

# **Toward Coherent Object Recognition and Scene Layout** Understanding

### Silvio Savarese – University of Michigan at Ann Arbor



Byung Kim

Wongun Choi Min Sun



Ying-ze Bao

Ryan Tokola









Jingen Liu



28 Dicembre 2010

Yu Xiang

Anush Mohan Shyam Kumar

**Fisher Yu** 



## **Geometrically "rich" scene understanding**







Object categorization





Scene classification



Pattern recognition





**Object categorization** 





#### Scene classification



#### Pattern recognition

0	1	2	3	4
0	1	2	3	J
0	1	X	3	ч

# Joint semantic interpretation and geometric reconstruction of the scene

#### Camera localization



### 30



### 3D object & scene modeling



### 3D geometry



**Object categorization** 





Scene classification



Pattern recognition

0	1	2	3	4
0	1	2	5	ч
0	1	X	3	4
-			-	

# Joint semantic interpretation and geometric reconstruction of the scene



- •Hoiem et al. 06-10 •Gould et al. 09 •Hedau et al. 09
- Gupta et al, 10
- Bao, Sun, Savarese 10
- Sun, Bao, Savarese 10

### Camera localization







3D object & scene modeling

3D geometry



#### **Object categorization**





Scene classification



#### Pattern recognition

0	1	Z	3	4
0	1	2	5	ч
0	1	λ	3	ч

# Joint semantic interpretation and geometric reconstruction of the scene



- Ess et al, 2009
- Pellegrini et al , 2010

- Choi, Shahid, Savarese, 2009
- Choi & Savarese, 2010



#### 3D object & scene modeling





3D geometry

## How can we achieve all of this?



**Intuition:** Objects in the 3D physical space show consistent geometrical properties within the same view and across views



Monitor x, y location in the image Bounding box/ scale Azimuth: 5 degrees Zenith: 15 degrees 3D shape

mouse

**Theorem**: supporting plane orientation, camera pose and focal length can be estimated from zenith pose of at least 3 (non-collinear) objects

keyboard

**Intuition:** Objects in the 3D physical space show consistent geometrical properties within the same view and across views

# **Our objectives**

- Multi-view models for object categorization and 3D attribute estimation
- Coherent object detection and scene layout estimation



- Detect objects under generic view points
- Estimate object pose
- General and work for any object category



- Detect objects under generic view points
- Estimate object pose
- General and work for any object category

# **Current paradigm**



- No information is shared
- •No sense of correspondences of parts under 3D transformations
- Non scalable to large number of categories/view-points

# A new recent paradigm

- •Thomas et al. '06
- Kushal, et al., '07
- Savarese et al, 07, 08
- Chiu et al. '07 Hoiem, et al., '07

• Yan, et al. '07

• Liebelt et al., '08

• Xiao et al.,'08

• Sun et al 09

- I., '08 Liebelt et al., '08
  - Xiao et al.,'08
  - Arie-Nachimson & Basri, '09



Sparse set of interest points or parts of the objects are linked across views.

# A new recent paradigm



- Canonical parts captures view invariant diagnostic appearance information
- 2d ½ structure linking parts via weak geometry
- Parts and relationship are modeled in a probabilistic fashion
  - Parameters are learnt so as to maximize detection accuracy

# **Key contributions**

### Representation:

### Leadeing: representation on the viewing sphere:

- Model object appearance and shape from any position on the viewing sphere
- Enable view synthesis from novel view points
- Multi-view generative part-based model
  - Object is represented by collections of parts
  - Parts are linked across views
  - Parts and relationships are probabilistic

### Semi-supervised learning

- No part or pose labels are required
- Incremental:
  - Training images can be provided sequentially

## **Dense representation on view-sphere**



- Triangle T
- Parameter S

## Multi-view generative part-based model





## Multi-view generative part-based model





# $P(X, Y, T, S, R, \pi) \propto P(\pi | \alpha_T)$

 $\prod_{n} \{ P(X_n | \theta_{TR_n}(S), A) P(Y_n | \eta_{TR_n}(S)) P(R_n | \pi) \}$ 

### Exact Inference is intractable! We use Variational EM:

 $\ensuremath{\alpha}$  = Part Prop. Prior  $\pi \sim Dir(\ensuremath{\alpha})$   $R \sim Mult(\ensuremath{\pi})$   $Y_n \sim Mult(\ensuremath{\eta})$   $X_n \sim N(theta)$   $\ensuremath{\eta}$  = Part Appearance  $\ensuremath{ heta}$  = Part Location/shape Yn=Codeword

