# From Rigid Templates to Grammars: Object Detection with Structured Models

#### Ross B. Girshick

Dissertation defense April 20, 2012

#### The question



#### What objects are where?

# Why it matters

#### Intellectual curiosity

- How do we extract this information from the signal?

#### Applications

- Semantic image and video search
- Human-computer interaction (*e.g.*, Kinect)
- Automotive safety
- Camera focus-by-detection
- Surveillance
- Semantic image and video editing
- Assistive technologies
- Medical imaging

# Proxy task: PASCAL VOC Challenge

- Localize & name (*detect*) 20 basic-level object categories
  - Airplane, bicycle, bus, cat, car, dog, person, sheep, sofa, monitor, etc.



- 11k training images with 500 to 8000 instances / category
- Evaluation: bounding-box overlap; average precision (AP)





- Deformation
- Viewpoint
- Subcategory
- Variable structure
- Occlusion
- Background clutter
- Photometric

#### Deformation



Image credit: http://i173.photobucket.com/albums/w78/yahoozy/MultipleExposures2.jpg

Viewpoint



Image credits: PASCAL VOC



#### Subcategory — "airplane" images



0

Variable structure









# PASCAL VOC Challenges 2007-2011

#### 2007 Challenge

- Winner: Deformable part models & Latent SVM [FMR'08]
- 21% mAP

   Baseline for dissertation
   Prior work

   Winners of 2008 & 2009 Challenges

   This work
   ↓

   Fast forward to the 2011 Challenge
  - Our system (voc-release4): 34% mAP
  - Top system (NLPR): 41% mAP
  - NLPR method: voc-release4 + LBP image features + richer spatial model (GMM) + more context rescoring
  - Second (MIT-UCLA) and third place (Oxford) also based on voc-release4

#### Contributions — By area

- Object representation\*
  - Mixture models (in PAMI'10); Latent orientation; Person grammar model
- Efficient detection algorithms\*
  - Cascaded detection for DPM (oral at CVPR'10)
- Learning\*
  - Weak-label structural SVM (spotlight at NIPS'11)
- Detection post-processing
  - Bounding box prediction & context rescoring
- Image representation
  - Enhanced HOG features; features for boundary truncation & small objects
- Software
  - voc-release  $\{2,3,4\}$  currently the "go to" object detection system

#### **Object representation**

# Model lineage – Dalal & Triggs



- Histogram of Oriented Gradients (HOG) features
- Scanning window detector (linear filter)
- w learned by SVM

# Model lineage – Latent SVM DPM



- Dalal & Triggs + Parts in a deformable configuration z
- Scanning window detection: max over z at each p<sub>0</sub>
- w learned by latent SVM

# Superposition of views





#### Mixture of DPMs



#### Training (component labels are hidden)

- Cluster training examples by bounding-box aspect ratio
- Initialize root filters for each component (cluster) independently
- Merge components into mixture model and train with latent SVM

#### Mixtures with latent orientation

("pushmi-pullyu" instead of horse)







#### Learning without latent orientation









Learning with latent orientation

[GFM voc-release4]

#### Unsupervised orientation clustering

Online clustering with a hard constraint



Assign i<sup>th</sup> example to nearest cluster Flipped example *must* go to the other cluster

#### Latent orientation improves performance



#### Results – Mixture models and latent orientation

- Mixture models boost mAP by 3.7 points
- Latent orientation boost mAP by 2.6 AP points

	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	$\mathbf{tv}$	mAP
LSVM CVPR [34]	18.0	41.1	9.2	9.8	24.9	34.9	39.6	11.0	15.5	16.5	11.0	6.2	30.1	33.7	26.7	14.0	14.1	15.6	20.6	33.6	22.3
					voc-	relea	se4.	01: E	HOG	$+ \mathbf{T}$	RUNG	C + 1	BBOX	+ OP'	TIM +	MAX-	REG				
	aero	bike	bird	boat	bottle	bus	car	$\mathbf{cat}$	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	$\mathbf{tv}$	mAP
star model	21.5	49.7	9.3	9.3	27.3	39.8	53.3	11.4	13.5	17.8	28.6	3.7	42.1	38.3	33.1	13.4	15.1	20.4	28.7	42.7	26.0
MIX	33.2	58.0	2.3	15.7	26.5	49.3	55.5	13.9	18.1	19.1	24.6	12.6	47.3	42.6	40.5	14.4	17.0	27.4	37.9	39.0	29.7
$\cdots$ + <b>ORIENT</b>	28.9	59.5	10.0	15.2	25.5	49.6	57.9	19.3	22.4	25.2	23.3	11.1	56.8	48.7	41.9	12.2	17.8	33.6	45.1	41.6	32.3
$\cdots$ + CNTXT	31.2	61.5	11.9	17.4	27.0	49.1	59.6	23.1	23.0	26.3	24.9	12.9	60.1	51.0	43.2	13.4	18.8	36.2	49.1	43.0	34.1

