



ISR INSTITUTO DE SISTEMAS E ROBÓTICA
UNIVERSIDADE DE COIMBRA


universität**bonn**
Rheinische
Friedrich-Wilhelms-
Universität Bonn

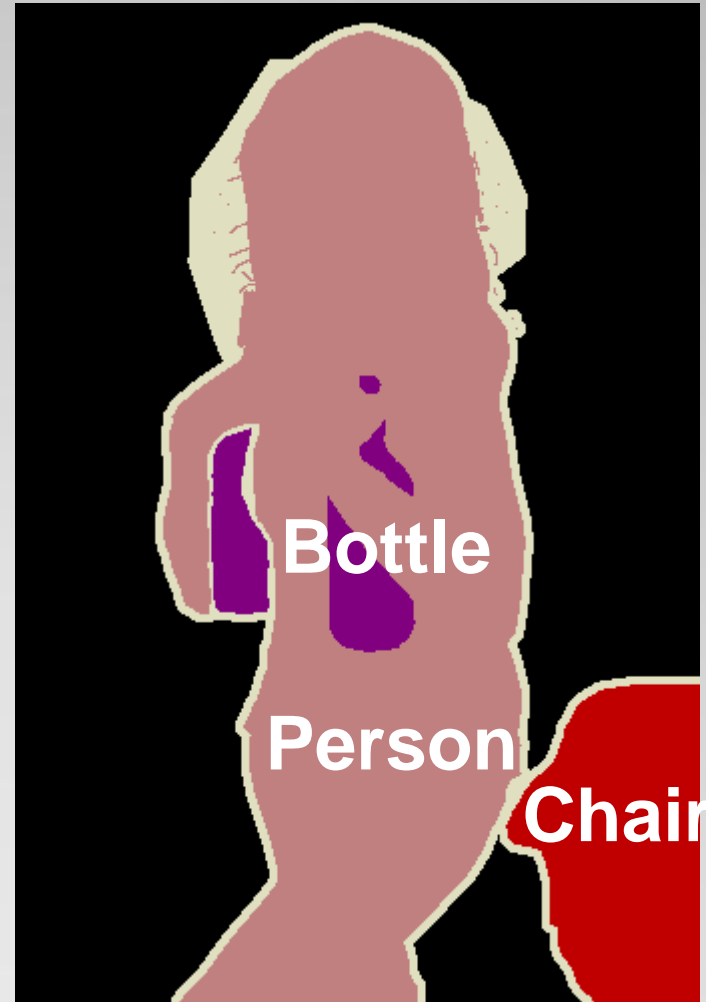
Semantic Segmentation with Second-Order Pooling

João Carreira^{1,2}, Rui Caseiro¹, Jorge Batista¹, Cristian Sminchisescu²

¹ *Institute of Systems and Robotics,*
University of Coimbra

² *Faculty of Mathematics and Natural Science,*
University of Bonn

Semantic Segmentation



Example from Pascal VOC segmentation dataset

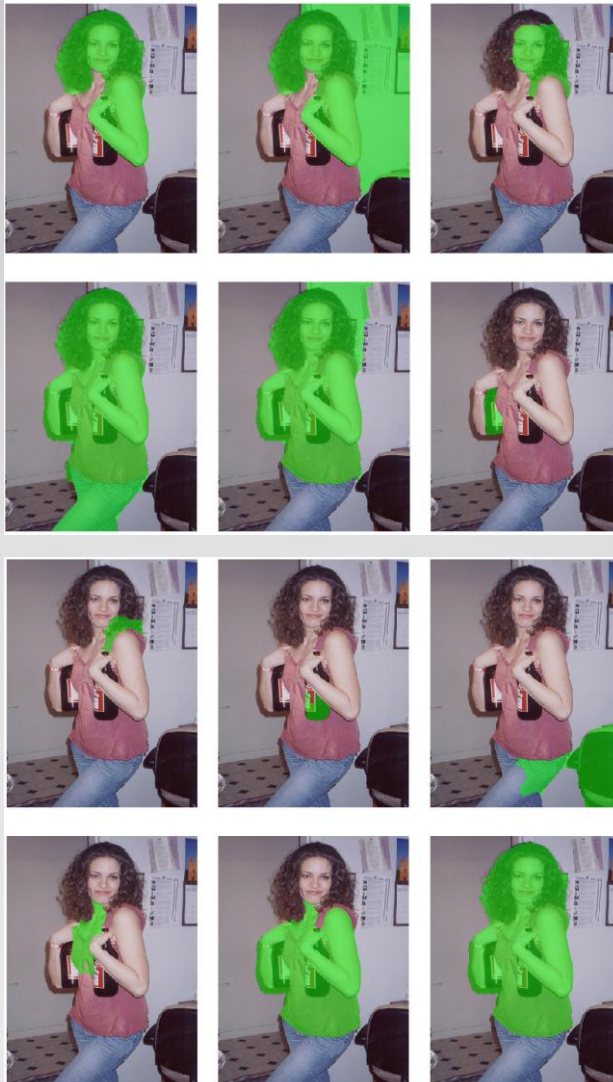
Semantic Segmentation

Our bottom-up pipeline:

Li, Carreira, Sminchisescu, CVPR 2010, IJCV2011

1. Sample candidate object regions (figure-ground)
2. Region description and classification
3. Construct full image labeling from regions

Semantic Segmentation

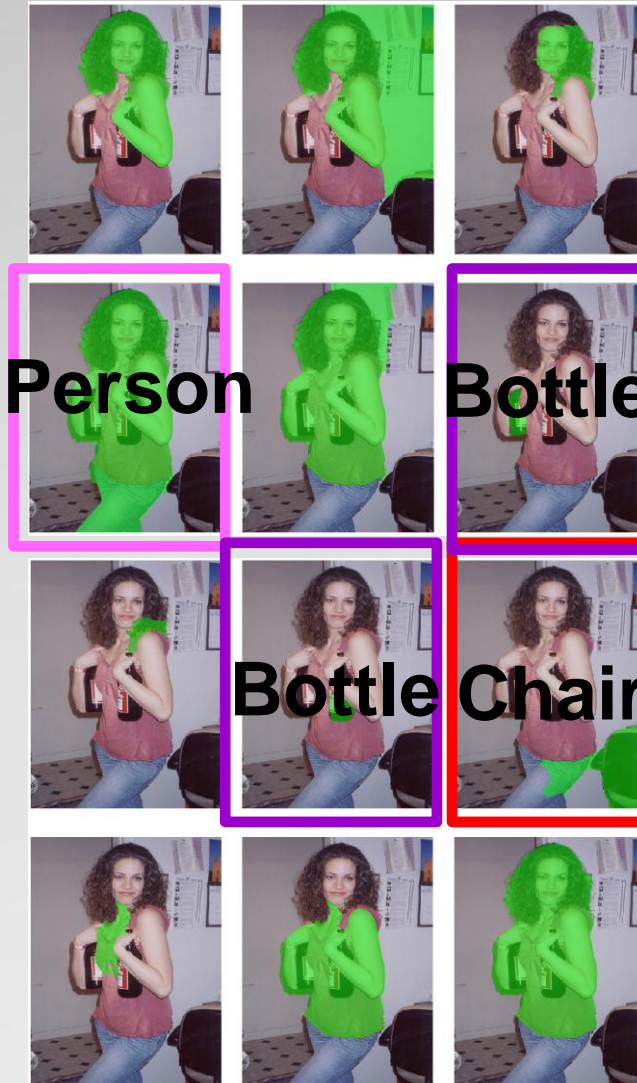


Key: generate good object candidates, not superpixels

CPMC: Constrained Parametric Min-Cuts for Automatic Object Segmentation, Carreira and Sminchisescu, CVPR 2010, PAMI 2012

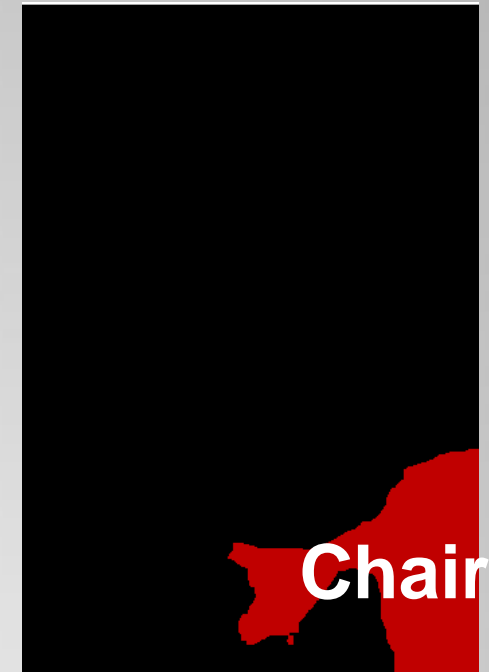
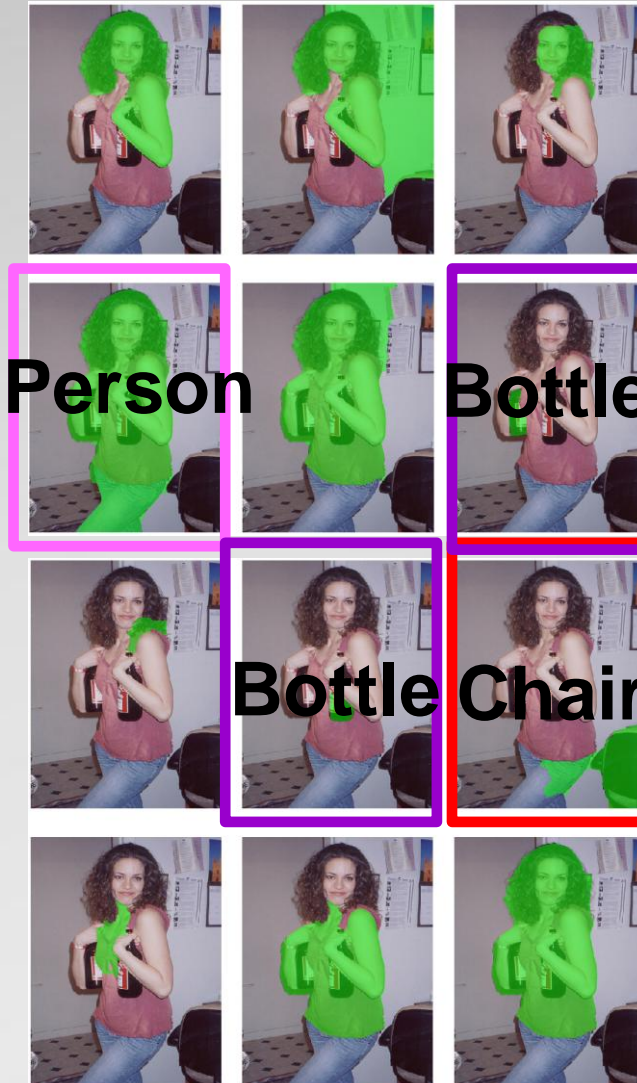
1. Sample candidate object regions (figure-ground)
2. Region description and classification
3. Construct full image labeling from regions

