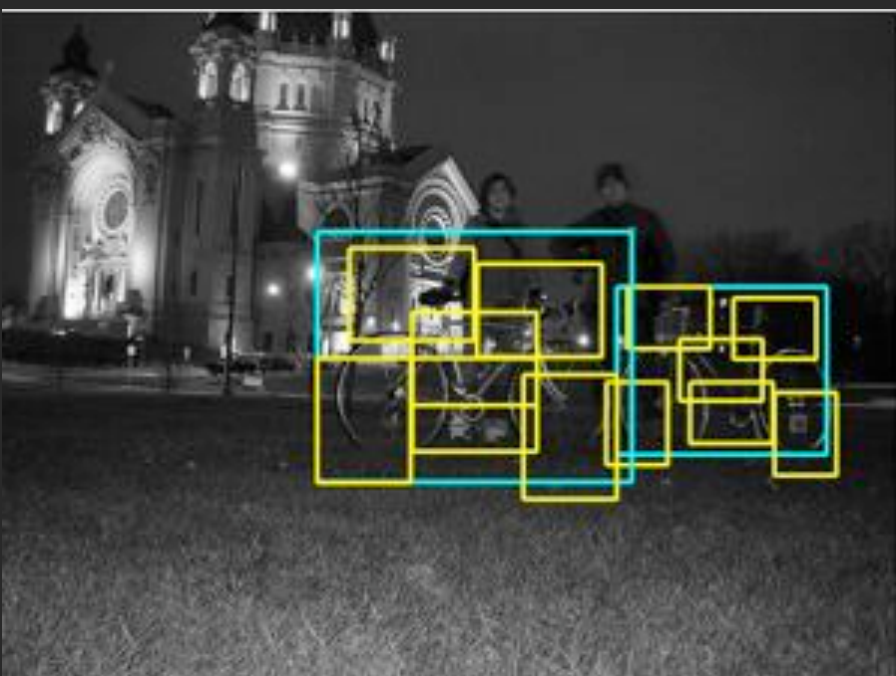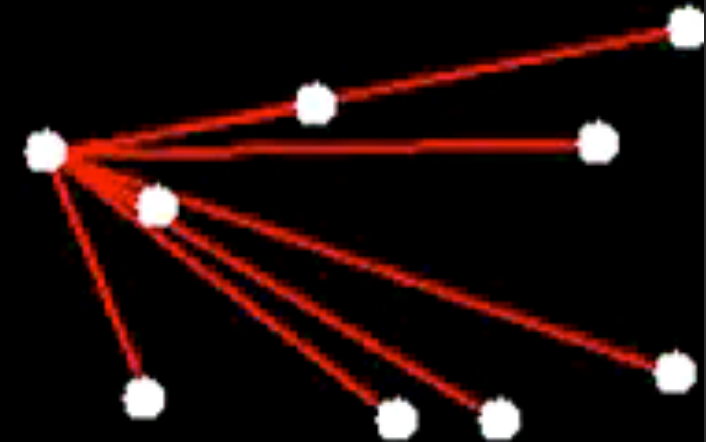# Treat ground-truth labels as partially latent



