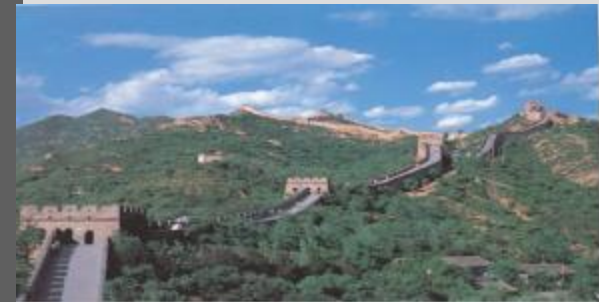


Methods for Representing and Recognizing 3D objects

part 2



1st Sino-USA Summer School in Vision,
Learning, and Pattern Recognition

VLPR 2009 • July 20-27, 2009 • Peking University, Beijing, China



Silvio Savarese

**University
of Michigan
at Ann Arbor**

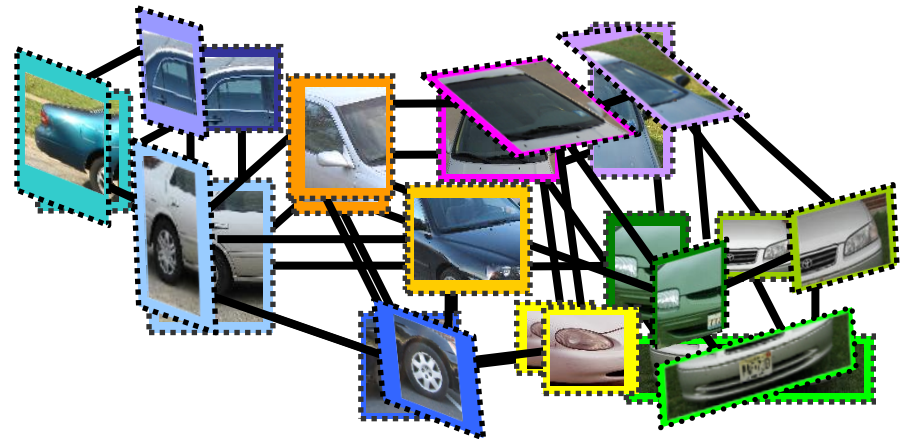
Part-based Multi-view models

Savarese, Fei-Fei, ICCV 07

Savarese, Fei-Fei, ECCV 08

