# Object and Activity Recognition Grounded on Midlevel Image Representations

Sinisa Todorovic joint work with: N. Payet, M. Amer Oregon State University September 20, 2011

# **Marr's Vision**



Slide by SC Zhu

# **Open Basic Problems**

- Semantic gap between visual features and categories
  - General vs. task-specific representations

# **Open Basic Problems**

- Semantic gap between visual features and categories
  - General vs. task-specific representations
  - An observation: recent research mostly task-specific

# **Open Basic Problems**

- Semantic gap between visual features and categories
  - General vs. task-specific representations
  - An observation: recent research mostly task-specific

- What are successful general representations?
  - Grammars and logic are back!
  - SIG-II workshop at ICCVII
  - U Grenander, D Mumford, SC Zhu, A Yuille, L Davis, R Chellapa, N Ahuja, A Leonardis, P Felzenszwalb, S Todorovic ...

# Goal

A unified computational framework capable of: Discovery Detection Segmentation Summarization . . . of visual categories

in images and video

# What is an Object?

compositionality

Spatial arrangement of its parts



# What is an Object?

compositionality

Spatial arrangement of its parts

#### Parts = Objects in their own right

Part discovery = Suspicious coincidences

# What is an Object?

#### compositionality

### Spatial arrangement of its parts

#### and

#### context

Spatial and semantic constraints with other objects

![](_page_8_Picture_6.jpeg)

# What is an Activity?

#### compositionality

Spatiotemporal arrangement of its parts

and

#### context

Spatiotemporal constraints with other activities

![](_page_9_Picture_6.jpeg)

# What is an Activity?

#### compositionality

Spatiotemporal arrangement of its parts

and

#### context

Spatiotemporal constraints with other activities

![](_page_10_Picture_6.jpeg)

![](_page_10_Picture_7.jpeg)

AND-OR graph: a hierarchy of random fields

Brendel & Todorovic CVPRII, ICCVII

### Issues

![](_page_11_Figure_1.jpeg)

Grounding grammars on pre-selected low-level features

# **Example Low-Level Features**

![](_page_12_Picture_1.jpeg)

**STIPs** 

![](_page_12_Figure_3.jpeg)

![](_page_12_Figure_4.jpeg)

optical flow

#### frame 36 hand shaking frame 195 kicking frame 310 punching

video segmentation

![](_page_13_Figure_0.jpeg)

Grounding requires (probabilistic) repeatability of: features and their spatiotemporal placement

Satisfied only in particular settings/tasks (e.g., single-view object recognition)

# **Example Issues**

low-level —	→ mid-level
point-based features	regions, contours
stable, repeatable	unstable, poor repeatability
poor informativeness	informative

# Hypothesis

![](_page_15_Figure_1.jpeg)

# Grounding grammars on summaries of low-level features

# Hypothesis

![](_page_16_Figure_1.jpeg)

# Grounding grammars on summaries of low-level features where the summarization is guided top-down

# Objective

# Formalize a mid-level feature that will be: informative like regions (e.g., encode structure) and

repeatable like points (e.g., view-invariant)

# **Bags of Right Features/Detections**

![](_page_18_Picture_1.jpeg)

#### If the category occurs,

#### it has to be in the spotlight of many BORDs

so they can jointly support the occurrence hypothesis

Perina & Jojic CVPRII

# **Example Problem: Activity Recognition**

![](_page_19_Picture_1.jpeg)

# Given a video with noisy people detections

Amer & Todorovic ICCVII

# **Example Problem**

![](_page_20_Picture_1.jpeg)

Given a video with noisy people detections

Detect and localize: all activity instances & actors

Amer & Todorovic ICCVII

# **Example Problem**

![](_page_21_Picture_1.jpeg)

Grounding the grammar on a space-time grid of Bags of Right Detections (BORDs)

Structure is encoded in the overlap of BORDs

# **Bags of Right Detections**

![](_page_22_Picture_1.jpeg)

![](_page_22_Picture_2.jpeg)

iBORD

![](_page_22_Figure_4.jpeg)

# **Bags of Right Detections**

![](_page_23_Picture_1.jpeg)

