

# Towards Scalable Representations of Object Categories: Learning a Hierarchy of Parts

Sanja Fidler and Ales Leonardis

*Faculty of Computer and Information  
Science  
University of Ljubljana, Slovenia*



# Overview

1. Goal
2. Motivation
3. Hierarchical object description
4. Learning part compositions
5. Results



# Goal

- Detection & Recognition of a large number of object categories

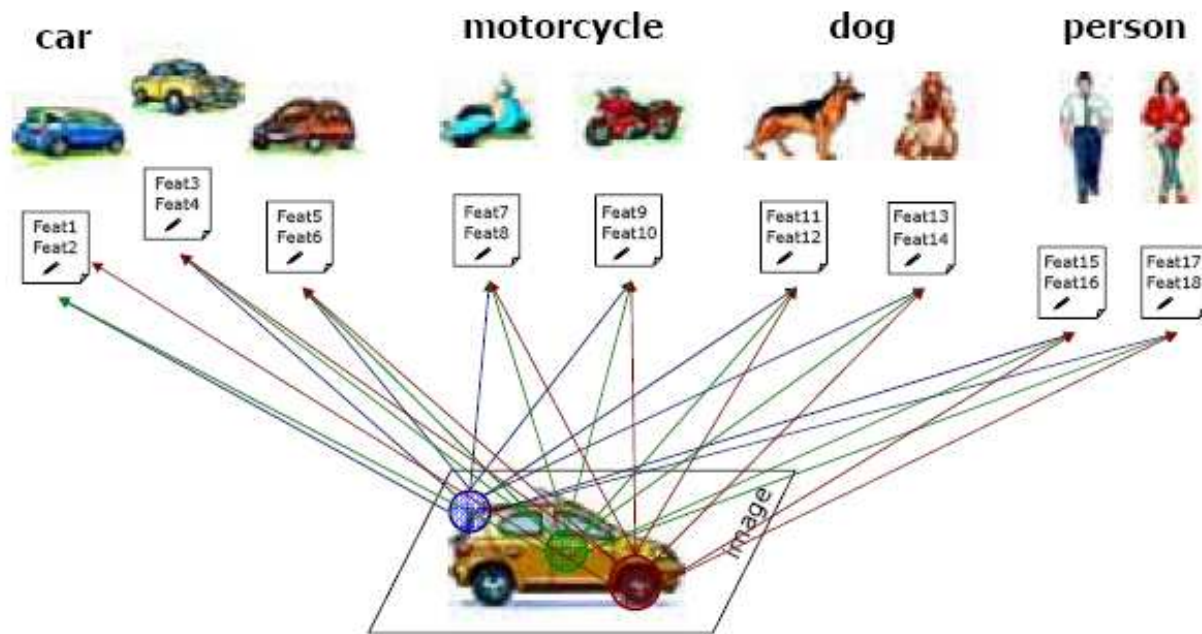



# Desired Properties

- *Computational Plausibility*: Fast indexing & matching
- *Statistics driven learning*: Unsupervised learning of object parts for compact & concise representation
- *Robust detection*: Flexible, yet accurate models
- *Fast, Incremental Learning*: Easy addition of new object categories