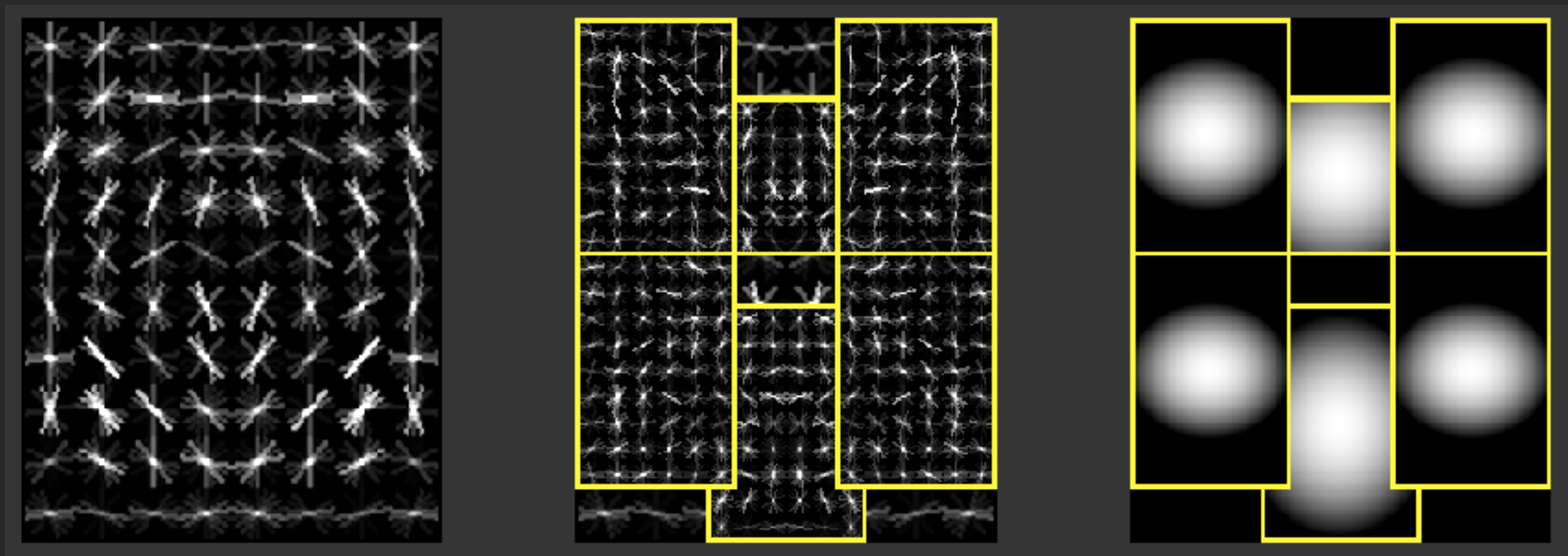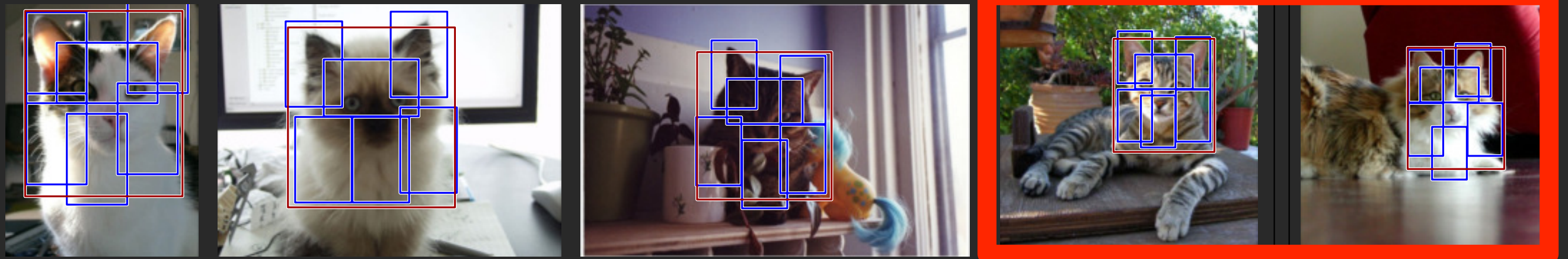Allows for "cleaning up" of noisy labels (in blue) during iterative learning