Semantic Segmentation



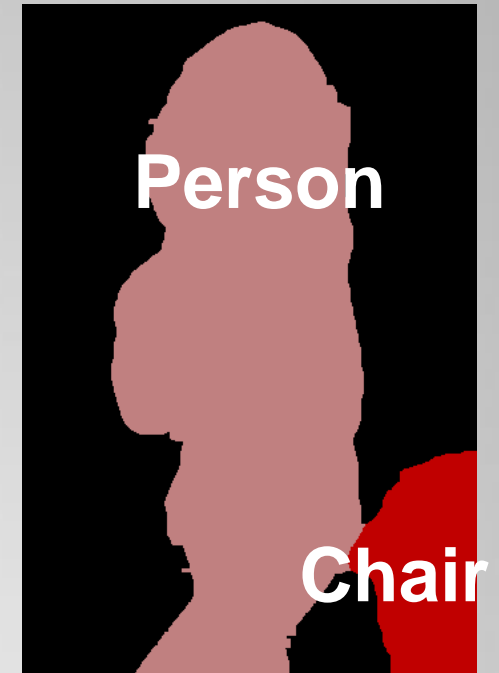
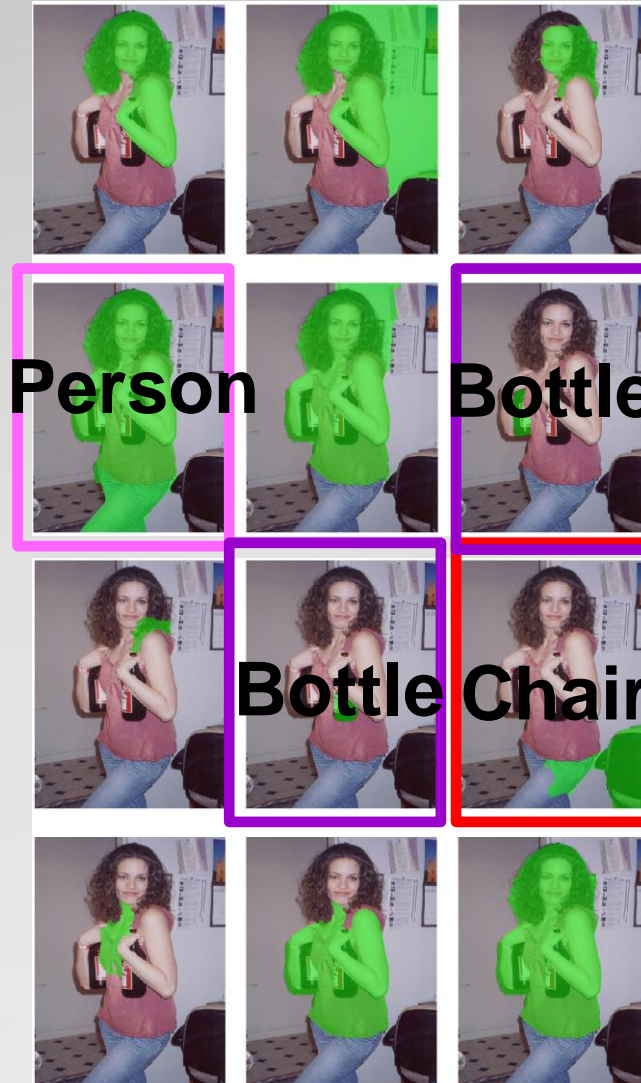
1. Sample candidate object regions (figure-ground)
2. Region description and classification
3. Construct full image labeling from regions

Semantic Segmentation



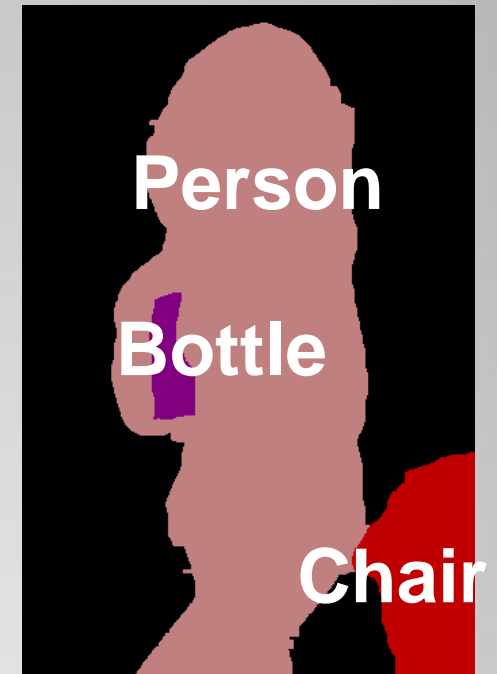
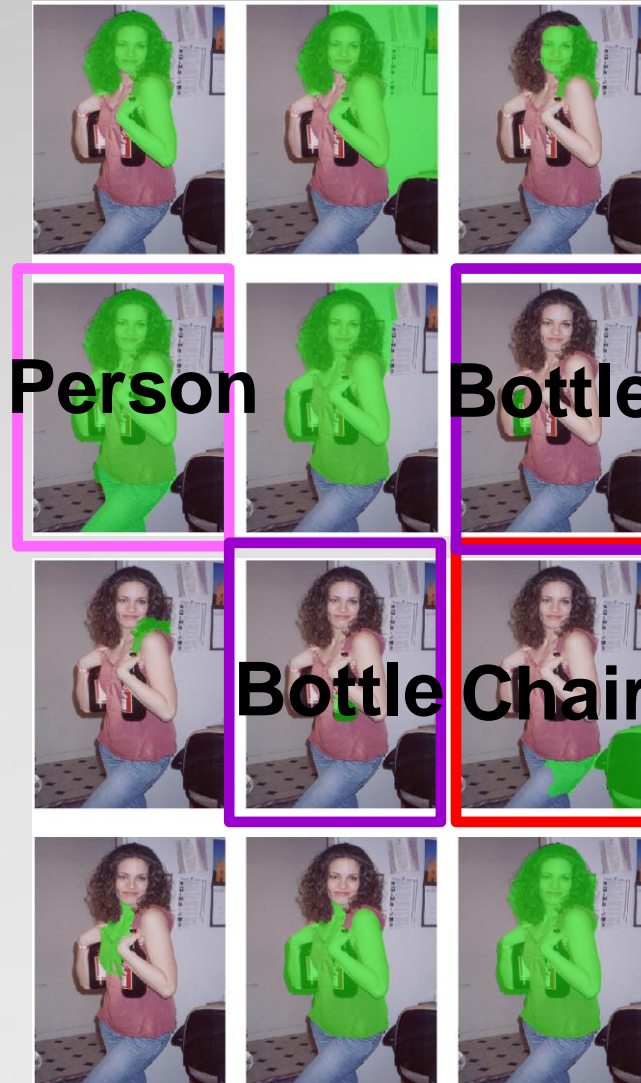
1. Sample candidate object regions (figure-ground)
2. Region description and classification
3. Construct full image labeling from regions

Semantic Segmentation



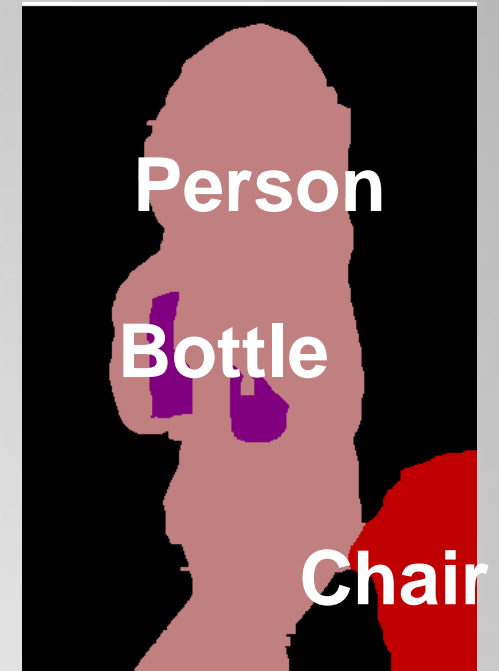
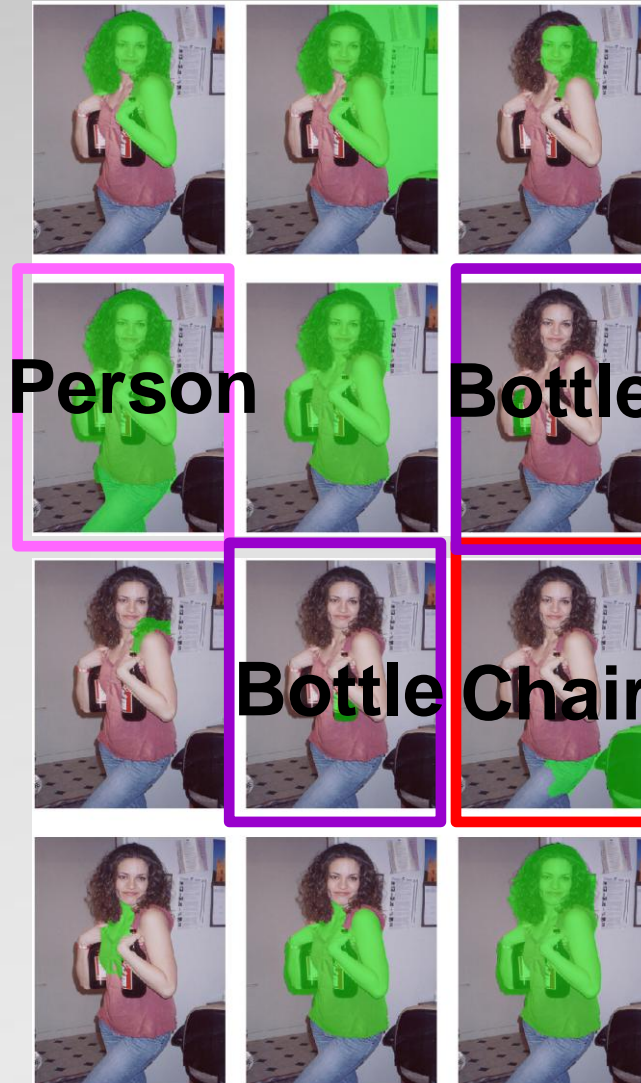
1. Sample candidate object regions (figure-ground)
2. Region description and classification
3. Construct full image labeling from regions

Semantic Segmentation



1. Sample candidate object regions (figure-ground)
2. Region description and classification
3. Construct full image labeling from regions

Semantic Segmentation



1. Sample candidate object regions
2. Region description and classification
3. Construct full image labeling from regions

Semantic Segmentation

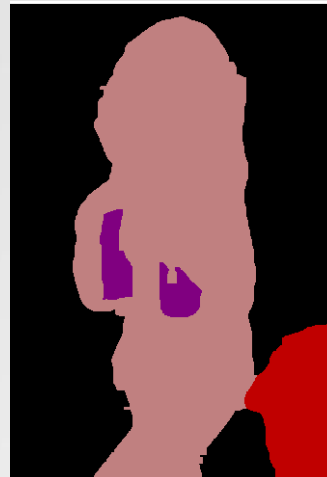
Bottom-up formulation:

1. Sample candidate object regions (figure-ground)
2. Region description and classification
3. Construct full image labeling from regions

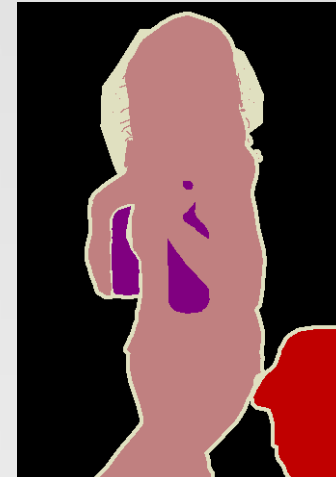
Image



Predicted



Ground Truth



Semantic Segmentation

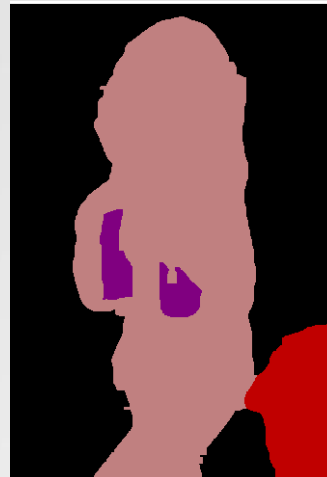
Bottom-up formulation:

1. Sample candidate object regions (figure-ground)
2. Region description and classification **This work!**
3. Construct full image labeling from regions

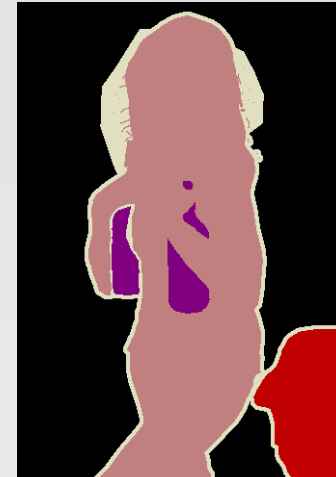
Image



Predicted



Ground Truth



Describing Free-form Regions

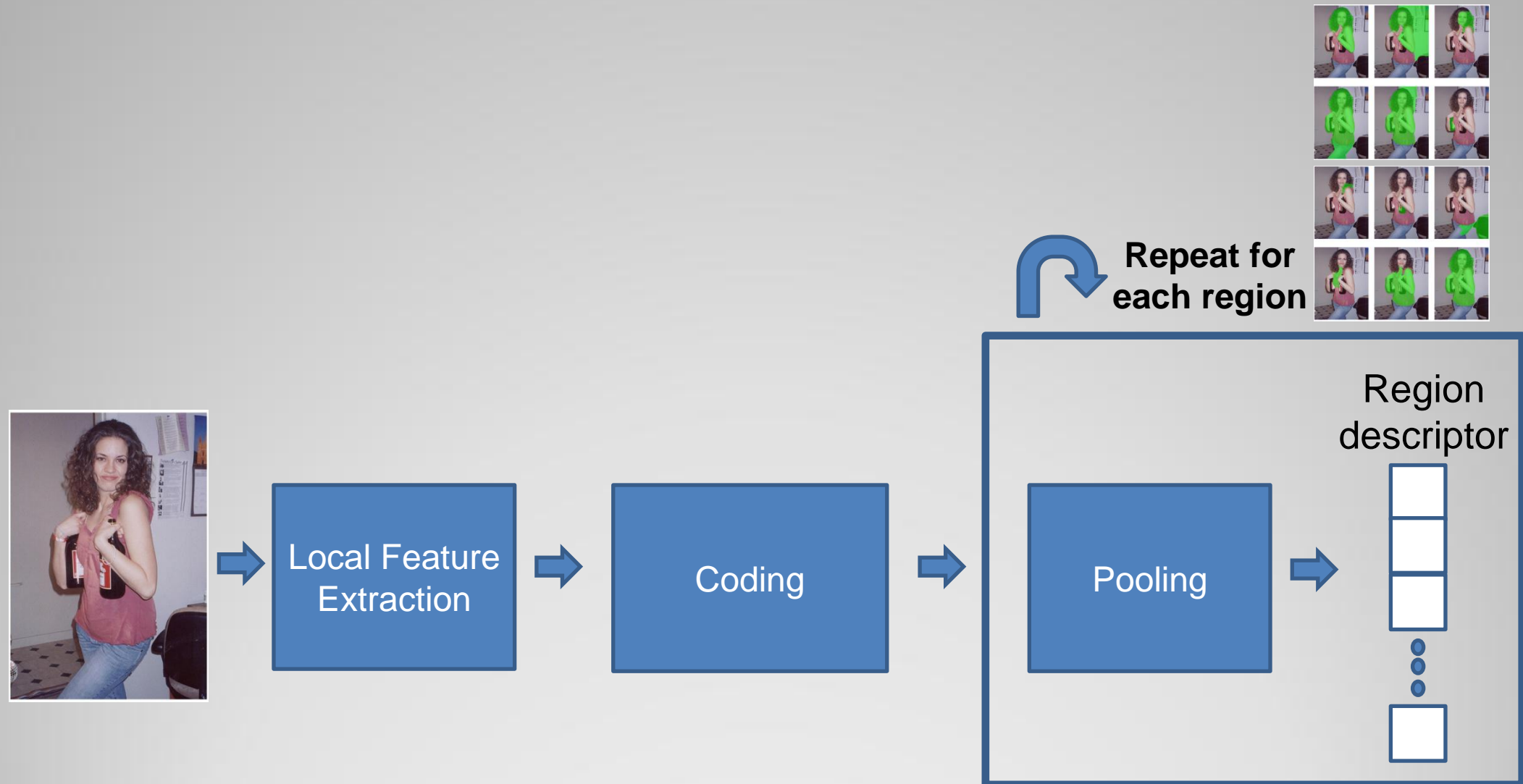
Currently, most successful approaches use variations of Bag of Words (BOW) and HOG

- Require **expensive** classifiers with **non-linear kernels**
- Used in sliding-window detection and in image classification

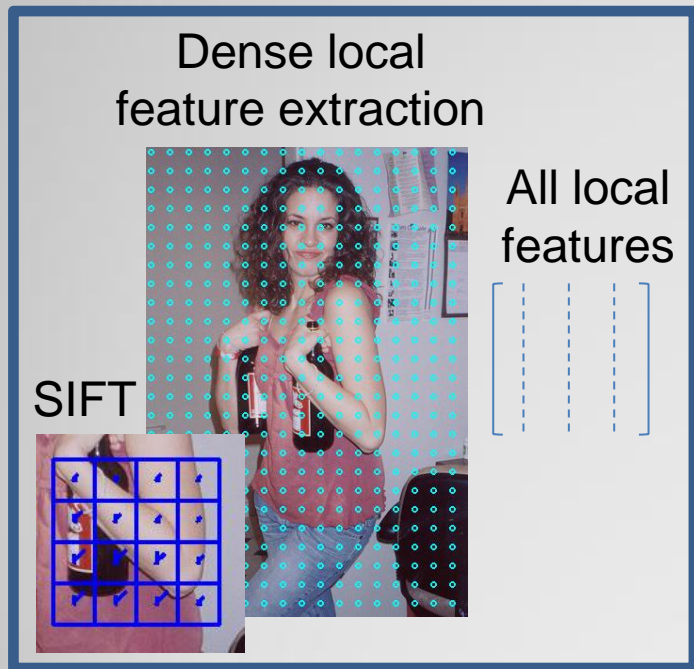
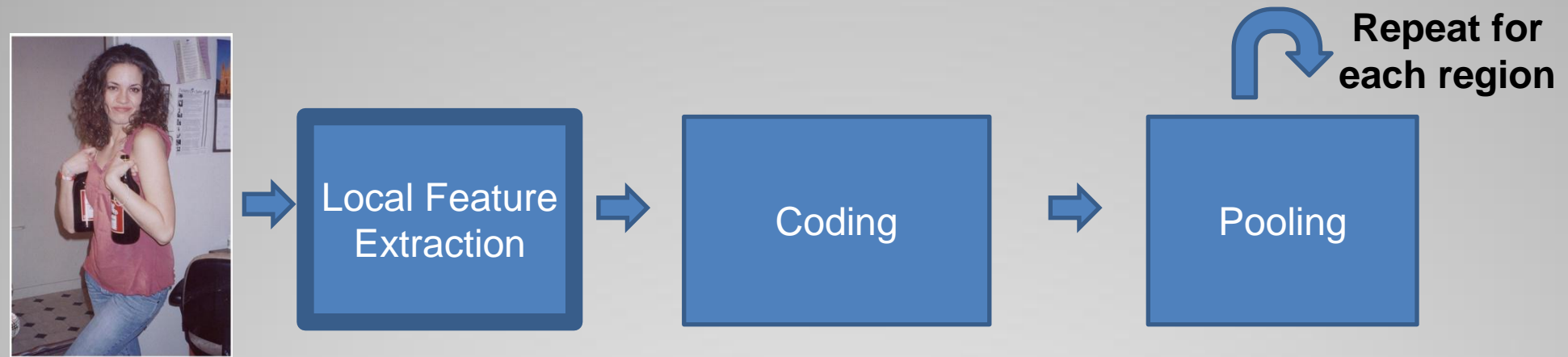


Are there descriptors better suited for regions (segments)?

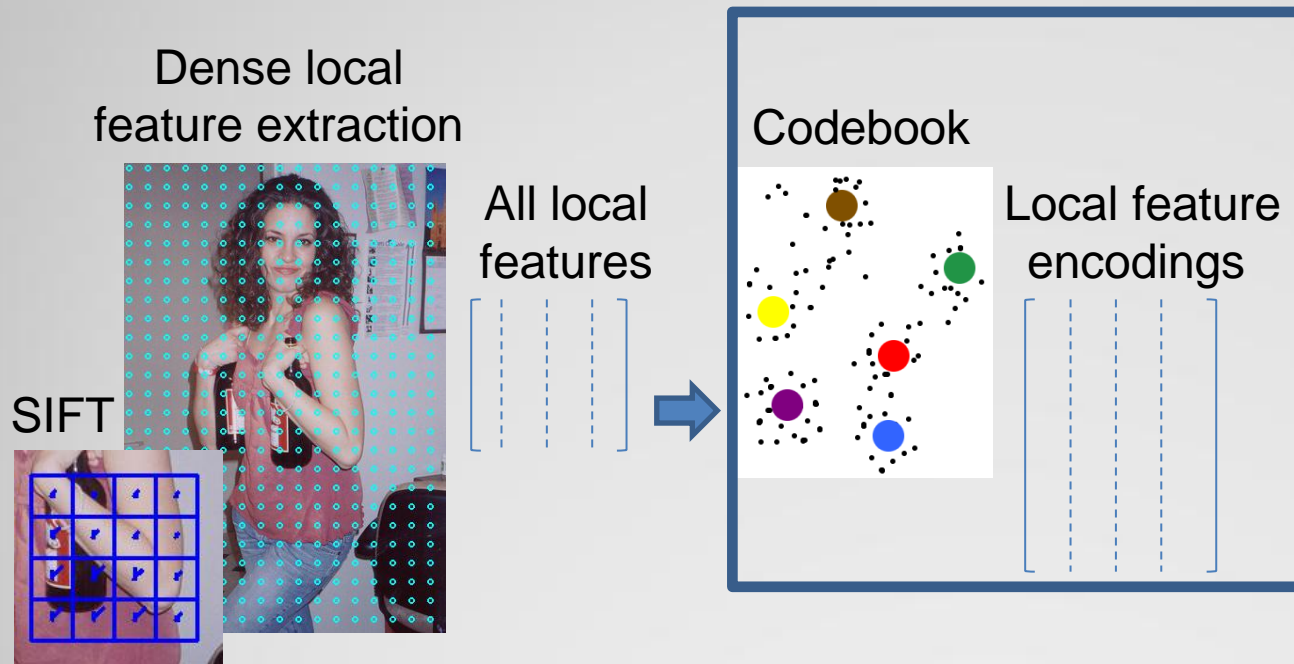
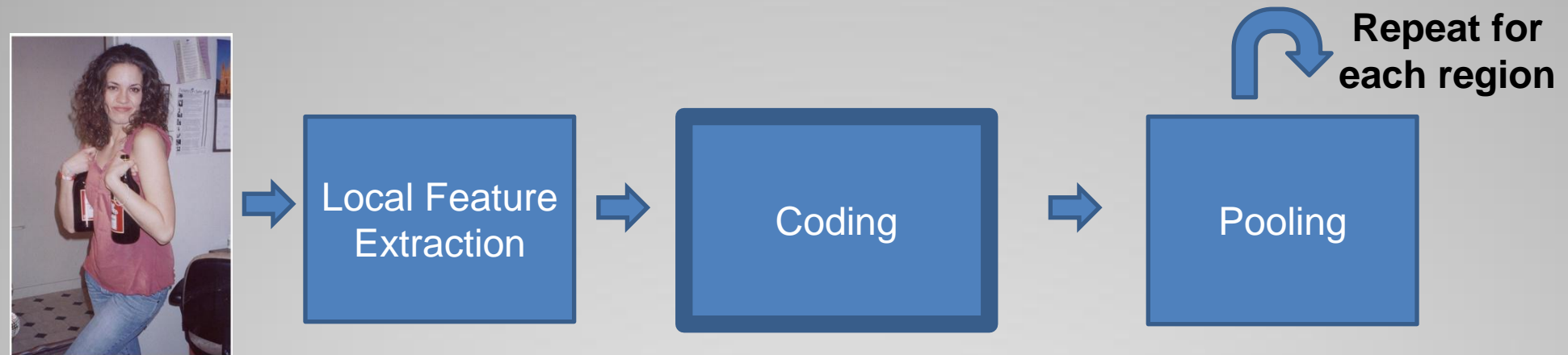
Aggregation-based Descriptors