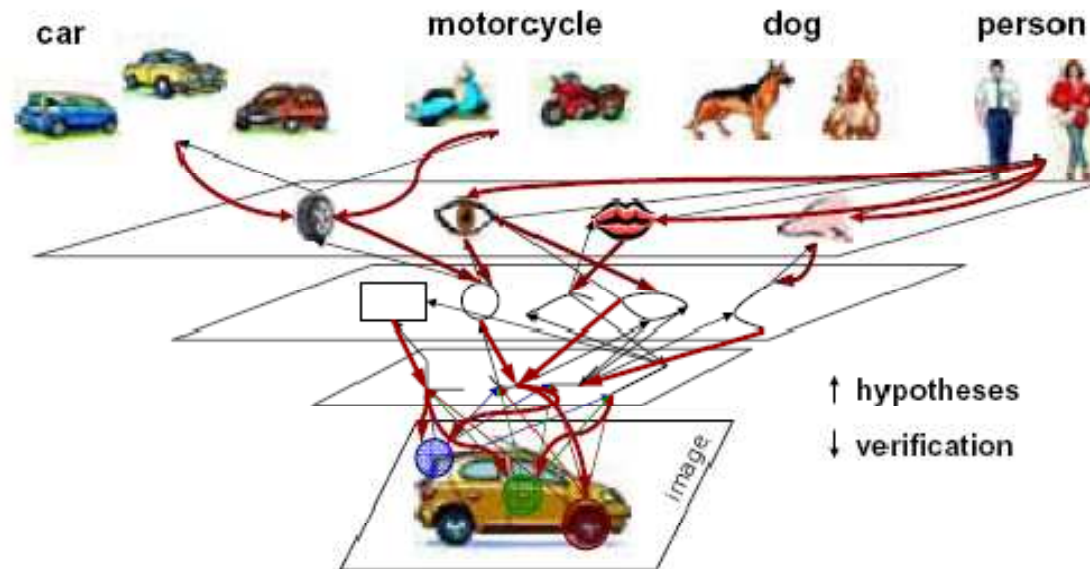
# Flat Representations

- Match each set of features to all object in the collection to find a good match
- Computationally demanding



# Hierarchical Representations

- A natural framework for indexing & matching

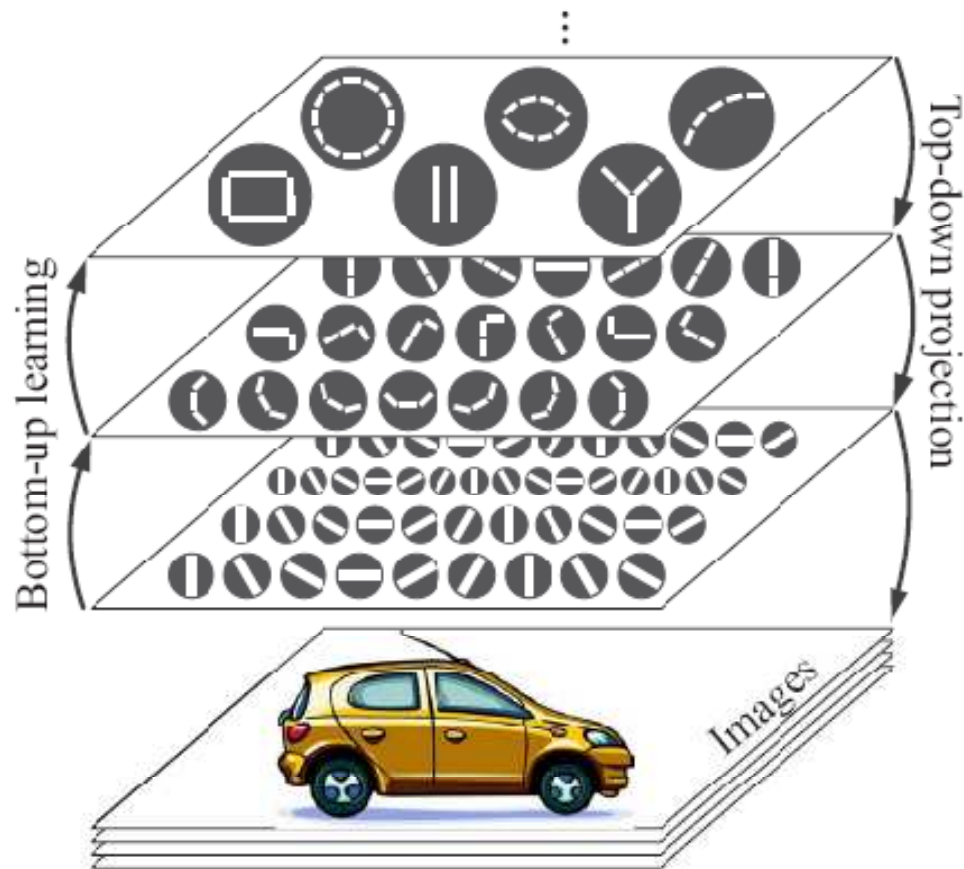




# Object Composition Hierarchy

- We wish to learn a hierarchical representation for the objects in an ***unsupervised*** manner
- Each object is made up of “parts” (*compositionability*)
- Parts appear in all levels of the hierarchy, where subsequent layers’ parts are compositions of parts from previous layers
  - All but the most basic parts are composed of parts