# Example models
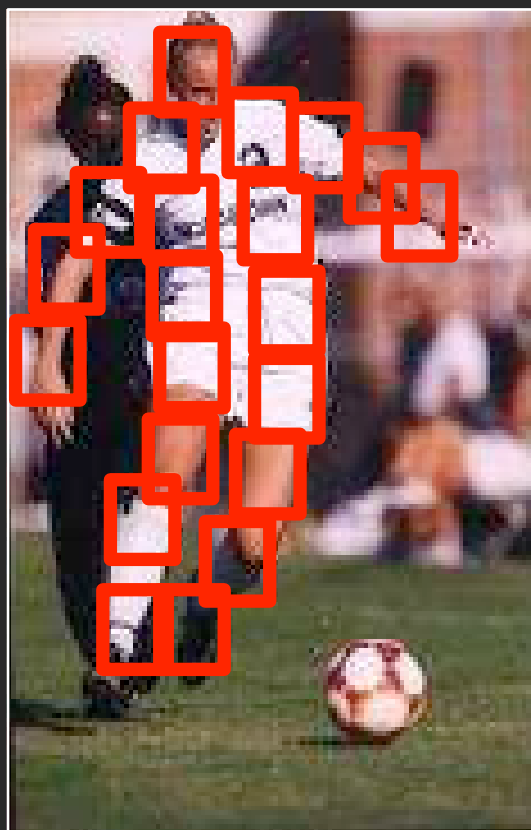
# Example models

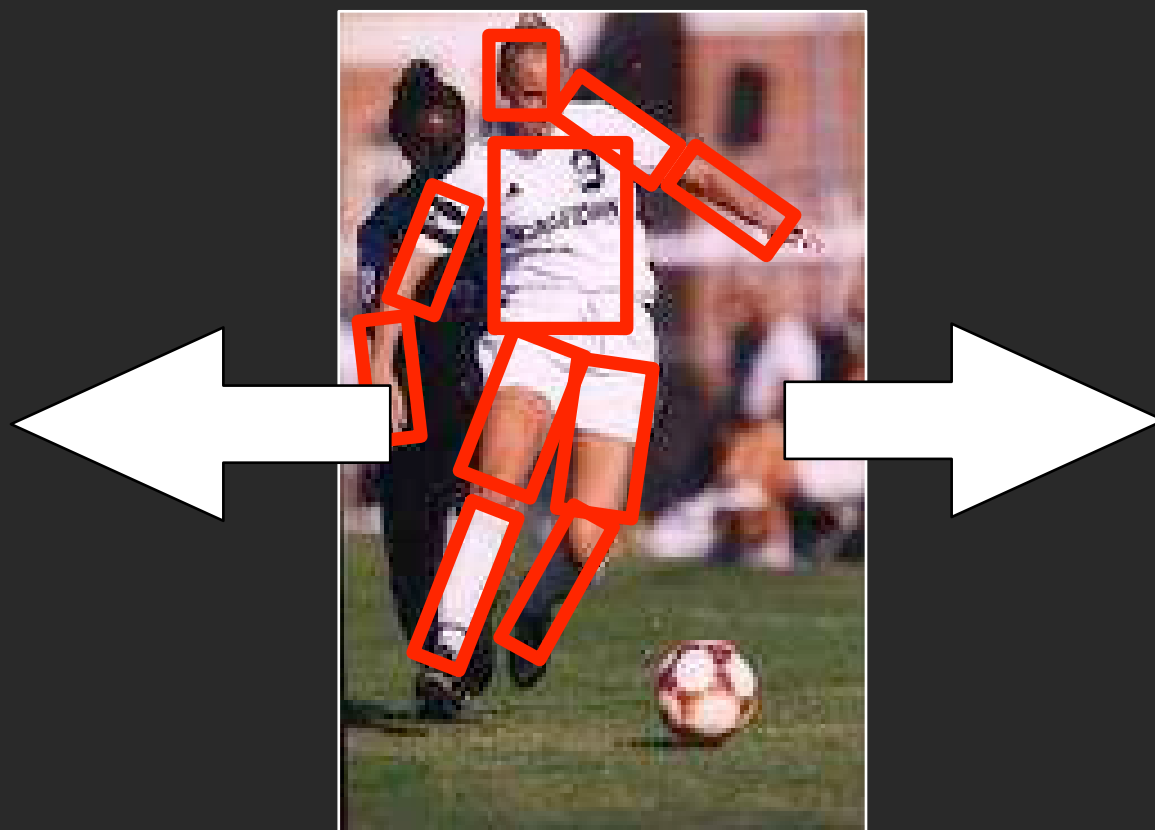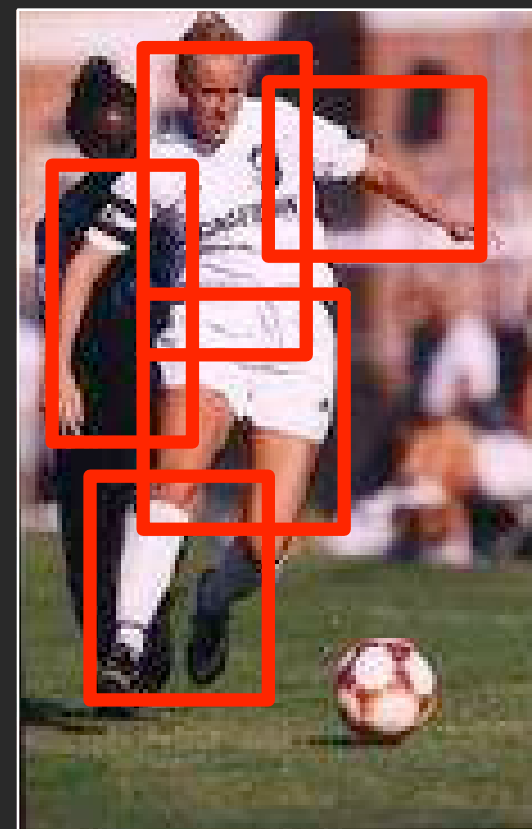False positive due to imprecise bounding box