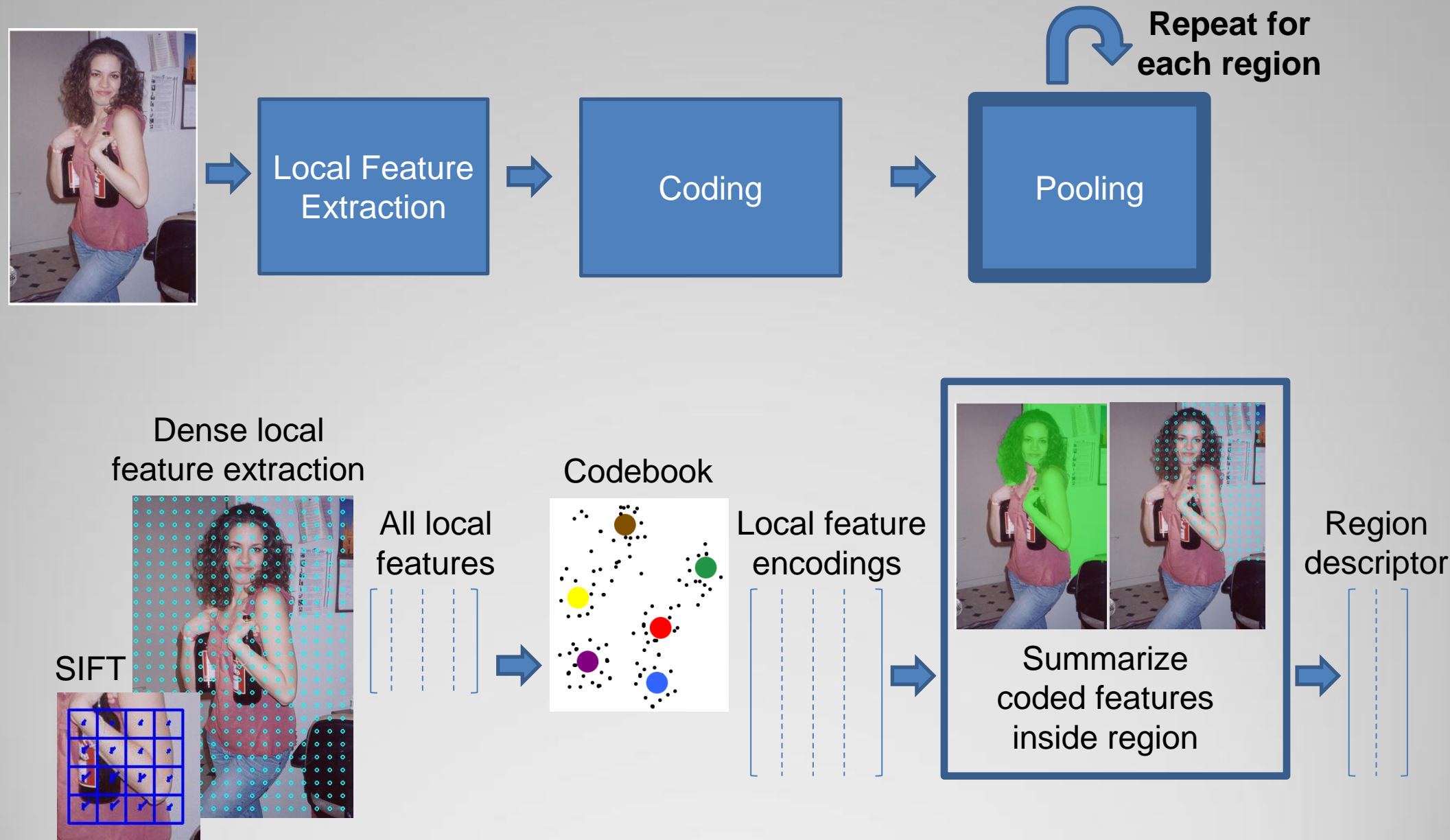
Aggregation-based Descriptors



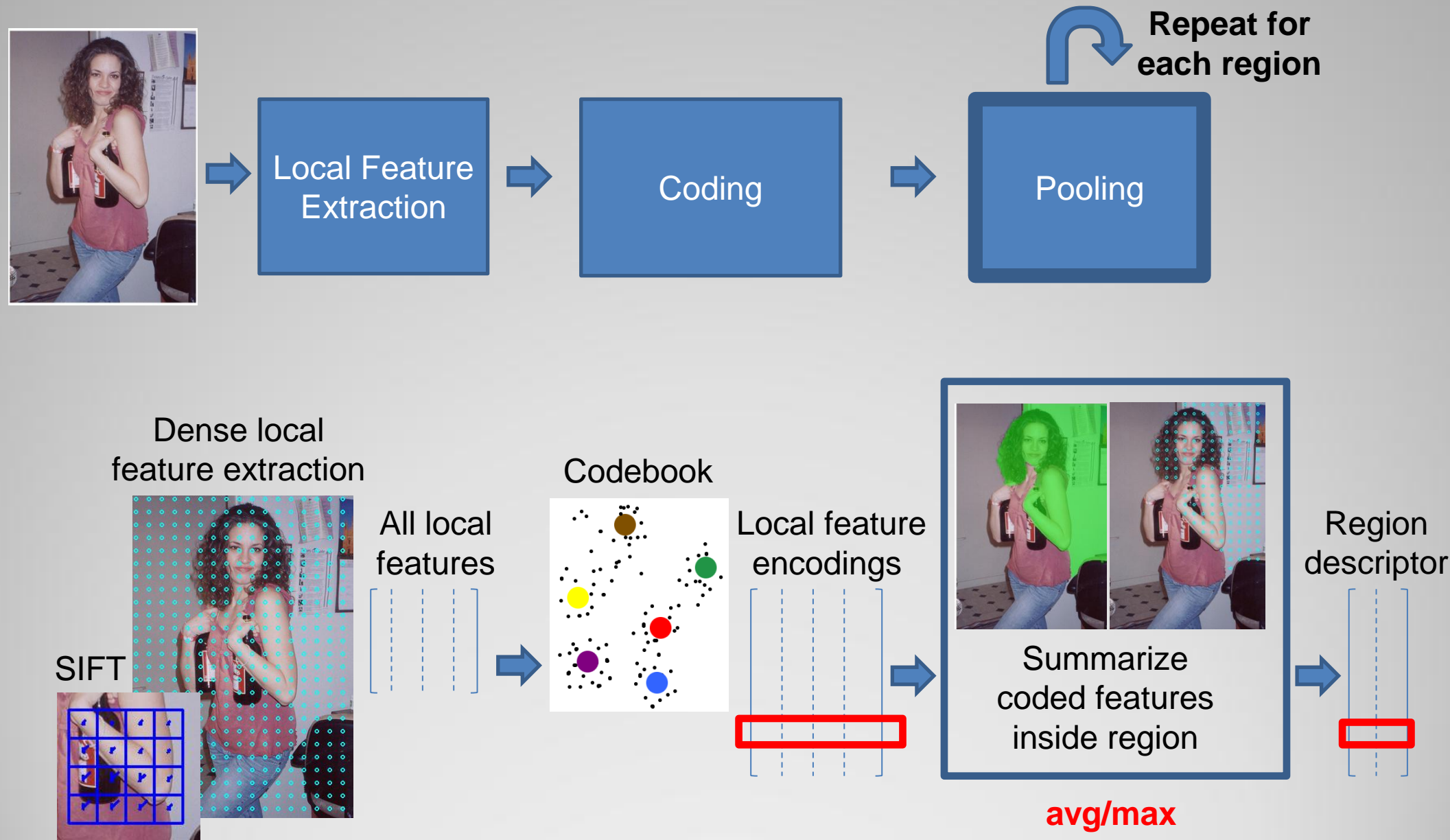
Aggregation-based Descriptors



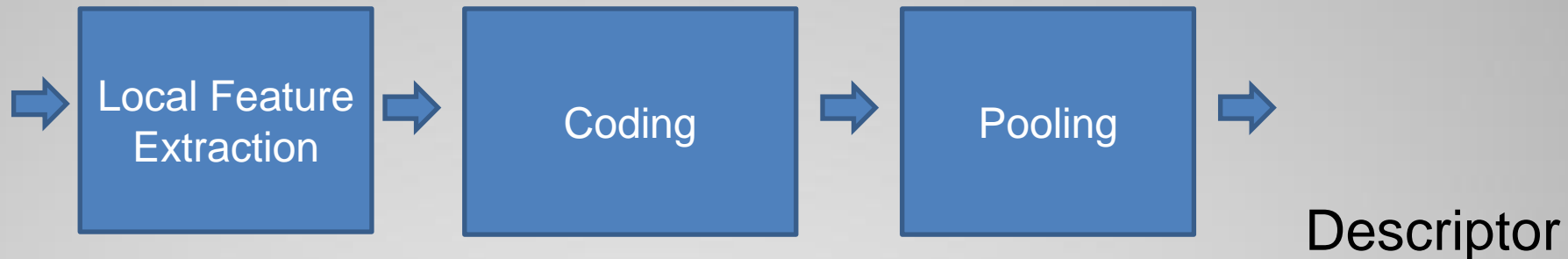
Aggregation-based Descriptors



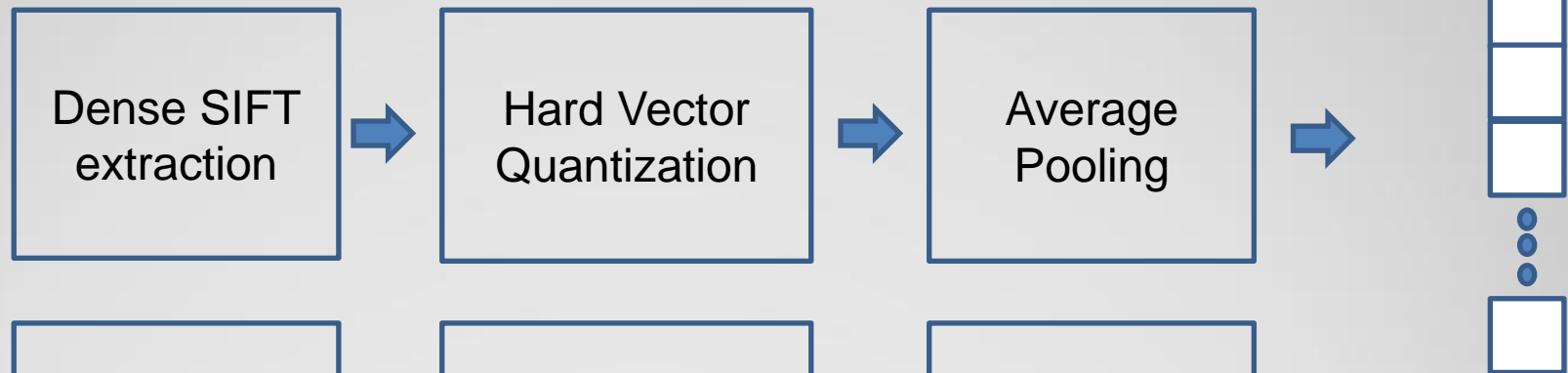
Aggregation-based Descriptors



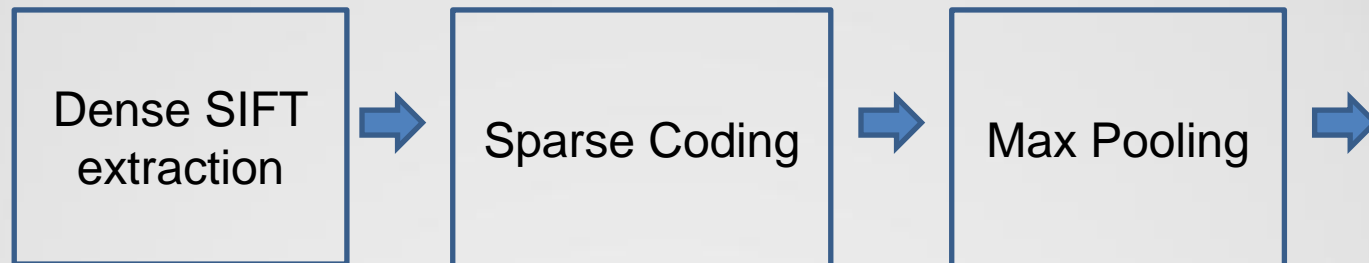
Aggregation-based Descriptors



Bag of Words



Yang et al 09



Aggregation-based Descriptors

Most research so far focused on **coding**

Hard Vector Quantization, Kernel Codebook encoding, Sparse Coding, Fisher encoding, Locality-constrained Linear Coding...

Sivic03, Csurka04, Philbin08, Gemert08, Yang09, Perronnin10, Wang10, (...)

Pooling has received far less attention

Given N local feature descriptors $\mathbf{x}_1, \dots, \mathbf{x}_N$ extracted inside region

Max

$$\mathbf{g}_{max} = \max_i \mathbf{x}_i$$

Average

$$\mathbf{g}_{avg} = \frac{1}{N} \sum_i^N \mathbf{x}_i$$

Second-Order Pooling

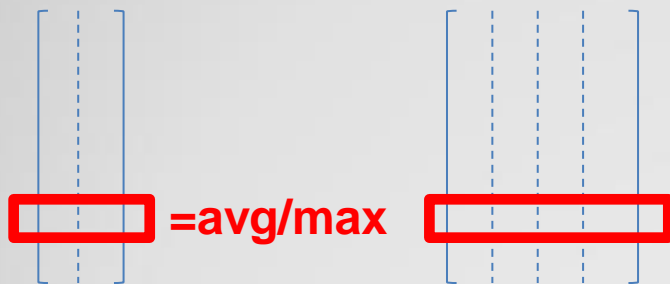
Can we pursue richer statistics for pooling ?

Second-Order Pooling

Can we pursue richer statistics for pooling ?

$$\mathbf{g}_{avg} = \frac{1}{N} \sum_i^N \mathbf{x}_i$$

$$\mathbf{g}_{max} = \max_i \mathbf{x}_i$$



Given N local feature descriptors $\mathbf{x}_1, \dots, \mathbf{x}_N$ extracted inside region

Second-Order Pooling

Can we pursue richer statistics for pooling ?

Capture correlations

$$\mathbf{g}_{avg} = \frac{1}{N} \sum_i \mathbf{x}_i$$

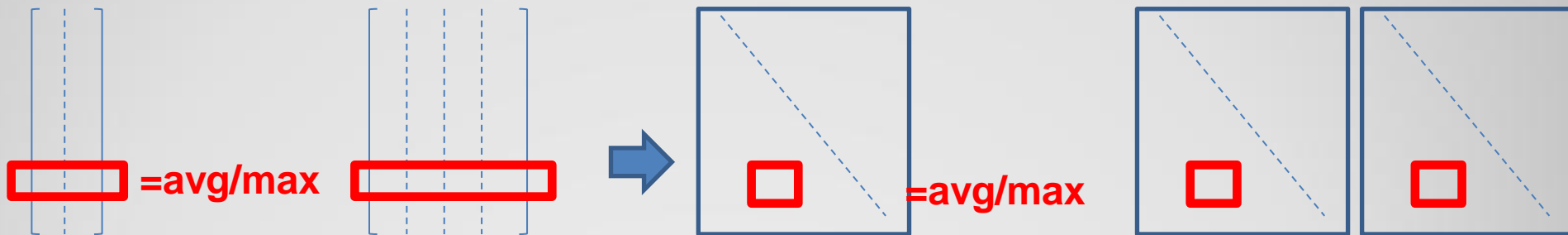


$$\mathbf{G}_{avg} = \frac{1}{N} \sum_i \mathbf{x}_i \cdot \mathbf{x}_i^T$$

$$\mathbf{g}_{max} = \max_i \mathbf{x}_i$$



$$\mathbf{G}_{max} = \max_i \mathbf{x}_i \cdot \mathbf{x}_i^T$$



$$\mathbf{g} = \frac{1}{N} \sum_i \mathbf{x}_i$$

Second-Order Pooling

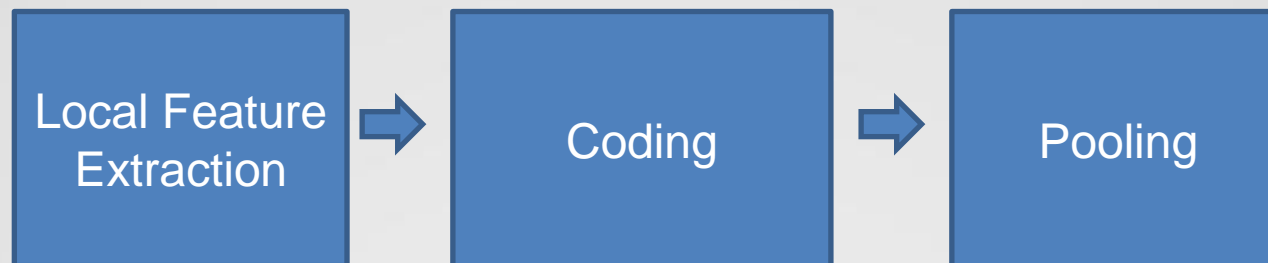
Can we pursue richer statistics for pooling ?

Capture correlations

$$\mathbf{g}_{avg} = \frac{1}{N} \sum_i^N \mathbf{x}_i \quad \longrightarrow \quad \mathbf{G}_{avg} = \frac{1}{N} \sum_i^N \mathbf{x}_i \cdot \mathbf{x}_i^T$$

$$\mathbf{g}_{max} = \max_i \mathbf{x}_i \quad \longrightarrow \quad \mathbf{G}_{max} = \max_i \mathbf{x}_i \cdot \mathbf{x}_i^T$$

Dimensionality = (local descriptor size)²