AP scores using the PASCAL 2007 evaluation

• 12 AP point improvement (>50% relative) over the baseline

#### Efficient detection

#### Cascaded detection for DPM

Add in parts one-by-one and prune partial scores



Sparse dynamic programming tables (reuse computation!)



#### Threshold selection & PCA filters

- Data-driven threshold selection
  - Based on statistics of partial scores on training data
  - Provably safe ("probably approximately admissible" thresholds)
  - Empirically effective
- 2-stage cascade with simplified appearance models
  - Use PCA of HOG features (or model filters)
  - Stage 1: place low-dimensional filters; Stage 2: place original filters



# Results — 15x speedup (no loss in mAP)

#### High recall



#### Lower recall $\Rightarrow$ faster



#### Towards richer grammar models

#### People are complicated



Helmet, occluded left side

Ski cap, no face, truncated

Pirate hat, dresses, long hair

Truncation, holding glass, heavy occlusion

# Objects from visually rich categories have diverse *structural* variation

#### Compositional models



#### More mixture components?



There are too many combinations! Instead...

... compositional models defined by grammars

#### **Object detection grammars**

- A modeling language for building object detectors [FM'10]
  - Terminals (model image data appearance)
  - Nonterminals (objects, parts, ...)
  - Weighted production rules (define composition, variable structure)
- Composition
  - Objects are recursively composed of other objects (parts)
- Variable structure
  - Expanding different rules produces different structures
- Person  $\rightarrow$  Head, Torso, Arms, Legs
- Head  $\rightarrow$  Eye, Eye, Mouth
- Mouth  $\rightarrow$  Smile OR Mouth  $\rightarrow$  Frown

#### **Object detection grammars**

Object hypothesis = derivation tree T



Linear score function

$$\operatorname{score}(I, T) = \mathbf{w} \cdot \boldsymbol{\psi}(I, T)$$

Detection with DP

$$T^*(p) = \operatorname*{argmax}_{T \in \mathcal{T}_p} \mathbf{w} \cdot \boldsymbol{\psi}(I, T)$$

#### Build on what works



Can we build a better person detector?

# Case study: a person detection grammar

Subtype 1 Subtype 2





Example detections and derived filters



Parts 1-6 (no occlusion)

Parts 1-4 & occluder

Parts 1-2 & occluder

- Fine-grained occlusion
- Sharing across all derivations
- Model of the stuff that causes occlusion
- Part subtypes and multiple resolutions
- Parts have subparts (not pictured)





# Training models



- PASCAL data: bounding-box labels
- No derivation trees given! (weakly-supervised learning)
- Learn the parameters w

# Defining examples



- Each bounding box is a foreground example
- All locations in background images are background examples
- From these examples, learn the prediction rule



#### Parameter learning

Richer models, richer problems



#### One good output... and many bad ones!

• Which learning framework should we use?

#### **Classification training**

$$E_{\text{lsvm}}(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max\left[0, 1 - y_i \max_{s \in \mathcal{S}(x_i)} \mathbf{w} \cdot \boldsymbol{\psi}(x_i, s)\right]$$
  
**Training:**  
LSVM objective:  
"score +1 here"  

$$LSVM \text{ objective:}$$
"score +1 here"

**Testing:** 







Who wins? Both derivations were trained to score +1.

#### Structured output training



#### Training:

"outscore *all* other outputs by a margin"

"score lower by a margin"



A "good" output should win.