**Xn=Location** 

 $X_n \leftarrow A \cdot X$ Image

# **Key contributions**

### •Representation:

- Dense representation on the viewing sphere:
  - Model object appearance and shape from any position on the viewing sphere
    - Enable view synthesis

### • Multi-view generative part-based model [Sun et al cvpr 09]

- Object is represented by collections of parts
- Parts are linked across views
- Parts and relationships are probabilistic

### •Learning:

- Semi-supervised learning
  - no part or pose labels are required
- Incremental:
  - Training images can be provided sequentially



## Incorporating geometrical constraints



- Parts are linked across views
- Part topology is preserved under view point transformations

# Incorporating geometrical constraints



## **Incremental learning**



- Enable unorganized and on-line collection training images
- Increase efficiency in learning (no need large storage space)

## **Incremental learning**



- Assign new training image to a triangle of the view sphere
- Evidence of training image is used to update model parameters
- Re-estimate sufficient statistics in a iterative fashion

# **Evolution of learnt parts**



## Examples of learnt part-based models



## Examples of learnt part-based models



Bicycle

## Examples of learnt part-based models



# Travel iron

# **Detection and pose estimation**

- **Detect** objects from any viewing angles
- Accurate pose estimation
- Synthesize (generate) object shape and appearance from novel views

- PASCAL 2006 dataset
- 3D Object Dataset

# Bicycle



# **Travel Iron**



# Car


## **Detection -** Pascal 2006 dataset



#### **3D Object Dataset** [Savarese et al 07]



Poses





8 azimuth angles
3 zenith
3 distances

#### ~ 7000 images!



## Detection

#### **3D object dataset**

[Savarese et al 07]



- Sun et al, ICCV 2009
- Su et al, CVPR
- Savarese et al, ICCV '07

## **Viewpoint Classification**



# **Viewpoint Classification**

#### **3D object dataset**

[Savarese et al 07]



#### Failure example

Category: car Azimuth = 225<sup>o</sup> Zenith = 30<sup>o</sup>





# Predicting object appearance from novel views



## Predicting object appearance from novel views

[For natural scenes, see Hoiem et al 07; Saxena et al 07] Thomas et al 08 Cremer et al 09









## Predicting object appearance from novel views



# Outline

- Multi-view models for object categorization and 3D attribute estimation
  - 3D Pose
  - 3D shape
- Coherent object detection and scene layout estimation

#### Depth-Encoded Hough Voting for coherent object detection and shape recovery

M. Sun, B. Xu, G. Bradski, S. Savarese, ECCV 2010



#### A new class of detectors:

input:

- single image;
- rough 3D location (optional)

#### Output:

- localize object;
- 3D shape



LH



















# Outline

- Multi-view models for object categorization and 3D attribute estimation
  - 3D Pose
  - 3D shape
- Coherent object detection and scene layout estimation

# Toward coherent object detection and scene layout understanding

#### Y. Bao, M. Sun, S. Savarese, CVPR 2010





Min Sun

Ying-ze Bao

•Recognizing objects can help estimate the layout

•Estimating the layout can help reinforce the existence of the objects

## **Our Contributions**

- Only use object detections to estimate 3D surface
  - No scene appearance model needed

•Hoiem et al. 06-10 •Gould et al. 09

- -Estimate focal length from single image
- -General camera pose model
  - camera tilt, camera in-planar rotation
- Hoiem et al. 06-10Gould et al. 09Hedau et al. 09
- No assumptions on camera height





**1. Layout:** most likely supporting surface 3D locations & orientations • camera pose & focal length • 3D objects location



2. Improve detection: remove false alarms, discover missed detections





## Joint inference process





**Theorem**: supporting plane orientation, camera pose and focal length can be estimated from zenith pose of at least 3 (non-collinear) objects

### **3D Encoded Detector**



Modified from Sun et. al. ECCV'10











#### Labelme dataset









# Outline

- Multi-view models for object categorization and 3D attribute estimation
- Coherent object detection and scene layout estimation
  - From single image
  - From videos



## Joint multi-target tracking and camera motion estimation from videos W. Choi & K. Shahid & S. Savarese WMC 2010

W. Choi & K. Shahid & S. Savarese WMC 201 W. Choi & S. Savarese , ECCV 2010

Wongun Choi

- Monocular cameras
- Un-calibrated cameras
- Arbitrary motion
- Highly cluttered scenes
   Occlusion
  - Background clutter