Second-Order Pooling

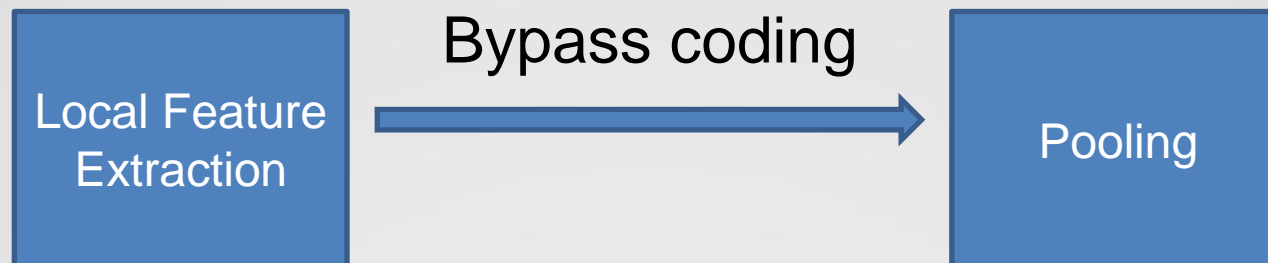
Can we pursue richer statistics for pooling ?

Capture correlations

$$\mathbf{g}_{avg} = \frac{1}{N} \sum_i^N \mathbf{x}_i \quad \longrightarrow \quad \mathbf{G}_{avg} = \frac{1}{N} \sum_i^N \mathbf{x}_i \cdot \mathbf{x}_i^T$$

$$\mathbf{g}_{max} = \max_i \mathbf{x}_i \quad \longrightarrow \quad \mathbf{G}_{max} = \max_i \mathbf{x}_i \cdot \mathbf{x}_i^T$$

Dimensionality = (local descriptor size)²



Second-Order Pooling

What can we say about these matrices ?

$$\mathbf{G}_{avg} = \frac{1}{N} \sum_i^N \mathbf{x}_i \cdot \mathbf{x}_i^T$$

$$\mathbf{G}_{max} = \max_i \mathbf{x}_i \cdot \mathbf{x}_i^T$$

Second-Order Pooling

What can we say about these matrices ?

$$\mathbf{G}_{avg} = \frac{1}{N} \sum_i^N \mathbf{x}_i \cdot \mathbf{x}_i^T \quad \text{Symmetric}$$

$$\mathbf{G}_{max} = \max_i \mathbf{x}_i \cdot \mathbf{x}_i^T \quad \text{Symmetric}$$

... so we can simply keep upper triangle

Second-Order Pooling

What can we say about these matrices ?

$$\mathbf{G}_{avg} = \frac{1}{N} \sum_i^N \mathbf{x}_i \cdot \mathbf{x}_i^T \quad \text{Symmetric Positive Definite (SPD)}$$

$$\mathbf{G}_{max} = \max_i \mathbf{x}_i \cdot \mathbf{x}_i^T \quad \text{Symmetric}$$

SPD matrices have rich geometry: they form a **Riemannian manifold**

Second-Order Pooling

What can we say about these matrices ?

$$\mathbf{G}_{avg} = \frac{1}{N} \sum_i^N \mathbf{x}_i \cdot \mathbf{x}_i^T \quad \text{Symmetric Positive Definite (SPD)}$$

$$\mathbf{G}_{max} = \max_i \mathbf{x}_i \cdot \mathbf{x}_i^T \quad \text{Symmetric}$$

SPD matrices have rich geometry: they form a **Riemannian manifold**

- Linear classifiers ignore this additional geometry

Embedding SPD Manifold in Euclidean Space

Usual solution is to flatten the manifold by projecting to local tangent spaces

- Only valid in a local neighborhood, in general

Embedding SPD Manifold in an Euclidean Space

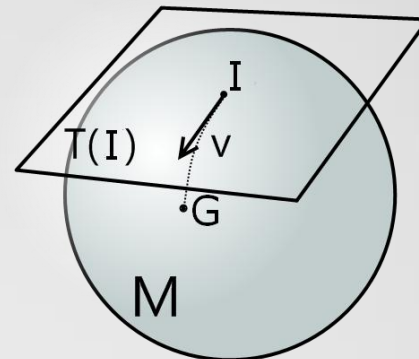
Usual solution is to flatten the manifold by projecting to local tangent spaces

- Only valid in a local neighborhood, in general

By using special, **Log-Euclidean metric**, it is possible to directly embed entire manifold