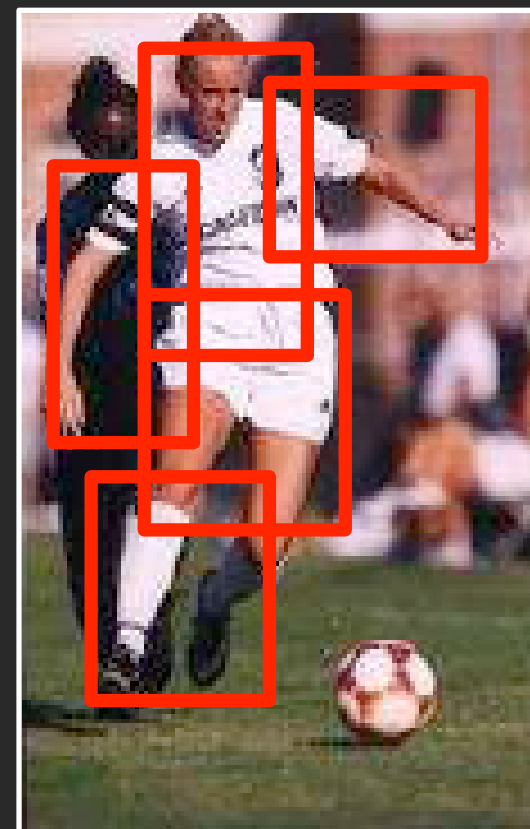# Challenge 1:

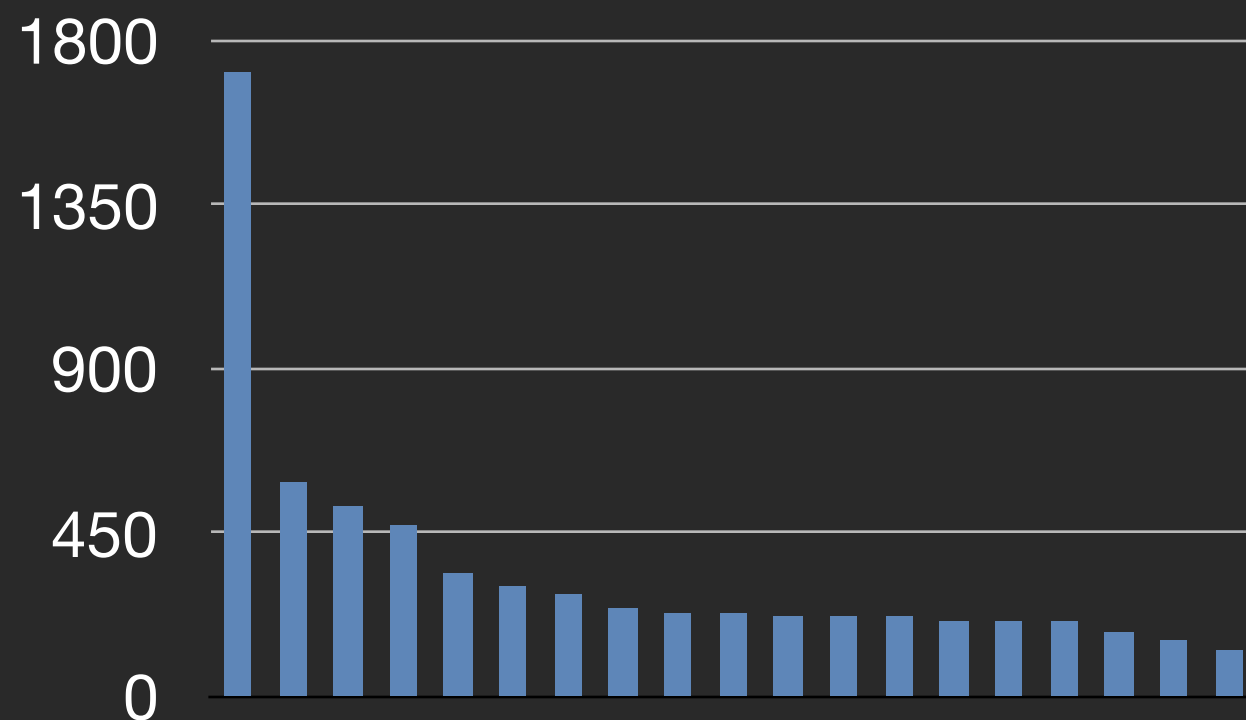## What are the right parts?
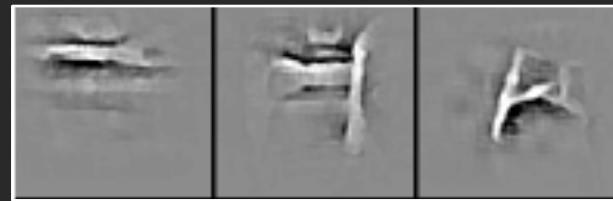
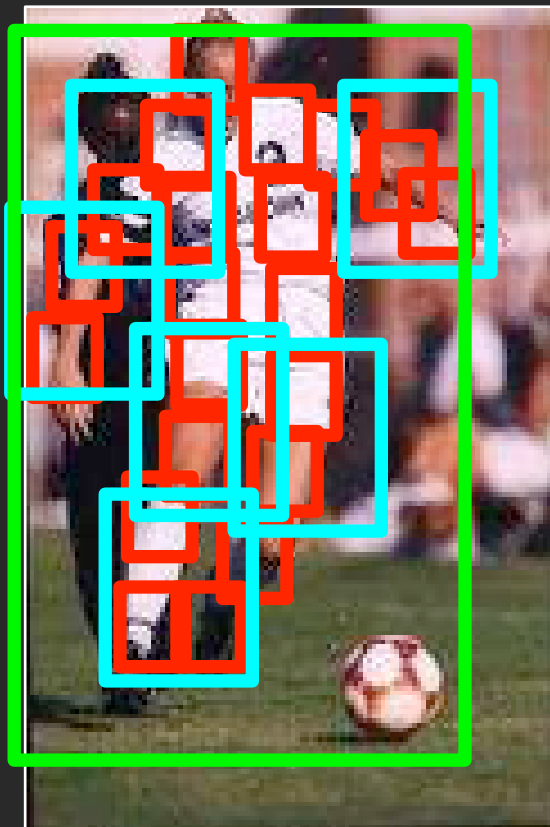

Patches

Skeleton

"Poselets"

Are "computer graphics" primitives the right choice?
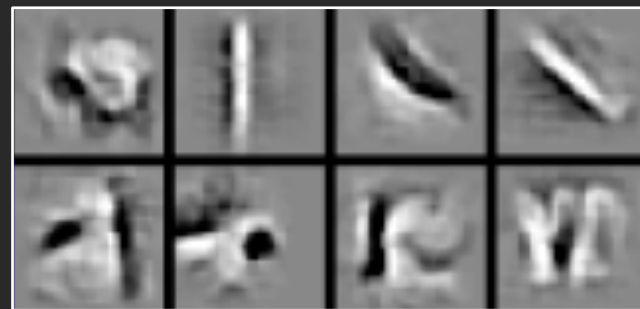
# Challenge 2:

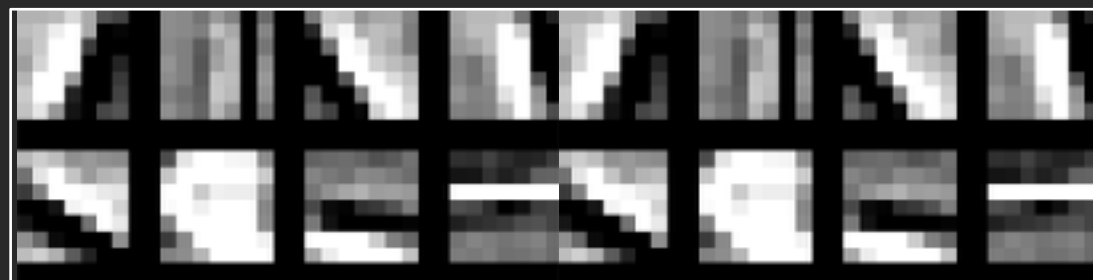How to deal with long-tail distribution of part types?

# Solution 1+2:
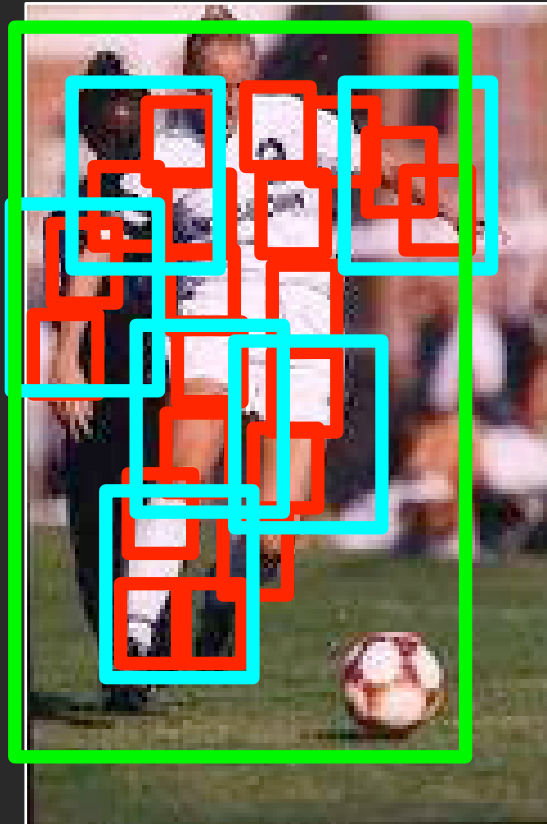# latent hierarchical (or "deep") models



objects

parts

subparts

Inference on such models requires layers of convolution and max-pooling

(we've *almost* derived a convolutional neural net)
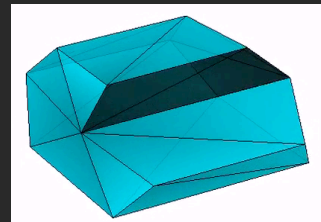
# Latent hierarchical models



$$S(x, z) = \sum_{i \in V} w_i \cdot \phi(x, z_i) + \sum_{ij \in E} w_{ij} \cdot \psi(z_i, z_j)$$
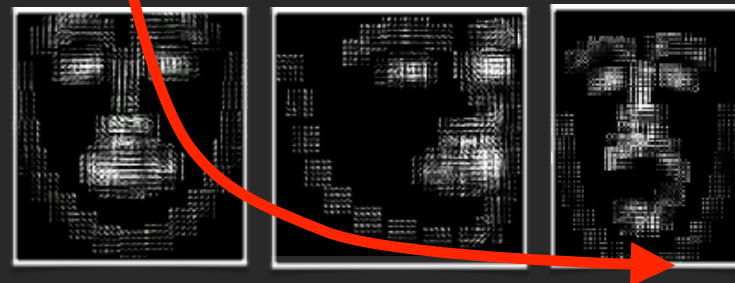
Next lecture (deep models): define $z_{in}$ to be binary variable that specify if (sub)part i is found at location n

# Parts: a look back

Recognition through reconstruction: latent-variable classification



Sharing + synthesis: zero & one-shot learning for tails



Representation learning: part discovery