#### Joint camera and object track estimation

•Choi & Savarese, ECCV 2010



 $P(\Omega_t | \chi^t) \propto P(\Omega_t, \chi_t | \chi^{t-1}) = P(\chi_t | \Omega_t) \int P(\Omega_t | \Omega_{t-1}) P(\Omega_{t-1} | \chi^{t-1}) d\Omega_{t-1}$ 

 $\Omega$  : set of state variables X : set of observationse

## Interaction between Targets

- Interaction is modeled as a pair-wise MRF.
- Two exclusive models
  - Group motion vs. repulsion
  - Hidden variable to find the "mode" (hypothesis for the underlying interaction)








#### Object detection and tracking



**Object categorization** 



# Modeling human-human & human-object interaction



- Gupta et al 2009
- Yao and Fei-Fei 2010
- Desai et al 2010
- Lan et al 2010
- Choi et al , 2009
- Choi & Savarese, 2010

#### 3D object & scene modeling





Object semantic



#### Object detection and tracking



**Object categorization** 



**Object semantic** 

# Modeling human-human & human-object interaction



- Choi et al , 2009
- Choi & Savarese, 2010



3D object & scene modeling



Image: Second Secon

#### Learning spatial-temporal relationship among humans W. Choi & K. Shahid & S. Savarese , CVPR 2011, under review W. Choi & S. Savarese , under preparation, 2011





### Learning spatial-temporal relationship among humans

W. Choi & K. Shahid & S. Savarese , CVPR 2011, under review W. Choi & S. Savarese, under preparation, 2011



#### Learning spatial-temporal relationship among humans W. Choi & K. Shahid & S. Savarese WMC 2009 W. Choi & K. Shahid & S. Savarese , CVPR 2011, under review

Crossing – Talking – Queuing – Dancing – jogging



#### Learning spatial-temporal relationship among humans W. Choi & K. Shahid & S. Savarese WMC 2009 W. Choi & K. Shahid & S. Savarese , CVPR 2011, under review

Crossing – Talking – Queuing – Dancing – jogging



#### Learning spatial-temporal relationship among humans W. Choi & K. Shahid & S. Savarese WMC 2009 W. Choi & K. Shahid & S. Savarese , CVPR 2011

W. Choi & K. Shahid & S. Savarese WMC 2009W. Choi & K. Shahid & S. Savarese , CVPR 2011, under reviewW. Choi & S. Savarese, under preparation, 2011



## Conclusions

- Joint recognition and recognition is critical
- Leverage multi-view models for object categorization and 3D attribute estimation.
  - Pose estimation from a single view
  - 3D shape from a single view
- Models for coherent object detection/ tracking and scene 3D layout estimation
  - From single images
  - From videos

## Thank you



FORD



GigaScale Systems Research Center (GSRC)







ARO ARMY

### **Probabilistic Formulation**

 $P(\Omega_t | \chi^t) \propto P(\Omega_t, \chi_t | \chi^{t-1}) = P(\chi_t | \Omega_t) \int P(\Omega_t | \Omega_{t-1}) P(\Omega_{t-1} | \chi^{t-1}) d\Omega_{t-1}$ 





## **Probabilistic Formulation**



## Interaction between Targets

- Interaction is modeled as a pair-wise MRF.
- Two exclusive models
  - Group motion vs. repulsion
  - Hidden variable to find the "mode". (hypothesis for the underlying interaction)



## **Posterior Probability**

 $P(\Omega_t | \chi^t) \propto P(\Omega_t, \chi_t | \chi^{t-1}) = P(\chi_t | \Omega_t) \int P(\Omega_t | \Omega_{t-1}) P(\Omega_{t-1} | \chi^{t-1}) d\Omega_{t-1}$   $P(\chi_t | \Omega_t) = P(X_t, Y_t | Z_t, \Theta_t) P(\tau_t | G_t, \Theta_t)$   $P(\Omega_t | \Omega_{t-1}) = P(Z_t | Z_{t-1}) P(\Theta_t | \Theta_{t-1}) P(G_t | G_{t-1})$ 

- Challenges
  - Nonlinear projection.
  - Pair-wise MRF for interaction. (non-gaussian motion model)
  - Hypothesis variables for detection and ground features
- Not able to apply analytical inference algorithm.