# Part Hierarchy Demonstrated







# Hierarchy Structure

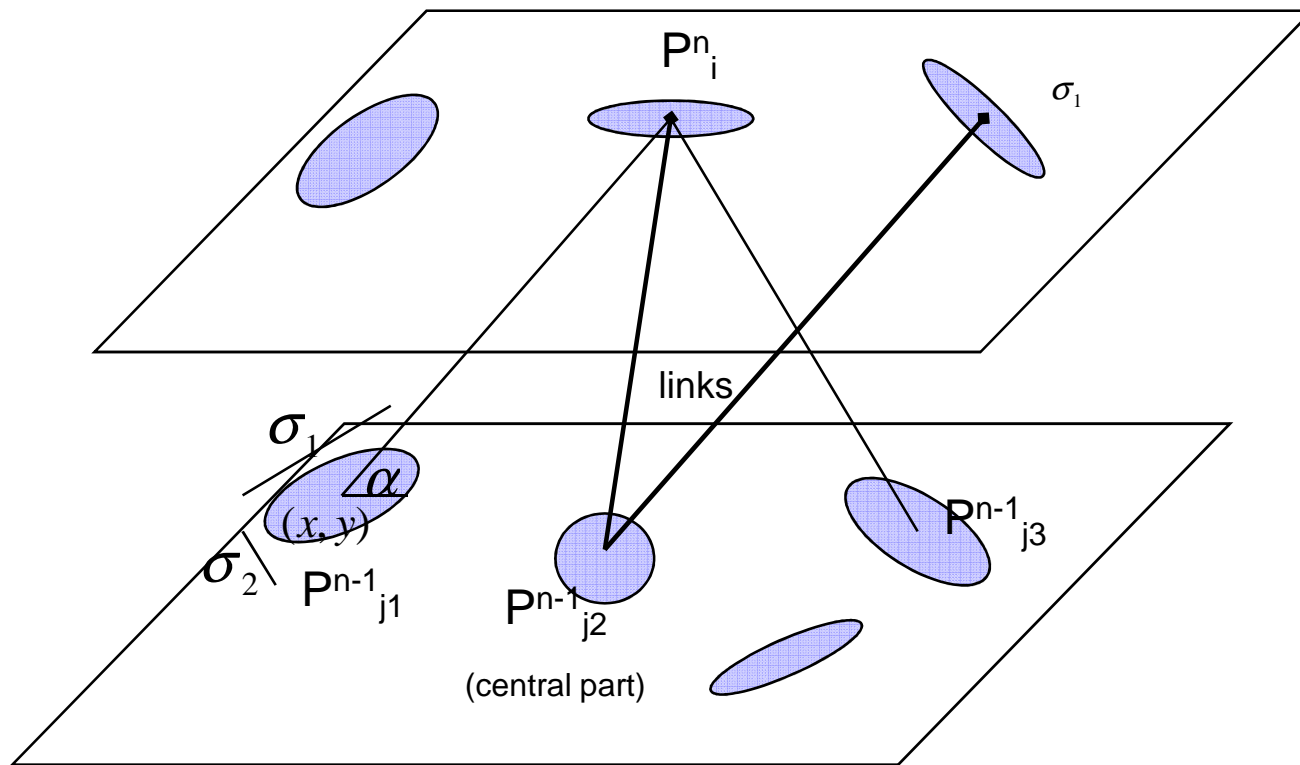
- $L_n$  – n'th Layer
- $P_i^n$  – i'th part of n'th layer, described by
  - Center of mass
  - Orientation
  - List of subparts from  $L_{n-1}$ , with position & orientation relative to  $P_i^n$ .