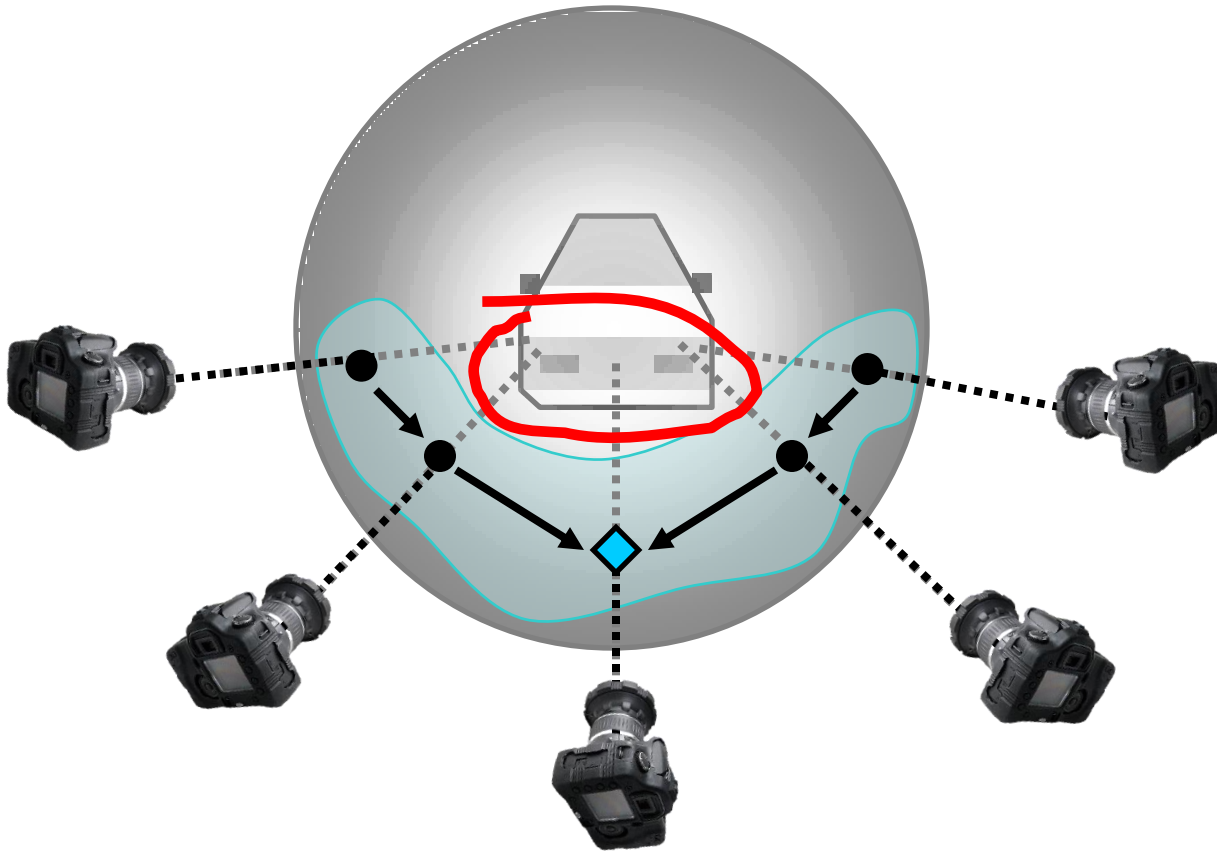
Sun, Su, Savarese, Fei-Fei, CVPR 09

Su, Sun, Fei-Fei, Savarese, ICCV 09



- Canonical parts captures diagnostic appearance information
- 2d $\frac{1}{2}$ structure linking parts via weak geometry

Canonical parts

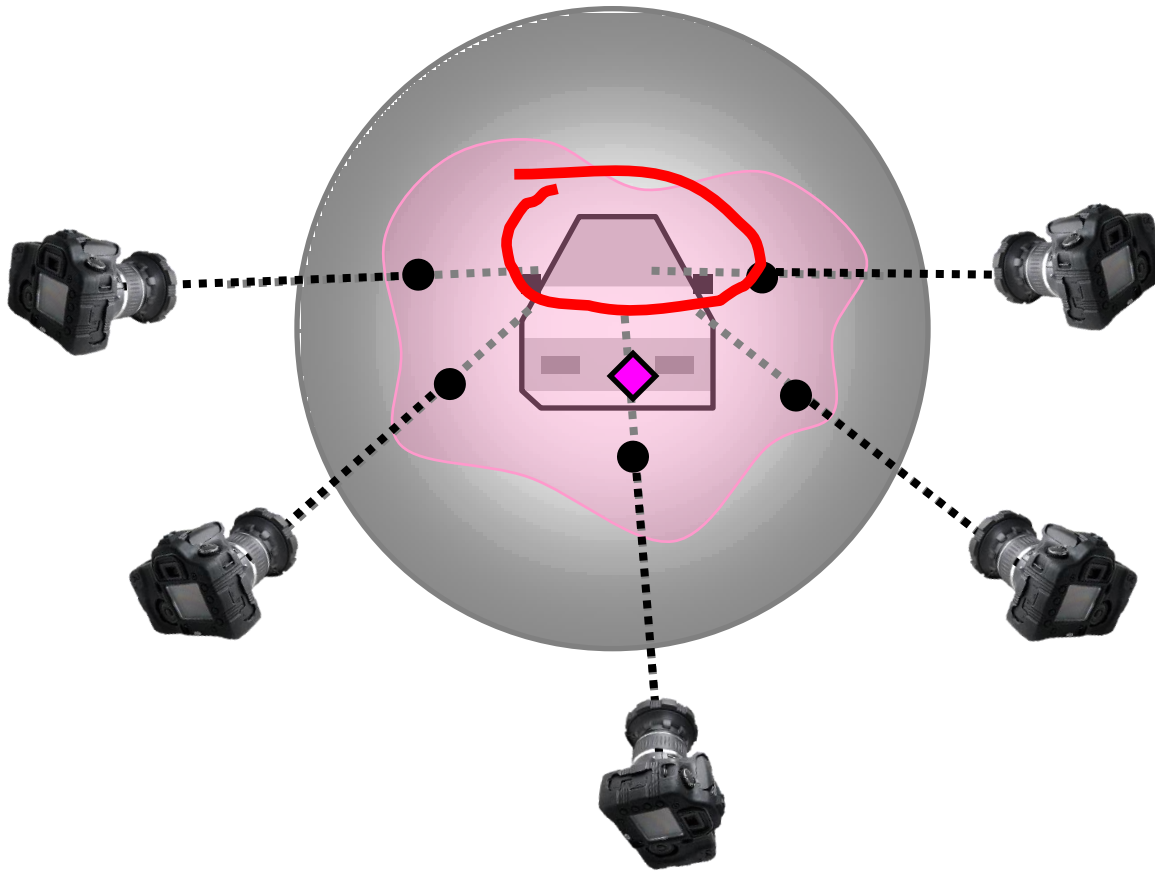


- If physical part is planar, canonical part is stable point on the manifold
- Canonical part can be computed from connected component of parts