 $\mathbf{X}$  - latent variable constrained by all BORDs

![](_page_23_Picture_3.jpeg)

# **Detection and Localization**

![](_page_24_Picture_1.jpeg)

#### BORDs jointly constrain the solution

Amer & Todorovic ICCVII

# **Example Problem: Object Recognition**

![](_page_25_Picture_1.jpeg)

# Given a set of edges in the image detect and localize all object instances and estimate their 3D pose

# **Bags of Right Detections**

![](_page_26_Figure_1.jpeg)

#### BORDs jointly constrain the solution

Payet & Todorovic ICCVII

# **Our Approach**

![](_page_27_Figure_1.jpeg)

Amer & Todorovic ICCVII

# **Our Approach**

Initial placement of BORDs on a regular grid Search for optimal features by warping the grid

![](_page_28_Picture_3.jpeg)

video frames

# MAP Inference:

 Warp the grid to expected locations
Select MAP BORDs

chains graphical model

# **The Chains Model**

![](_page_29_Figure_1.jpeg)

# **The Chains Model**

# $P(M, O, L_S, L_E, F)$

# $= P(M, O)P(L_S|M, O, F)P(L_E|M, O, F)$

# $\cdot \prod_{i} P(F_{O(i+1)}|F_{O(i)}) \prod_{i \in F-O} P_{bgd}(F_i)$

![](_page_31_Picture_1.jpeg)

 $P(M, O, L_S, L_E|F)$ 

# $P(L_S, L_E|F) \propto \sum_{M,O} P(M, O, L_S, L_E, F)$

 $P(L_S, L_E|F) \propto \sum_{N \in \mathcal{O}} P(M, O, L_S, L_E, F)$ M.O

 $= \sum_{M,O} P(M) \left| \prod_{i,j} P(F_j | F_i) \right| P(L_S, L_E | M, O, F)$ 

 $P(L_S, L_E|F) \propto \sum_{N \in \mathcal{O}} P(M, O, L_S, L_E, F)$ M,O

![](_page_34_Figure_2.jpeg)

# **MAP Inference = LP**

minimize  $\operatorname{tr} \{ \boldsymbol{C}_{\boldsymbol{b}}^{\mathrm{T}} \boldsymbol{X} \} + \alpha \| (\boldsymbol{I} - \boldsymbol{W}) \boldsymbol{X} \boldsymbol{Q} \|_{1} + \beta \| (\boldsymbol{I} - \boldsymbol{X}) \boldsymbol{Q} \|_{1}$ subject to  $\boldsymbol{X} \ge 0, \ \boldsymbol{X} \boldsymbol{1}_{n} = \boldsymbol{1}_{n}, \ \boldsymbol{b} \ge 0, \ \| \boldsymbol{b} \|_{2}^{2} = 1.$ 

> Searching for optimal features under non-rigid shape deformations of the grid of BORDs

![](_page_35_Picture_3.jpeg)

#### video frames

![](_page_36_Picture_2.jpeg)

Correct detection and localization of: kicking and pushing and actors involved

#### video frames

![](_page_37_Picture_2.jpeg)

Correct detection and localization of: handshaking and hugging and actors involved

#### video frames

![](_page_38_Picture_2.jpeg)

Failure example correct: handshaking and hugging wrong: actors involved

# **Example Problem: Object Recognition**

#### MAP inference = LP

![](_page_39_Picture_2.jpeg)

Searching for optimal features under non-rigid shape deformations of the grid of BORDs

Payet & Todorovic ICCVII

![](_page_40_Picture_1.jpeg)

#### Correct detection, localization, and pose estimation

![](_page_41_Picture_1.jpeg)

#### Correct detection, localization, and pose estimation

# Conclusion

Prior work: pre-selected features, typically low-level for repeatability

- Proposed mid-level features:
  - Allow abstraction of low-level features
  - Reduce the semantic gap
  - Enable addressing multiple tasks
  - Repeatable, and jointly encode structure

![](_page_42_Picture_7.jpeg)

# Acknowledgment

NSF IIS 1018490