



# Motivation

Our goal: **distributed coding** of local features in a **hierarchical model** that would allow full inference.

## Distributed Coding

Considering local patch gradient energy histograms by orientation and grid location.

- Traditionally, vector quantized as *visual words*.
- Coding as a mixture of components has been shown to be empirically better.
- May represent additive image formation.



+ Decent combinatoria capacity (~N<sup>K</sup>) picture from Bruno Olshausen



### Low combinatorial capacity (N)

# Hierarchical Models

- Biological evidence for increasing spatial support and complexity of visual pathway.
- Higher layers can help resolve local ambiguities.
- Possibly fewer parameters due to sharing of lower-layer components.
- Lots of past work in vision: HMAX, Convolutional Deep Belief Nets, Hyperfeatures, Hierarchies of Parts.

## Bayesian Inference

- Visual systems deal with inherently ambiguous data, which Bayesian inference helps disambiguate.
- Hierarchical priors allow unsupervised learning of complex visual structures.



# **A Probabilistic Model for Recursive Factorized Image Features**

Sergey Karayev<sup>1</sup>, Mario Fritz<sup>2</sup>, Sanja Fidler<sup>3</sup>, Trevor Darrell<sup>1</sup> <sup>1</sup>UC Berkeley, <sup>2</sup>MPI Informatics, <sup>3</sup>University of Toronto



Setting in line with previous experiments on image coding:

- Evaluation on Caltech101
- Spatial Max Pooling
- Linear SVM classification
- Three conditions:
- Feed-forward two-layer LDA (FLDA)
- Recursive two-layer LDA (RLDA)

	A	Caltech-101			
	Model	Basis size	Layer(s) used	15	30
128-dim models	LDA	128	"bottom"	$52.3\pm0.5\%$	$58.7 \pm 1.1\%$
	RLDA	128t/128b	bottom	$55.2\pm0.3\%$	$62.6\pm0.9\%$
	[ LDĀ -	$\overline{128}$		$5\bar{3}.\bar{7}\pm\bar{0}.\bar{4}\%$	$\left  \begin{array}{c} \bar{6}0.5 \pm 1.0 \% \end{array} \right $
	FLDA	128t/128b	top	$55.4\pm0.5\%$	$61.3\pm1.3\%$
	RLDA	128t/128b	top	$59.3\pm0.3\%$	$66.0\pm1.2\%$
	FLDĀ	$\overline{128t/128b}$	both	$5\bar{7}.\bar{8}\pm\bar{0}.\bar{8}\%$	$64.2 \pm 1.0\%$
	RLDA	128t/128b	both	$61.9\pm0.3\%$	$68.3 \pm 0.7\%$

## **Results**:

- Full inference increases performance (RLDA > FLDA).
- Using both layers increases performance.

	Approach	Caltech-101		
	Model	Layer(s) used	15	30
	RLDA (1024t/128b)	bottom	$56.6\pm0.8\%$	$62.7\pm0.5\%$
Our Model	RLDA (1024t/128b)	top	$66.7\pm0.9\%$	$72.6 \pm 1.2\%$
	RLDA (1024t/128b)	both	$67.4 \pm 0.5$	$73.7 \pm \mathbf{0.8\%}$
	Sparse-HMAX [21]	top	51.0%	56.0%
	CNN [15]	bottom	_	$57.6 \pm 0.4\%$
	CNN [15]	top	_	$66.3 \pm 1.5\%$
Hierarchical	CNN + Transfer [2]	top	58.1%	67.2%
Models	CDBN [17]	bottom	$53.2\pm1.2\%$	$60.5 \pm 1.1\%$
	CDBN [17]	both	$57.7 \pm 1.5\%$	$65.4\pm0.4\%$
	Hierarchy-of-parts [8]	both	60.5%	66.5%
	Ommer and Buhmann [22]	top		$61.3 \pm 0.9\%$

Various co-authors of this work were supported in part by the Air Force Office of Scientific Research, National Defense Science and Engineering Graduate (NDSEG) Fellowship, 32 CFR 168a; a Feodor Lynen Fellowship granted by the Alexander von Humboldt Foundation; by awards from the US DOD and DARPA, including contract W911NF-10-2-0059; by NSF awards IIS-0905647 and IIS-0819984; by EU FP7-215843 project POETICON; and by Toyota and Google.

Results

# • Single-layer LDA, on small and large patches ("bottom" and "top")

## • Additional layer increases performance (FLDA > LDA).

### • Using only RLDA's lower layer is still better than LDA.

## Acknowledgements