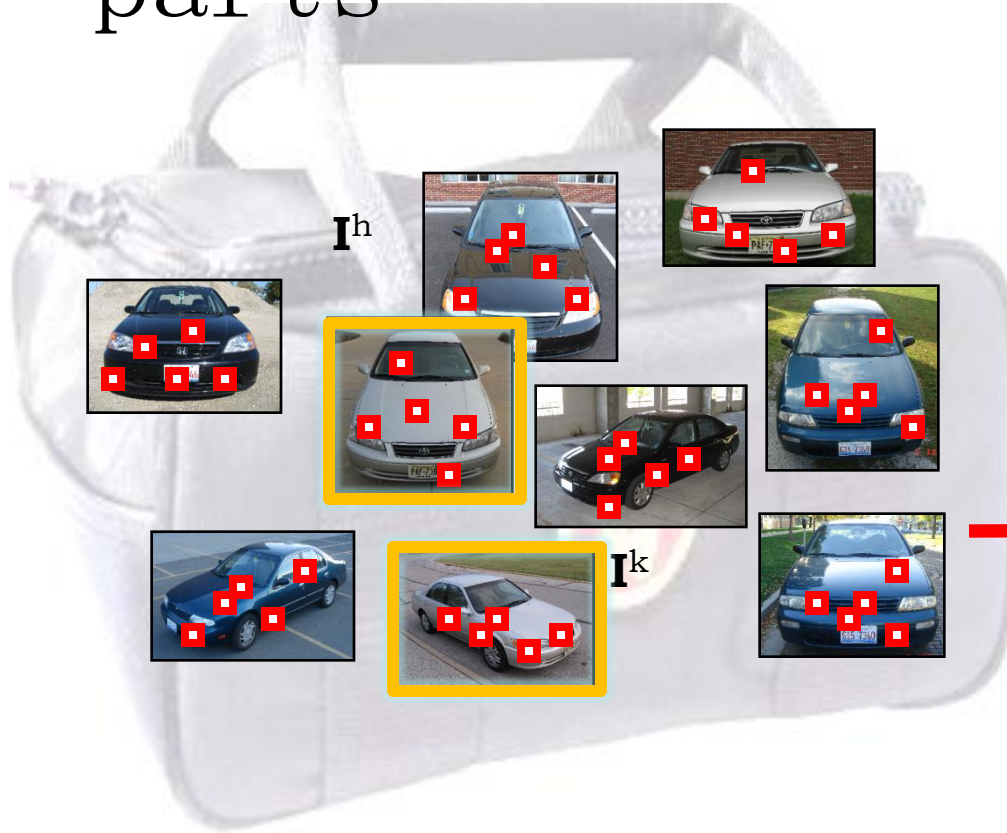
connected
component of parts

Canonical parts



connected
component of parts

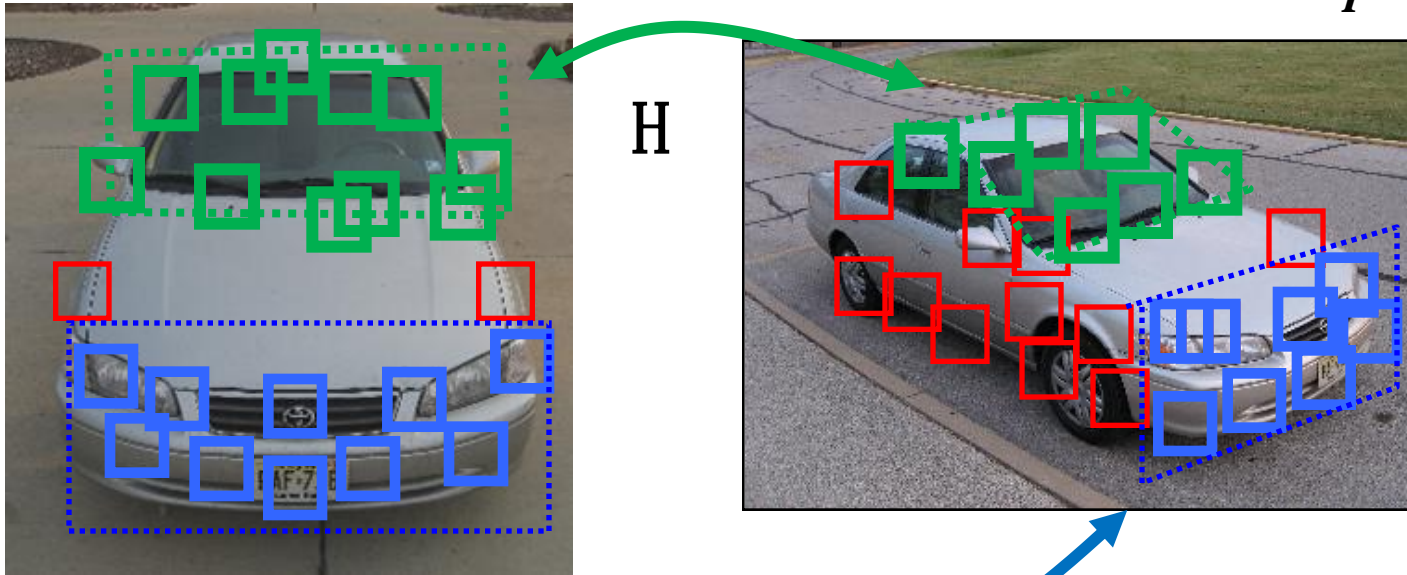
Connected components of parts



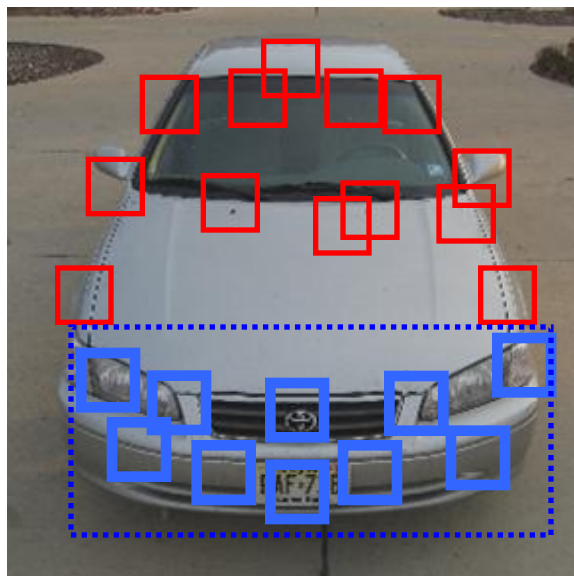
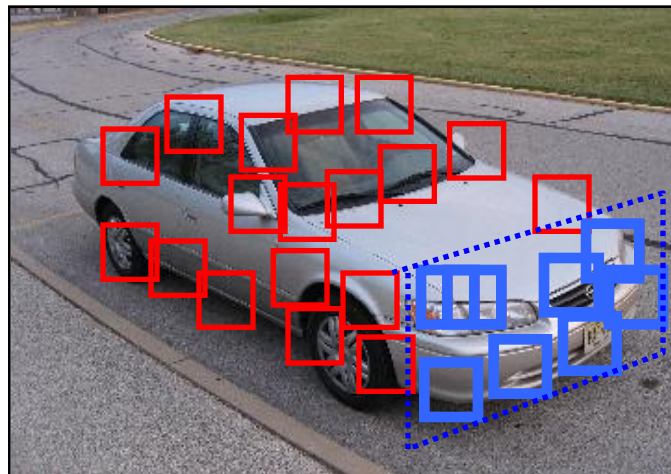
Unlabeled mix of images:

- category labels
- no pose labels;
- images of same instance from multiple views

$$\mathbf{I}^h = [x_1, x_2, \dots, x_M]$$

 \mathbf{I}^k  H

$$\mathbf{I}^h = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M]$$


 \mathbf{I}^k


$$\mathbf{x}_{2,1}, \mathbf{x}_{2,2}, \dots, \mathbf{x}_{2,n}$$

$$\pi : \mathbf{I}^h \rightarrow \{ \mathbf{P}_1^h, \mathbf{P}_2^h, \dots, \mathbf{P}_N^h, \mathbf{O}^h \}$$

$$\tau : \mathbf{I}^k \rightarrow \{ \mathbf{P}_1^k, \mathbf{P}_2^k, \dots, \mathbf{P}_N^k, \mathbf{O}^k \}$$