#### Latent structural SVM

[Yu and Joachims]

$$E(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n L_{\text{lssvm}}(\mathbf{w}, x_i, y_i)$$

$$L_{\text{Issvm}}(\mathbf{w}, x, y) = \max_{\substack{(\hat{y}, \hat{z}) \in \mathcal{Y} \times \mathcal{Z} \\ \hat{z} \in \mathcal{Z}}} \left[ \mathbf{w} \cdot \boldsymbol{\psi}(x, \hat{y}, \hat{z}) + L_{\text{margin}}(y, \hat{y}) \right]$$

Objective and task loss (L<sub>margin</sub>) might be inconsistent

- Many outputs with zero loss — LSSVM "requires" the training label

#### LSSVM requires label space = output space

- A simple example where label space != output space
  - Label space is *all pixel-accurate* bounding boxes
  - Outputs are bounding boxes on a low-res. grid at some scales
  - Does not naturally fit the LSSVM framework



# Structured learning desiderata

- Model can make any low-loss prediction
  - Many outputs might be compatible with one label
  - The model is free to choose between them
- Label space and output space can be different
  - *E.g.*, bounding boxes labels and derivation tree outputs
- Generalize frameworks that work well
  - Structural SVM
  - Latent structural SVM
  - Latent SVM



#### Label space != output space

Allowing different label spaces and output spaces



 $\mathsf{label} \in \mathcal{Y}$ 



Connect the spaces with loss functions of the form

$$L: \mathcal{Y} \times \mathcal{S} \to \mathbb{R}_{\geq 0}$$

#### Weak-label structural SVM

$$E(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n L(\mathbf{w}, x_i, y_i)$$

$$L(\mathbf{w}, x, y) = \max_{s \in S(x)} \left[ \mathbf{w} \cdot \boldsymbol{\psi}(x, s) + L_{\text{margin}}(y, s) \right] - \max_{s \in S(x)} \left[ \mathbf{w} \cdot \boldsymbol{\psi}(x, s) - L_{\text{output}}(y, s) \right]$$

- Allows different label spaces and output spaces
- Not "required" to predict the training label
  - Many outputs may be compatible with a label labels are "weak"
  - The model can pick any output with low L<sub>output</sub>
- Generalizes many frameworks
  - SSVM, LSSVM, LSVM, structural ramp loss

Person grammar results

	Grammar	+bbox	+ context	<b>UoC-TTI</b> [26]	+bbox	+ context	Poselets [26]
AP	47.5	47.6	<b>49.5</b>	44.4	45.2	47.5	48.5

AP scores on PASCAL 2010

WL-SSVM vs. LSVM

Person grammar	1	AP scores on splits							
objective	mAP	1	2	3	4	5			
LSVM WL-SSVM	46.0 48.0	$47.5 \\ 49.0$	$\begin{array}{c} 46.3\\ 48.4\end{array}$	$\begin{array}{c} 45.3\\ 48.4\end{array}$	$\begin{array}{c} 45.3\\ 48.0\end{array}$	$\begin{array}{c} 45.6\\ 46.5\end{array}$			

AP scores on 5 PASCAL 2011 train+val splits

#### Example detections



# Summary of contributions

- Richer models + post-processing + features + learning
  - >50% relative improvement in the state-of-the-art
- Cascaded detection for DPM
  - 15x speedup of detection with no loss in performance
- Person detection grammar and WL-SSVM
  - Highest-performing person detector
  - More general & natural learning framework for many problems
- Improved image features
- Detection post-processing
- Software
  - voc-release5 will be available soon!

## **Open directions**

- Grammar structure learning
  - Perhaps from more detailed annotations
- Score compatibility and linear grammars
  - Nonlinearities to normalize score ranges neural grammars?
- Rethink our low-level features
  - HOG features are too coarse to model fine detail
  - We are likely saturating the performance of HOG features
- Optimizing nonconvex objectives with latent variables
  - How can we free ourselves from careful (often model specific) initialization?

Subcategory — "car" images







# Organizing principles

- Gradually build richer models
  - Central research methodology
- Compositional models
  - Object Detection Grammars [FM'10]
  - Deformation, viewpoint, subcategory, composition in a unified framework
- Efficient computation
  - Tree-structured models
  - Cascaded detection
- Train models from weakly-labeled data
  - New models, old annotations

#### Preliminaries — Object detection grammars

Linear score function

$$\begin{aligned} \operatorname{score}(I, T) &= \sum_{(r, z) \in \operatorname{int}(T)} \beta_r(z) + \sum_{A(\omega) \in \operatorname{leaf}(T)} \operatorname{score}(I, A, \omega) \\ &= \sum_{(r, z) \in \operatorname{int}(T)} \mathbf{w} \cdot \phi_r(z) + \sum_{A(\omega) \in \operatorname{leaf}(T)} \mathbf{w} \cdot \phi_A(I, \omega) \\ &= \mathbf{w} \cdot \left[ \sum_{(r, z) \in \operatorname{int}(T)} \phi_r(z) + \sum_{A(\omega) \in \operatorname{leaf}(T)} \phi_A(I, \omega) \right] \\ &= \mathbf{w} \cdot \psi(I, T). \end{aligned}$$