## Hierarchy Structure ctd.

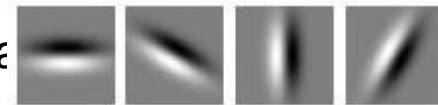
- *Central Part* – one specific subpart from  $L_{n-1}$  that indexes into  $P^n_i$ . Its location and orientation are defined as  $(0,0)$ ,  $0$  resp.
- Contains a list of
- $\{P^{n-1}_j, \alpha_j, (x_j, y_j), (\sigma_{1j}, \sigma_{2j})\}_j$ , denoting relative orientation, position, and position variance (via a gaussian) around  $(x_j, y_j)$ .
- Links – a list of all parts from  $L_n$  that this part indexes to.

# Hierarchy Structure



# Initialization

- $L_1$  is a set of local oriented filters:
  - 8 Odd **Gabor** filters, oriented at 45 degrees  $\alpha$
- At multiple scales
- $L_1$  parts are extracted from image via local maxima of filter responses (above threshold)
- Parts are denoted as  $\{\pi^1_i\}_i$
- $\pi^n_i = \{P_i, \alpha_i, x_i, y_i\}$  is a realization of part  $i$  from layer  $n$ , with the orientation & position at which it was found in the image
- $\Lambda_n(\pi^n_i)$  – List of image locations to contribute to part  $\pi^n_i$



# Indexing & Matching

---

**Algorithm 1** : Indexing and matching

---

```
1: INPUT:  $\{\{\pi_i^{n-1}\}_i, \Lambda_{n-1}\}_{scale=1}^{n_{scales}}$ 
2: for each scale do
3:    $\Pi_{scale} = \{\}$ 
4:   for each  $\pi_i^{n-1} = \{\mathcal{P}_{i_k}^{n-1}, \alpha_i, x_i, y_i\}$  do
5:     Rotate the neighborhood of  $\pi_i^{n-1}$  by angle  $-\alpha_i$ 
6:     for each part  $\mathcal{P}^n \in Links(\mathcal{P}_{i_k}^{n-1})$  do
7:       Check for subparts of  $\mathcal{P}^n$  according to their relative
       positions and spatial variance
8:       if subparts found then
9:         add  $\pi^n = \{\mathcal{P}^n, \alpha_i, x_i, y_i\}$  to  $\Pi_{scale}$ ,
         set  $\Lambda_n(\pi^n) = \bigcup \Lambda_{n-1}(\pi_j^{n-1})$ , where  $\pi_j^{n-1}$  are
         the found subparts of  $\mathcal{P}^n$ .
10:      end if
11:    end for
12:  end for
13: end for
14: Perform local inhibition over  $\{\pi_i^n\}$ 
15: return  $\{\{\pi_i^n\}_i, \Lambda_n\}_{s=1}^{n_{scales}}$ 
```

---



# Learning Part Hierarchy

- We'd like to reduce computational complexity, by:
  - Choosing parts with few occurrences (reduces the subsequent matching process)
  - Create simple models (limit overall number of parts)
  - Perform local inhibition to remove part redundancy
- Learn layers and links sequentially:
  - Perform voting for each layer
  - Choose best composition of parts for higher layer
- In addition: Choose parts to cover images well



# Incremental Learning of Layers

- $L_1$  : Oriented Gabor filters
- Subsequent layers: Learn compositions with *increasing complexity* (no. of parts), called *s-compositions*. Limit  $s$  to 4;
- An *s-composition*  $C^n_s$  is made up of  $s+1$  parts ( $s$  parts + 1 central)