(Arsigny et al. 07)

$$\mathbf{G}_{log} = \log(\mathbf{G})$$



Sequence of Operations

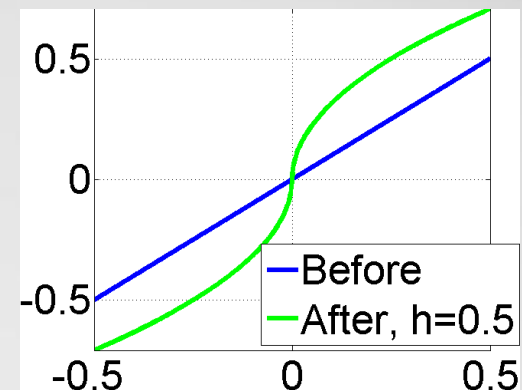
1. Second-Order **Avg** Pooling: $\log \left(\frac{1}{N} \sum_i^N \mathbf{x}_i \cdot \mathbf{x}_i^T \right)$

Second-Order **Max** Pooling: $\max_i \mathbf{x}_i \cdot \mathbf{x}_i^T$

2. Select upper triangle and convert to vector

3. Power normalize (Perronnin *et al* 2010)

$$x = \text{sign}(x) \cdot |x|^h, \text{ with } h \in [0,1]$$



Sequence of Operations

1. Second-Order **Avg** Pooling: $\log \left(\frac{1}{N} \sum_i^N \mathbf{x}_i \cdot \mathbf{x}_i^T \right)$

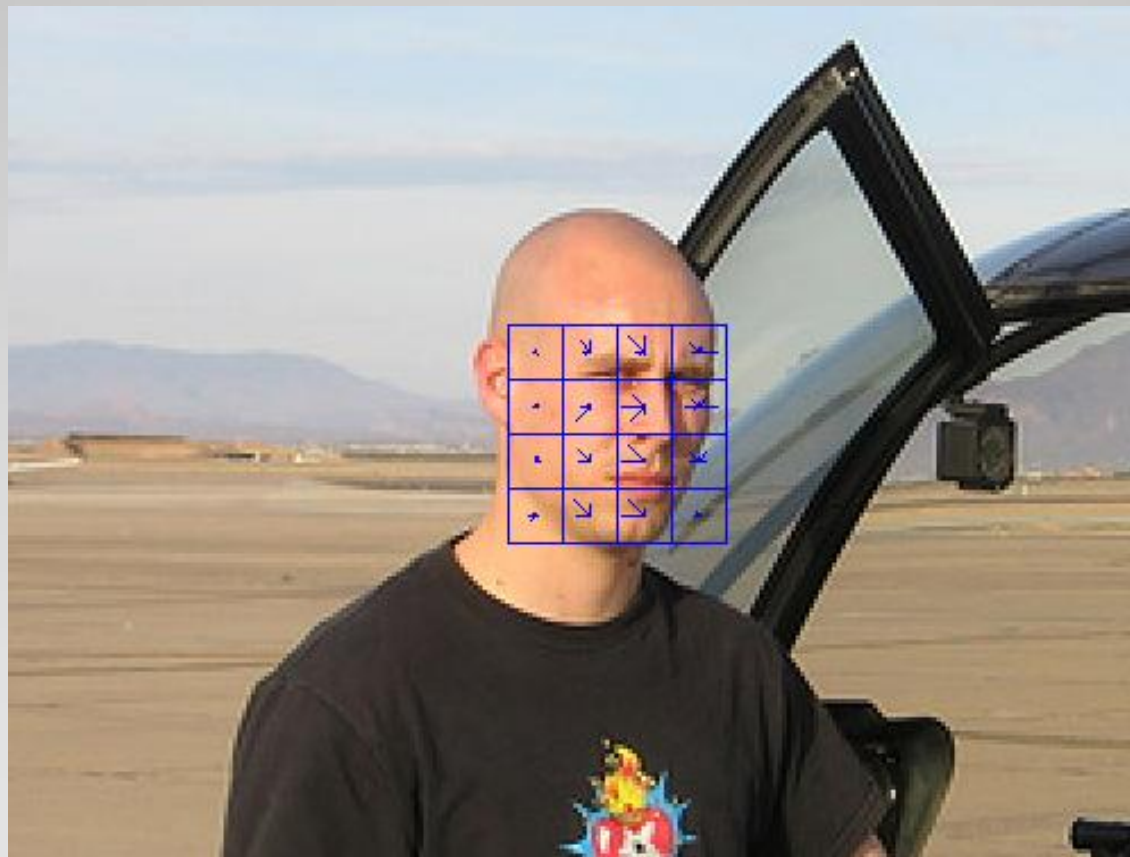
Second-Order **Max** Pooling: $\max_i \mathbf{x}_i \cdot \mathbf{x}_i^T$

2. Select upper triangle and convert to vector

3. Power normalize

Feed resulting descriptor to linear classifier

Additionally we use better local descriptors with pooling methods

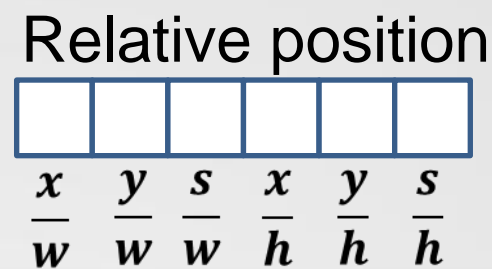
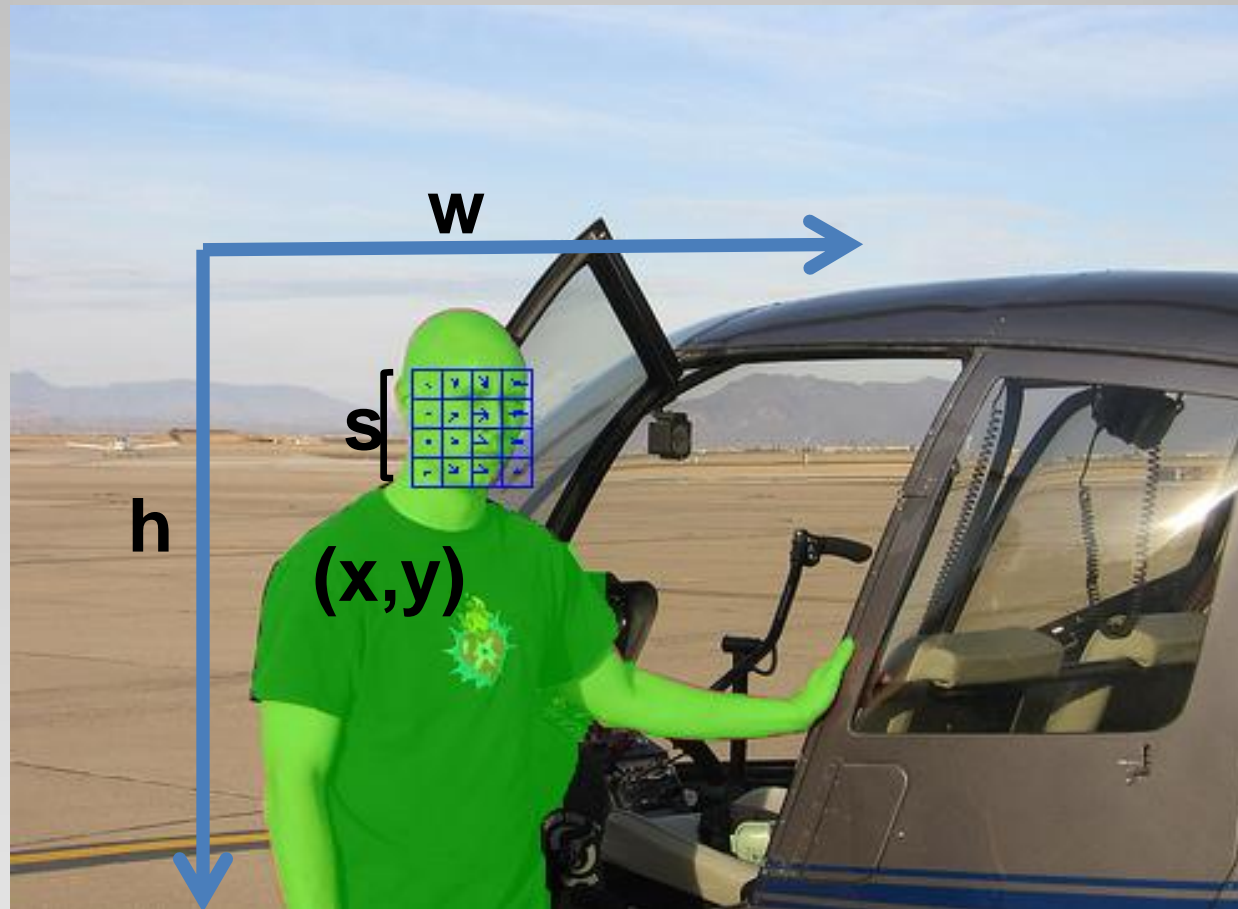


SIFT



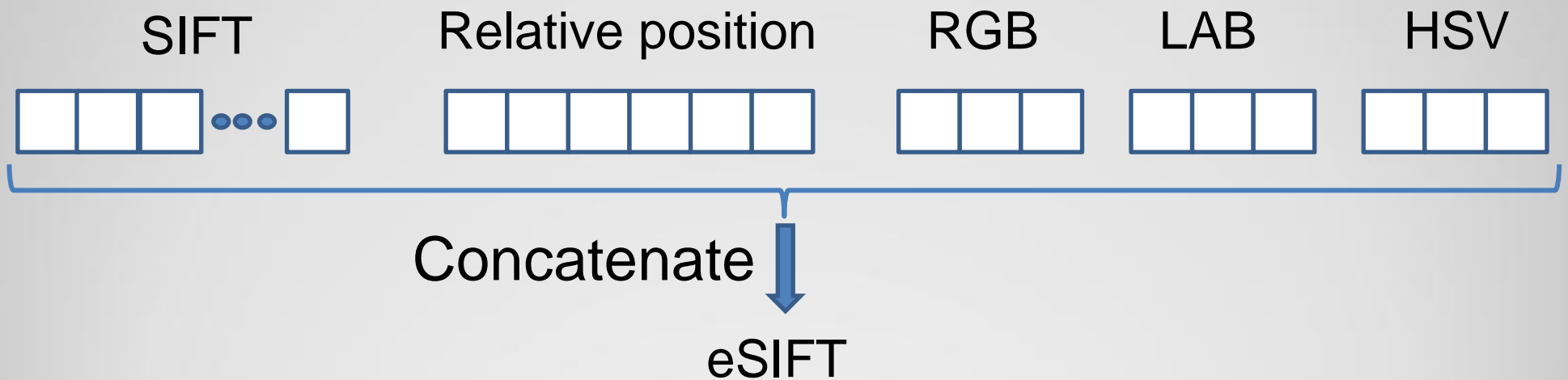
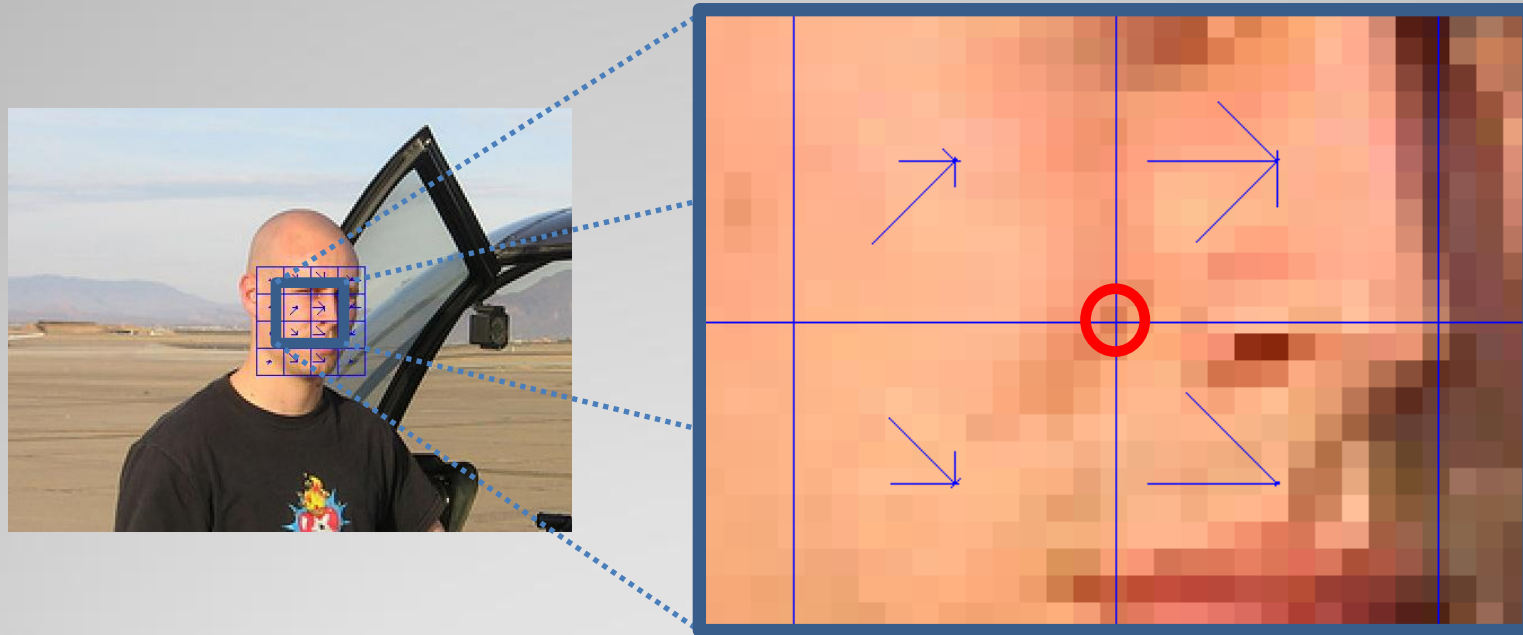
Local Feature Enrichment (1/2)

