Adding discriminative power to hierarchical compositional models for object class detection

University of Ljubljana Faculty of Computer and Information Science

Matej Kristan⁽¹⁾, Marko Boben⁽¹⁾, Domen Tabernik⁽¹⁾, Aleš Leonardis^(2,1) ⁽¹⁾CIS, University of Ljubljana, ⁽²⁾CN-CR Centre, University of Birmingham {matej.kristan}{marko.boben}{domen.tabernik}{ales.leonardis}@fri.uni-lj.si



1 Standard vs. hierarchical object detection

- Standard object category detection =
- sliding windows or oversegmentations + high-dimensional descriptors + SVM.
- Problems: slow inference time, large storage requirement, poor scaling.
- Hierarchical compositional models (e.g., [1,2,3,4]) can avoid these drawbacks.
- Hierarchical: organize the representations by level of granularity.
- Compositional: each part is a composition of parts from a lower layer.

A sliding-window inference

Hierarchical inference

Library of parts in IHop







2 Issues with generative hierarchies

- Layered hierarchy of parts (IHop) [1]
 - Fast inference, good scaling and low storage.
 - Due to significant sharing of parts acrosss categories.
- Problem: difficult to discriminate between visually-similar categories.
- Proposed solution: identify the subset of parts that differentiate pairs of

categories and combine them into **discriminative nodes**.





2 Which parts are active at detection?

• Detection of a horse (left) and of a cow (right) by IHop [1]:

• The frequency/strength of the activated parts can be sumarized by a **cumulative histogram** of part responses:

4. Training / detection

Training stage:

Detection stage:

5. Experimental results (cont.)

Experiment 2:

• Goal: Demonstrate the improvements in hypotheses rescoring. • ETHZ shape dataset (Ferrari et al., 2010):

• Half images of category for training and the other half + all images of remaining categories for testing, (5 random splits). • Evaluation: Detection rate at 1FP per image at 50% overlap • Hop trained with 7 layers (~525 parts per experiment).

3 Adding discriminative power

• We cast **identification of discriminative parts** and formation of discriminative nodes as a problem of finding a **sparse linear separation** between the cumulative histograms of parts responses.

5. Experimental results

Two experiments designed to evaluate improvement of the proposed discriminative IHop (dlHop) over the baseline IHop [1].

Experiment 1:

Goal: Analyze discrimination between two visually similar categories and the background.

Dataset composed of Weizman Horses (Borenstein&Ullman, 2008) and Leeds Cows (Leibe et al., 2008):

Results of hypothesis rescoring. N _{disc} denotes the number of selected library part.											
lHoP	dlHoP $[N_{\rm disc}]$	PSM	Hough	w_{ac}	M^2HT	PMK	PMK				
[1]	our work	[7]	[8]	[6]	[5]	[6]	[7]				
92.5	$92.5\ [5.2\ (1.3)]$	90.4	43.0	80.0	85.0	80.0	90.4				
79.6	$85.4 \ [7.4 \ (1.7)]$	84.4	64.4	92.4	67.0	89.3	96.4				
75.1	82.3 [13 (4.6)]	50.0	52.2	36.2	55.0	80.9	78.8				
85.9	86.5[13.2(6.9)]	32.3	45.1	47.5	55.0	74.2	61.4				
58.6	70.5~[6~(2.6)]	90.1	62.0	58.8	42.5	68.6	88.6				
78.3	$83.4 \ [9.0 \ (5.1)]$	69.4	53.3	63.0	60.9	78.6	83.2				
	lHoP [1] 92.5 79.6 75.1 85.9 58.6 78.3	thypothesis rescoring. N_a lHoPdlHoP $[N_{disc}]$ [1]our work92.592.5[5.2 (1.3)]79.685.4[7.4 (1.7)]75.182.3[13 (4.6)]85.986.5[13.2 (6.9)]58.670.5[6 (2.6)]78.383.4[9.0 (5.1)]	hypothesis rescoring. N_{disc} deno1HoPdlHoP $[N_{disc}]$ PSM[1]our work[7]92.592.5[5.2 (1.3)]90.479.685.4[7.4 (1.7)]84.475.182.3[13 (4.6)]50.085.986.5[13.2 (6.9)]32.358.670.5[6 (2.6)]90.178.383.4[9.0 (5.1)]69.4	IhopdlHoP $[N_{disc}]$ PSMHough[1]our work[7][8]92.592.5[5.2(1.3)]90.443.079.685.4[7.4(1.7)]84.464.475.182.3[13(4.6)]50.052.285.986.5[13.2(6.9)]32.345.158.670.5[6(2.6)]90.162.078.383.4[9.0(5.1)]69.453.3	IHoPdlHoP $[N_{disc}]$ PSMHough w_{ac} [1]our work[7][8][6]92.592.5[5.2(1.3)]90.443.080.079.685.4[7.4(1.7)]84.464.492.475.182.3[13(4.6)]50.052.236.285.986.5[13.2(6.9)]32.345.147.558.670.5[6(2.6)]90.162.058.878.383.4[9.0(5.1)]69.453.363.0	Thypothesis rescoring. N_{disc} denotes the number of selectIHoPdlHoP $[N_{disc}]$ PSMHough w_{ac} M^2HT [1]our work[7][8][6][5]92.592.5[5.2 (1.3)]90.443.080.085.079.685.4[7.4 (1.7)]84.464.492.467.075.182.3[13 (4.6)]50.052.236.255.085.986.5[13.2 (6.9)]32.345.147.555.058.670.5[6 (2.6)]90.162.058.842.578.383.4[9.0 (5.1)]69.453.363.060.9	Thypothesis rescoring. N_{disc} denotes the number of selected libra $ \text{HoP} $ $d \text{HoP}[N_{disc}]$ PSM Hough w_{ac} M^2HT PMK $[1]$ our work $[7]$ $[8]$ $[6]$ $[5]$ $[6]$ 92.5 92.5 $[5.2 (1.3)]$ 90.4 43.0 80.0 85.0 80.0 79.6 85.4 $[7.4 (1.7)]$ 84.4 64.4 92.4 67.0 89.3 75.1 82.3 $[13 (4.6)]$ 50.0 52.2 36.2 55.0 80.9 85.9 86.5 $[13.2 (6.9)]$ 32.3 45.1 47.5 55.0 74.2 58.6 70.5 $[6 (2.6)]$ 90.1 62.0 58.8 42.5 68.6 78.3 83.4 $[9.0 (5.1)]$ 69.4 53.3 63.0 60.9 78.6				