# 1-compositions

- Choose a part  $P^{n-1}_i$  with low avg. image frequency ( $N_i$ ), to be the central part.
- Choose  $P^{n-1}_j$  s.t.  $N_i \leq N_j$ . From the neighboring features (neighborhood size chosen to minimize information loss)
- Perform Local inhibition to disregard parts having low novelty over central part
- $\{C^n_{s=1}\} = \{P^{n-1}_i, \{P^{n-1}_j, map_j\}\}$  is the set of possible 1-compositions.
- $map_j$  – Spatial distribution of appearance of  $P^{n-1}_j$  conditioned on  $P^{n-1}_i$  being the central part.
- $Links(P^{n-1}_i)$  – set of all compositions with  $P^{n-1}_i$  as the central part



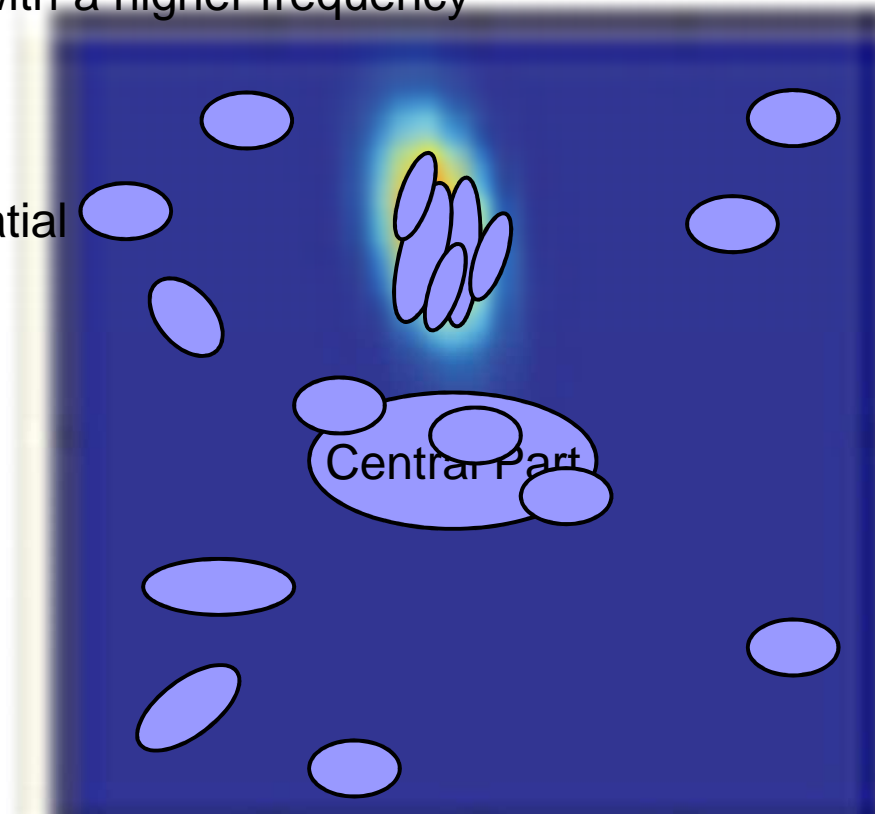
# Formation of spatial maps

Candidate Parts with a higher frequency

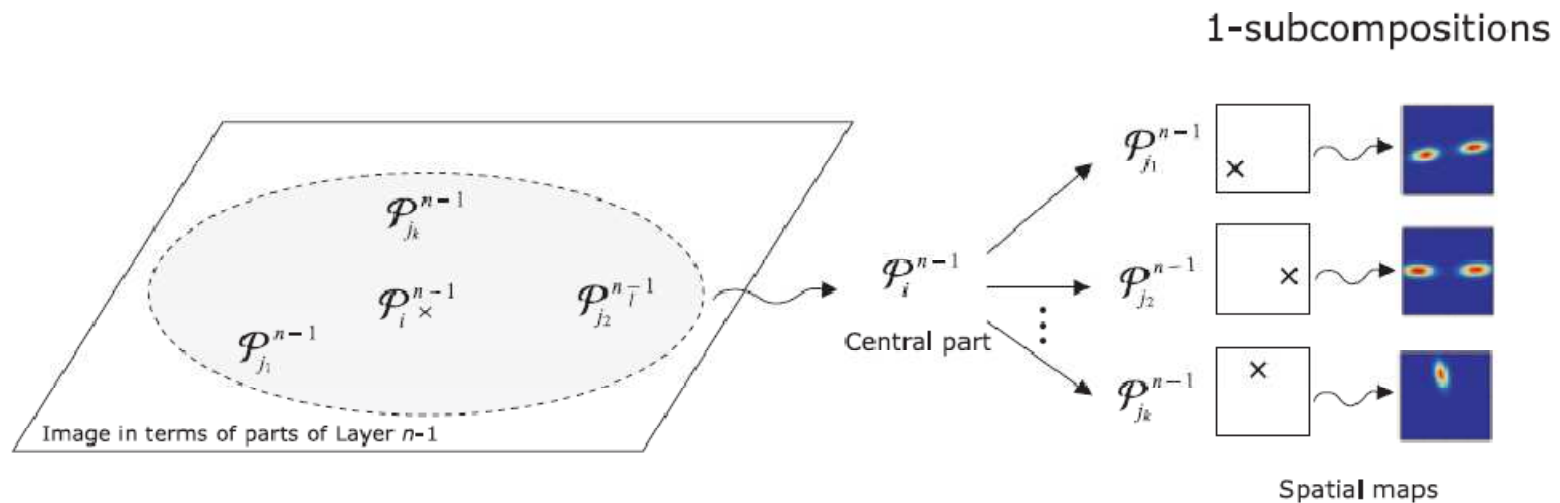