Relative Position

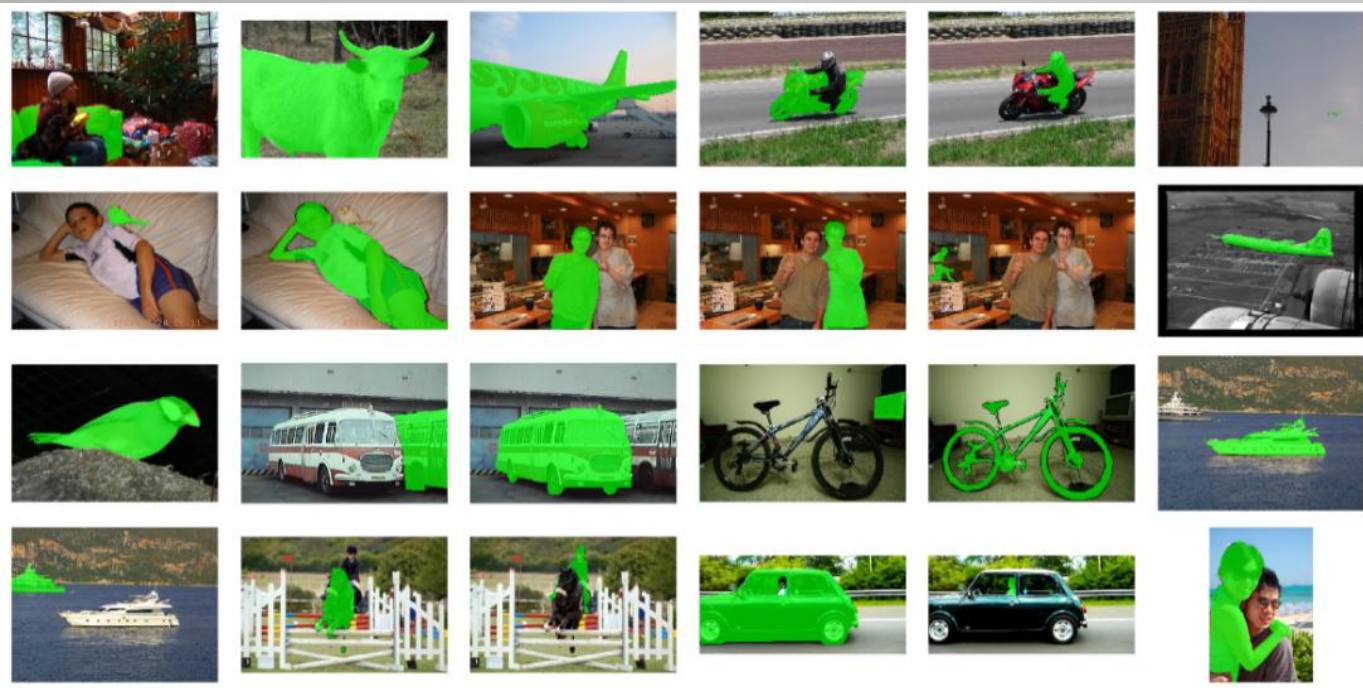


Local Feature Enrichment (2/2)