- Detection = find high scoring derivations
  - Efficient dynamic programming algorithm

#### Image representation & features

#### Enhanced HOG features

- 36D  $\rightarrow$  31D with more information
- Contrast sensitive and insensitive



#### Boundary truncation

- Nonzero scores outside the image





#### Small objects

- Scale-dependent score bias



#### Detection post-processing

#### Bounding box prediction



Contextual information

#### Results – Mixture model and latent orientation

- Mixture models boost mAP 3.7 points
- Latent orientation boost mAP 2.6 AP points



AP scores using the PASCAL 2007 evaluation

• 12 AP point improvement (>50% relative) over the baseline





Root filter shapes?



#### Number of parts?

Anchor positions?

Part shapes and sizes?

What are the grammar productions?

Heuristics, cross validation, insight (from humans)

#### Object representation — Summary



[Tsochantaridis et al., Taskar et al.]

No latent variables

$$E(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n L_{ssvm}(\mathbf{w}, x_i, y_i)$$

$$L_{\text{ssvm}}(\mathbf{w}, x, y) = \max_{\hat{y} \in \mathcal{Y}} \left[ \mathbf{w} \cdot \boldsymbol{\psi}(x, \hat{y}) + L_{\text{margin}}(y, \hat{y}) \right] \\ - \mathbf{w} \cdot \boldsymbol{\psi}(x, y)$$

- Objective and task loss (L<sub>margin</sub>) might be inconsistent
  - Two outputs with zero loss SSVM "requires" the training label

$$L_{\text{margin}}(\mathbf{w}, x, y) = 0$$
 and  $L_{\text{margin}}(\mathbf{w}, x, \hat{y}) = 0$ 



Use the convex-concave procedure

$$E(\mathbf{w}) = E_{\text{convex}}(\mathbf{w}) + E_{\text{concave}}(\mathbf{w})$$

$$E_{\text{convex}}(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max_{s \in \mathcal{S}(x_i)} \left[\mathbf{w} \cdot \boldsymbol{\psi}(x_i, s) + L_{\text{margin}}(y_i, s)\right]$$
$$E_{\text{concave}}(\mathbf{w}) = -C \sum_{i=1}^n \max_{s \in \mathcal{S}(x_i)} \left[\mathbf{w} \cdot \boldsymbol{\psi}(x_i, s) - L_{\text{output}}(y_i, s)\right]$$

Sequence of convex slave problems

$$\mathbf{w}_{t+1} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \underset{s \in \mathcal{S}(x_i)}{\max} \left[ \mathbf{w} \cdot \boldsymbol{\psi}(x_i, s) + \mathcal{L}_{\operatorname{margin}}(y_i, s) \right] - \mathbf{w} \cdot \boldsymbol{\psi}(x_i, s_i(\mathbf{w}_t))$$

#### **Object detection grammars**

- A modeling language for building object detectors [FM'10]
  - Terminals (model image data appearance)
  - Nonterminals (objects, parts, ...)
  - Weighted production rules (define compositions, subtypes)
- Composition

Subtypes (choice yields variable structure)

$$X \xrightarrow{\beta_1} \{ W_1, \dots, W_n \}$$
$$X \xrightarrow{\beta_2} \{ Y_1, \dots, Y_m \}$$

- X is a "smiling face" or a "frowning face"
- Symbols are placed Y(p)



#### **Object detection grammars**

- Symbols are placed
- Terminals model appearance (HOG filters)



# Definition of S(x) – foreground examples

Example (x, y) where x specifies an image I and y = B is a bounding box. Let S(x) be:

- 1. Derivations T with overlap $(box(T), B) \ge 0.1$  and overlap(box(T), B') < 0.5 for all B'in I such that  $B' \ne B$ .
- 2. The background output  $\perp$ .



# Definition of S(x) – background examples

Example (x, y) where x specifies an image I and position  $\omega$  and  $y = \bot$  is the background symbol. Let  $S(x) = T_{\omega} \cup \{\bot\}$ .





HOG feature pyramid