Local Inhibition

Find Peaks in Spatial  
Distribution

Model as  
Gaussian Density



# Spatial Maps



- $(\sigma_{1j}, \sigma_{2j})$  – represent the spatial variability of the distribution of  $P_j^{n-1}$  conditioned on the position of  $P_i^{n-1}$



# Spatial Maps ctd.

- probability for composition = sum of votes within area of variability / total inspected neighborhoods
- Keep only statistically significant 1-compositions:
  - $\Pr(C^n_1) \gg \Pr(P^{n-1}_i) \Pr(P^{n-1}_j)$
  - $N(C^n_1) > thresh_{n-1}$

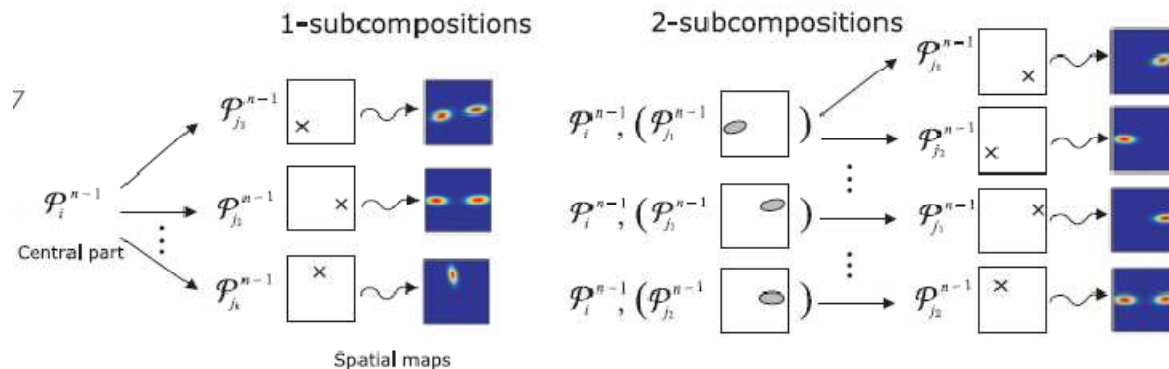
# S-Subcompositions

Central part

Additional part

$$\{C_s^n\} = \{P_i^{n-1}, \{P_{jm}^{n-1}, (x_{jm}, y_{jm}), ((\sigma_{1jm}, \sigma_{2jm}))_{m=1,2,3,4}\}, \{P_j^{n-1}, map_j\}\}$$