Pixel Color

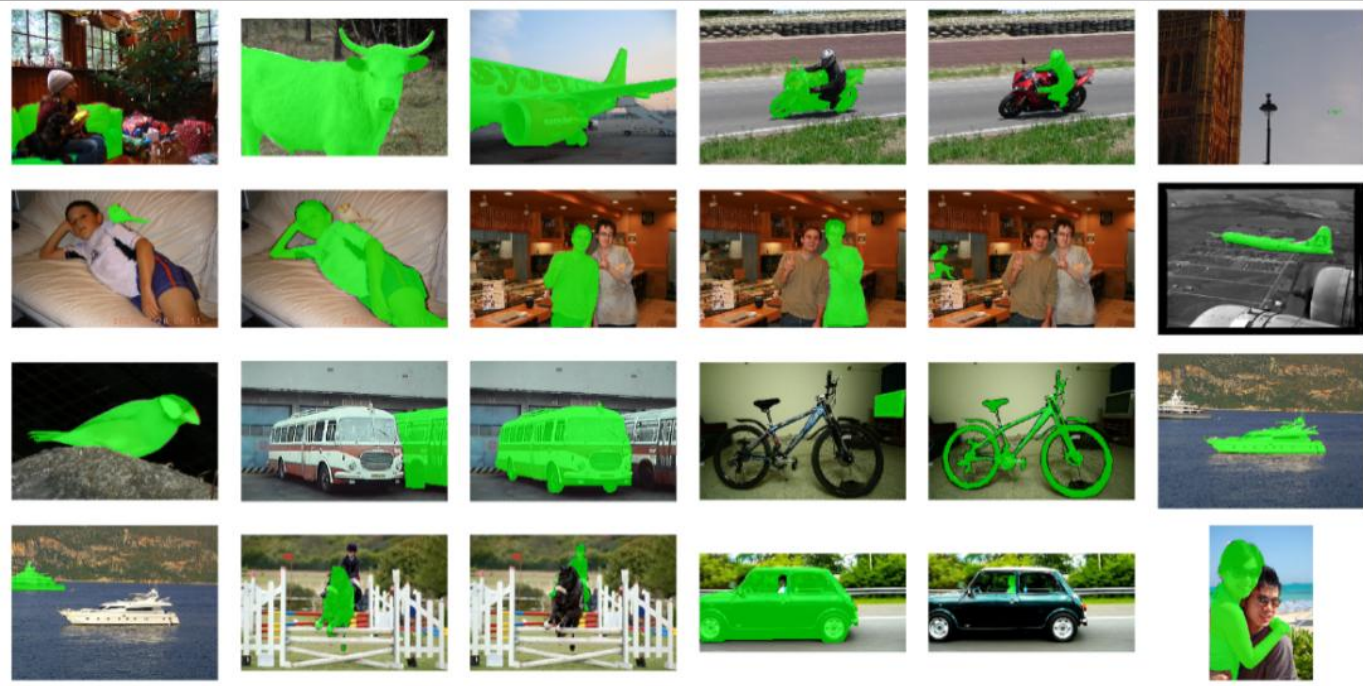


Region Classification VOC 2011



Ground truth regions

Region Classification VOC 2011

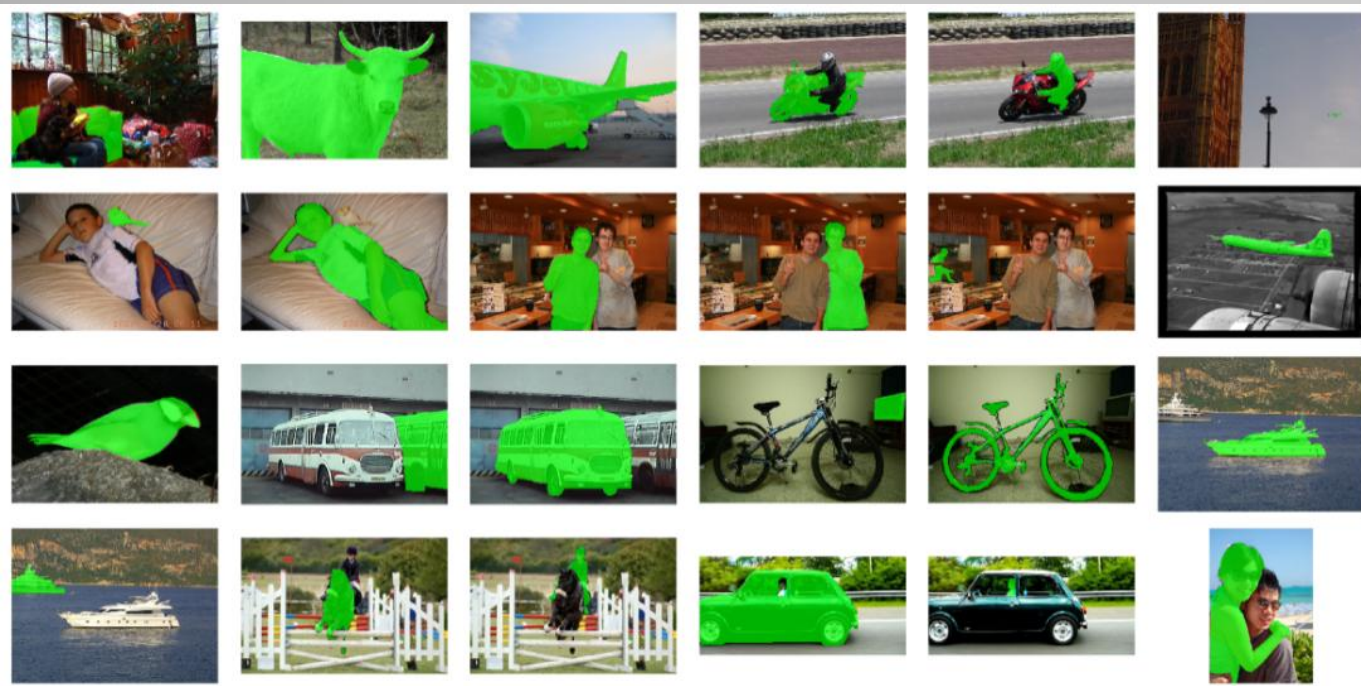


Ground truth regions

Linear classification accuracy

HOG: **41.79%**

Region Classification VOC 2011

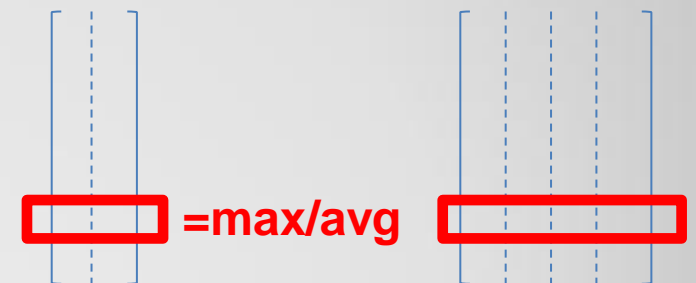


Ground truth regions

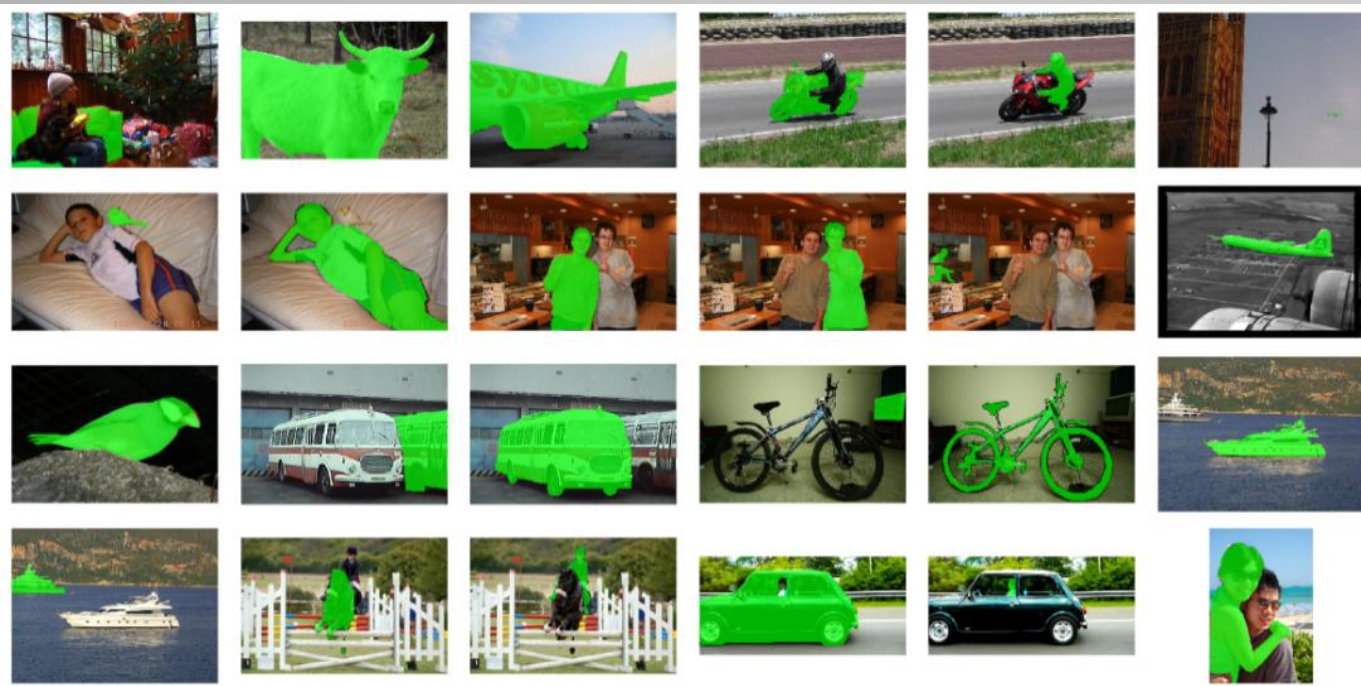
Linear classification accuracy

	1MaxP	1AvgP	2MaxP	2AvgP	Log 2AvgP
SIFT	16.61	33.92			
eSIFT	26.00	43.33			

HOG: 41.79%



Region Classification VOC2011

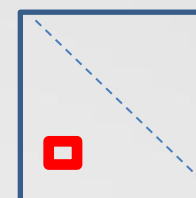


Ground truth regions

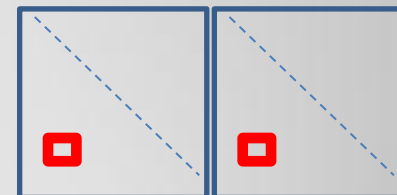
Linear classification accuracy

	1MaxP	1AvgP	2MaxP	2AvgP	Log 2AvgP
SIFT	16.61	33.92	38.74	48.74	
eSIFT	26.00	43.33	50.16	54.30	

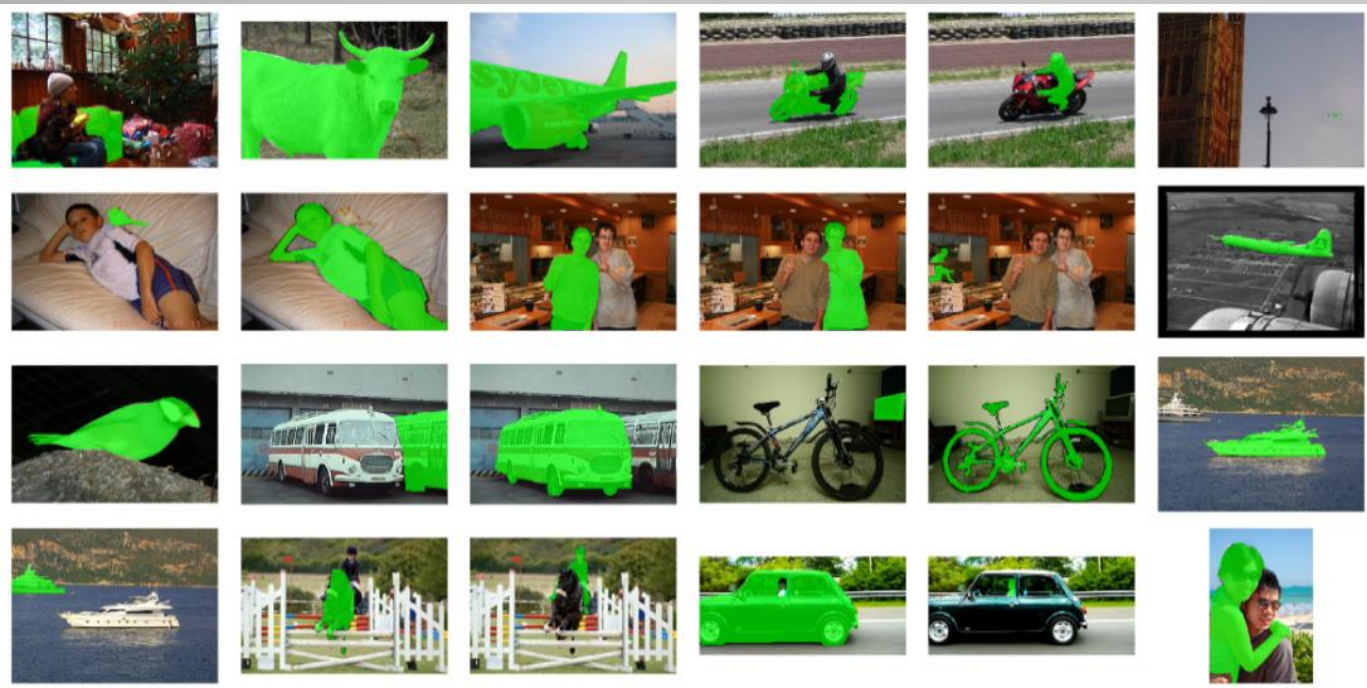
HOG: 41.79%



= max/avg



Region Classification VOC 2011



Ground truth regions

Linear classification accuracy

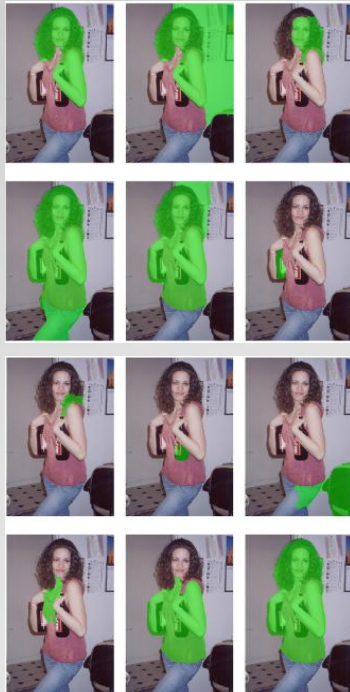
	1MaxP	1AvgP	2MaxP	2AvgP	Log 2AvgP
SIFT	16.61	33.92	38.74	48.74	54.17
eSIFT	26.00	43.33	50.16	54.30	63.85

HOG: 41.79%

$$\log \left(\frac{1}{N} \sum_i^N \mathbf{x}_i \cdot \mathbf{x}_i^T \right)$$

Semantic Segmentation in the Wild Pascal VOC 2011

CPMC



Thresholding + reverse
score overlaying

Local Descriptors:

- eSIFT
- Masked eSIFT
- eLBP

Region Descriptors:

- Second-Order
Average Pooling

Learning:

- LIBLINEAR



Semantic Segmentation in the Wild Pascal VOC 2011

This
work!!



comp6

comp5

	O ₂ P	Berkeley	BONN-FGT	BONN-SVR	BROOKES	NUS-C	NUS-S
Mean Score	47.6	40.8	41.4	43.3	31.3	35.1	37.7
N classes best	13	1	2	4	0	0	1

O₂P best on 13 out of 21 categories: background, aeroplane, boat, bus, motorbike, car, train, cat, dog, horse, potted plant, sofa, person

Semantic Segmentation in the Wild

Pascal VOC 2011

comp6

comp5

	O ₂ P	Berkeley	BONN-FGT	BONN-SVR	BROOKES	NUS-C	NUS-S
Mean Score	47.6	40.8	41.4	43.3	31.3	35.1	37.7
N classes best	13	1	2	4	0	0	1

Linear

Exp-Chi² kernels

Semantic Segmentation in the Wild

Pascal VOC 2011

comp6

comp5

	O ₂ P	Berkeley	BONN-FGT	BONN-SVR	BROOKES	NUS-C	NUS-S
Mean Score	47.6	40.8	41.4	43.3	31.3	35.1	37.7
N classes best	13	1	2	4	0	0	1



Linear



Exp-Chi² kernels

	Feature Extraction	Prediction	Learning
Exp-Chi ²	7.8s / image	87s / image	59h / class
O ₂ P	4.4s / image	0.004s / image	26m / class

20,000x faster **130x faster**

Caltech 101

Important testbed for coding and pooling techniques



Caltech 101

Important testbed for coding and pooling techniques



- No segments, spatial pyramid instead
- Linear classification

Caltech 101

Important testbed for coding and pooling techniques



	SIFT-O ₂ P	eSIFT-O ₂ P	SPM ¹	LLC ²	EMK ³	MP ⁴
Accuracy	79.2	80.8	64.4	73.4	74.5	77.3

1. Lazebnik et al. '06
2. Wang et al. '10
3. Bo & Sminchisescu '10
4. Boureau et al. '11

Conclusions

- Second-order pooling with Log-Euclidean tangent space mappings
- Practical aggregation-based descriptors without unsupervised learning stage (no codebooks)
- High recognition performance on free-form regions using linear classifiers
- Semantic Segmentation on VOC 2011 superior to state-of-the-art with models 20,000x faster

Code available online

Thank you!