• dlHop consistently outperforms lHop.

• On average, dlHop outperforms, or exhibits performance comparable to, the competing methods.

Example of selected discriminative parts along with the distribution of frequently chosen parts over several repetitions of the experiment.

- Parts were selected from different layers for each category. • Most parts selected from layers 3-5.
- Emphasizing more **global distinctive features** for categorization.

6. Conclusion and future work

• Sparsity of vector Θ results in selection of discriminative parts:

- Learn the posterior over scores $f(\mathbf{h}; \Theta)$ by logistic regression: $p(\text{horse}|\mathbf{h}, \Theta^{(\text{cow,horse})}) = \frac{1}{1 + \exp(-f(\mathbf{h}; \Theta^{(\text{cow,horse})}))}$
- Sparse Θ is computed via a variational approach from [9].

• Split: 75 images for training, 263 for testing. • IHop with 7 layers with [6, 33, 161, 180, 93, 104, 2] vocabulary parts. • PASCAL criterion at overlap >0.3 with ground truth (GT): False positives per experiment and average precision (AP)

• Confusion matrix: Classify by the strongest detection (after nonmax suppression) that overlaps >0.3 by GT.

Confusion matrix

	lHoP (199.8 fp)			dlHoP (166.4 fp)			
	COW	horse	back.	COW	horse	back.	
COW	32.44	40.36	27.19	62.57	15.0	22.43	
horse	12.25	68.63	19.11	6.50	78.45	15.05	
AP	0.31			0.51			

• dlHop significantly reduces confusion of cows and horses. • On average, 8% of all parts were selected for discrimination.

- Proposed a simple way to identify **discriminative parts** in a generative hierarchy.
- **Detection improved** by discriminative rescoring.
- The approach **does not hamper** the IHop's **advantages** (sharing, storage, scaling).
- Current approach takes into account **only frequency** of parts.
- Future work will explicitly take into account the positions as well.

7. References

[1] Fidler, S., Boben, M., Leonardis, A.: Evaluating multi-class learning strategies in a hierarchical framework for object detection. NIPS 2009 [2] Zhu, L.L., Chen, Y., Torralba, A., Freeman, W., Yuille, A.: Part and appearance sharing: Recursive compositional models for multi-view multi-object detection. CVPR 2010 [3] Poon, H., Domingos, P.: Sum-product networks: A new deep architecture. UAI 2011 [4] Si, Z., Zhu, S.: Unsupervised learning of stochastic and-or templates. WSIG 2011 [5] Maji, S., Malik, J.: Object detection using a max-margin hough transform. CVPR 2009 [6] Ommer, B., Malik, J.: Multi-scale object detection by clustering lines. CVPR 2009 [7] Riemenschneider, H., Donoser, M., Bischof, H.: Using partial edge contour matches for efficient object category localization. ECCV 2010 [8] Ferrari, V., Jurie, F., Schmid, C.: From images to shape models for object detection. IJCV 2010 [9] Yamashita, O., Sato, M., Yoshioka, T., Tong, F., Kamitani, Y.: Sparse estimation automatically

selects voxels relevant for the decoding of fMRI activity patterns. NeuroImage 2008

This work was supported in part by ARRS Research Programs P2-0214, P2-0095 and ARRS research projects J2-4284, J2-360, J2-2221 and FP7-ICT PACMAN