- i.e, build compositions using the previously learned s-1 compositions and one additional part.
- $map_j$  is updated whenever *all* parts forming a certain composition are found in the local image neighborhood.
- Prune possible combinations similarly to 1-compositions.
- *When no new decompositions pass the set statistical significance threshold, the layer learning ends.*



# Learning of S-subcompositions

---

**Algorithm 2** : Learning of  $s$ -subcompositions

---

```
1: INPUT: Collection of images
2: for each image and each scale do
3:   Preprocessing:
4:   process image with  $\mathcal{L}_1$  parts to produce  $\{\{\pi_i^1\}_i, \Lambda_1\}$ 
5:   for  $k = 2$  to  $n - 1$  do
6:      $\{\{\pi_i^k\}_i, \Lambda_k\} = \text{Algorithm 1}(\{\{\pi_i^{k-1}\}_i, \Lambda_{k-1}\})$ 
7:   end for

   Learning:
8:   for each  $\pi_i^{n-1} = \{\mathcal{P}^{n-1}, x_i, y_i\}$  do
9:     for each  $\mathcal{C}_s^n \in \text{Links}(\mathcal{P}^{n-1})$  do
10:      Find all parts  $\pi^{n-1}$  within the neighborhood
11:      Match the first  $(s - 1)$ -subparts contained within the
        subcomposition relative to the central part
12:      Perform local inhibition:  $\Lambda(\text{neigh. parts}) :=$ 
         $\Lambda(\text{neigh. parts}) \setminus \bigcup \Lambda(\text{found subparts})$ . Keep
        parts that have  $|\Lambda(\pi^{n-1})| \geq \text{thresh} \cdot |\Lambda(\pi_i^{n-1})|$ .
        We use  $\text{thresh} = 0.5$ .
13:      If all  $s - 1$  subparts are found and  $s$ -th subpart ap-
        pears anywhere in the neighborhood, update the spa-
        tial map for the  $s$ -th subpart.
14:     end for
15:   end for
16: end for
```

---



# Part Selection & grouping

- To control the complexity, compositions are removed if parts within them index too many parts in subsequent layers
- Usually 10-20 links per part
- Determined by computational resources
- Parts are deemed equal if average part overlap over set of images is large enough; this removes different yet perceptually similar parts



# Learning process

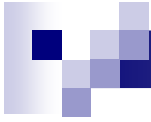
- Lower layers: Category independent, containing parts shared among many object classes  
→ Learn a union of image classes
- Higher layers: Number of part combinations increases rapidly. On the other hand, part combinations “specialize” for object categories  
→ learn for each category by itself



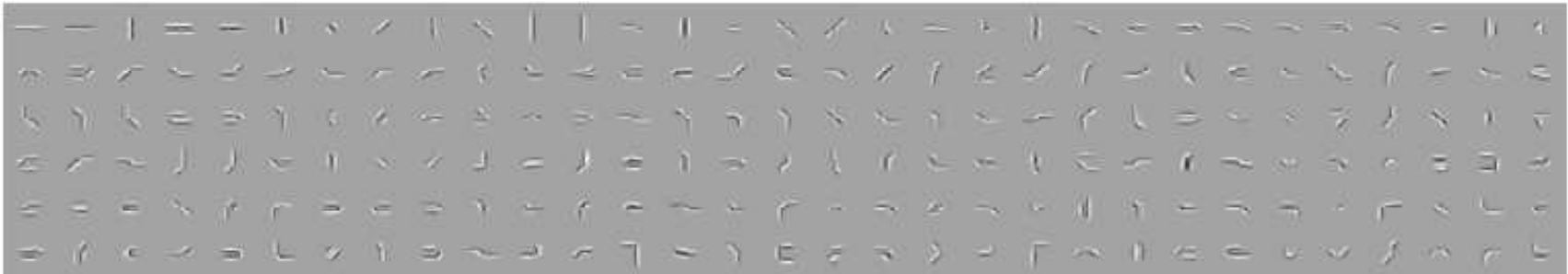
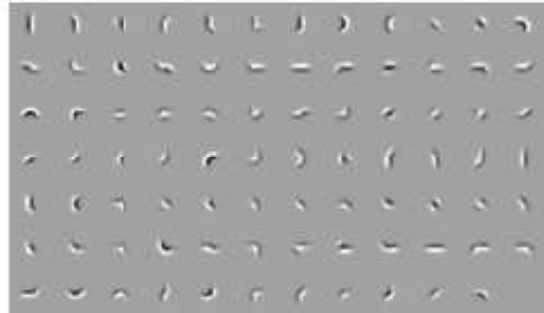
# Results

- Learned a collection of 3200 images from 15 categories (cars, faces, mugs, dogs...)
- Results are comparable with current approaches regarding object *Localization* for single-scale, and slightly better for multi-scale.





# $L_2, L_3$ (non-specific)

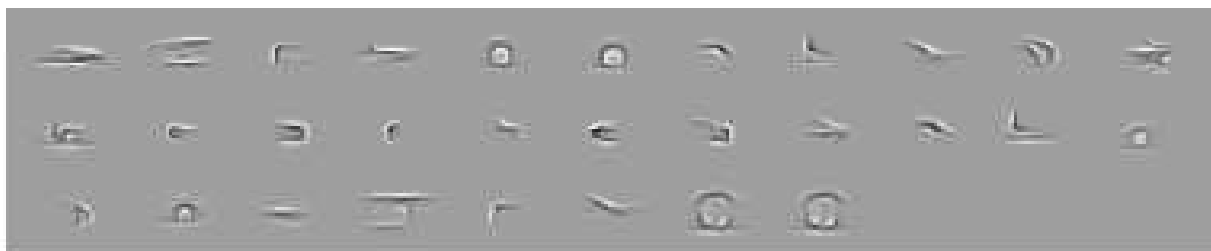


# L<sub>4</sub> (category specific)

faces



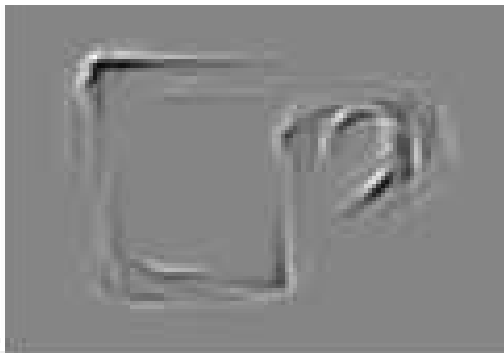
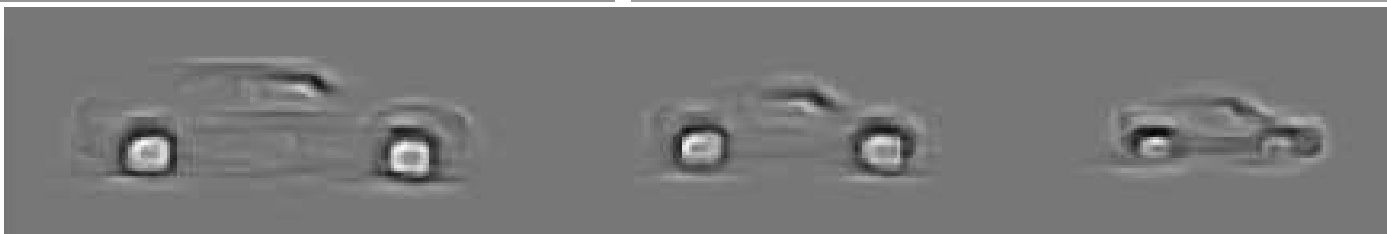
cars




mugs



L<sub>5</sub>





# Conclusions & Properties

- + Low-Level parts are mostly category independent
- + Mid-Level parts take on intuitive, familiar shapes (wheels, eyes, handles)
- + High levels still require supervision...
- + Number of indicative parts per image drops significantly for higher layers



# Summary

- A hierarchical representation for efficient indexing & matching
- High level sparseness allows for a large number of visual categories
- Adding new objects is easy since most low-level features are shared between objects



Questions?