$$\min \left| \mathbf{O}^h \cup \mathbf{O}^k \right| + c N$$

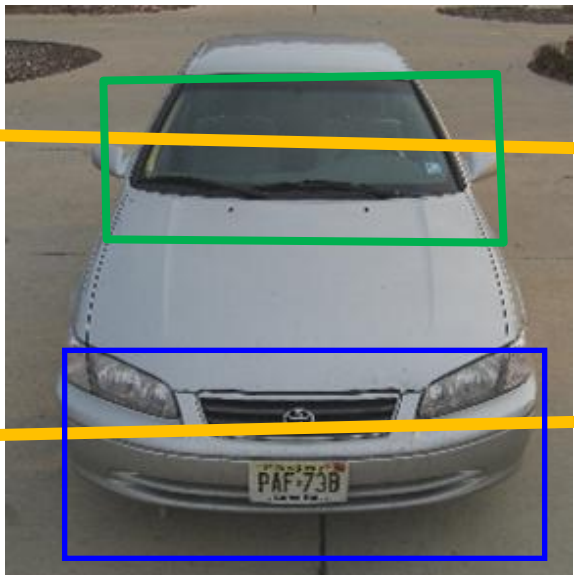
such that:

$$f(\mathbf{P}_i^h, \mathbf{P}_j^k, \beta_{i,j}) < \delta$$

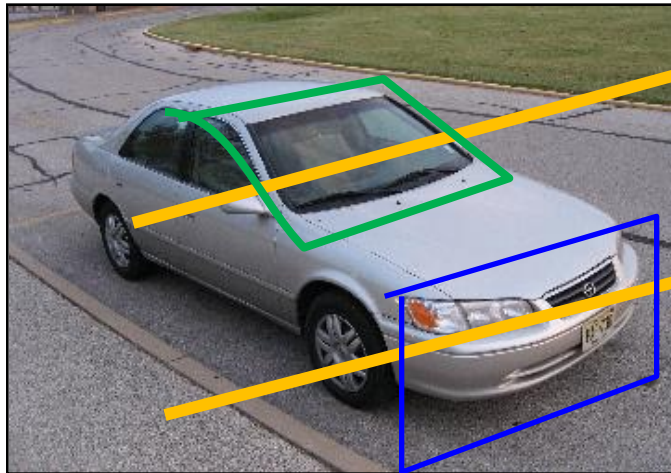
$$\forall i, j = 1 \dots N$$

$$f(\mathbf{P}_i^h, \mathbf{P}_j^k, \beta_{i,j}) = \left\| \mathbf{P}_i^h - \mathbf{H}_{i,j} \mathbf{P}_j^k \right\|$$

$$\mathbf{I}^h = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M]$$



$$\mathbf{I}^k$$



$$\pi : \mathbf{I}^h \rightarrow \{ \mathbf{P}_1^h, \mathbf{P}_2^h, \dots, \mathbf{P}_N^h, \mathbf{O}^h \}$$

$$\tau : \mathbf{I}^k \rightarrow \{ \mathbf{P}_1^k, \mathbf{P}_2^k, \dots, \mathbf{P}_N^k, \mathbf{O}^k \}$$

$$\min \left| \mathbf{O}^h \cup \mathbf{O}^k \right| + \mathbf{c} N$$

such that:

$$f(\mathbf{P}_i^h, \mathbf{P}_j^k, \beta_{i,j}) < \delta$$

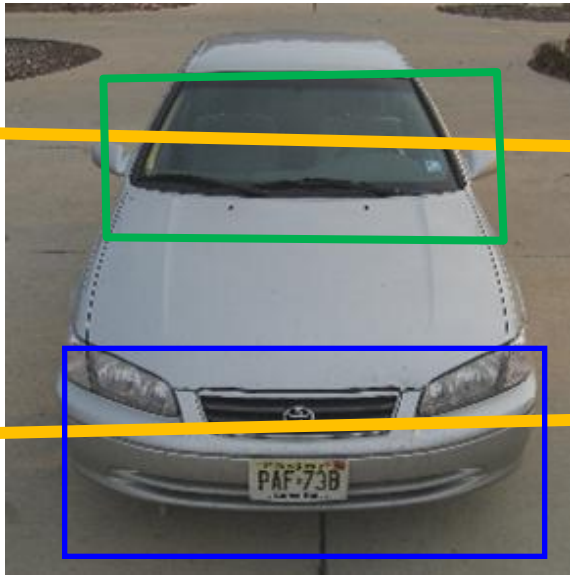
$$\forall i, j = 1 \dots N$$

$$f(\mathbf{P}_i^h, \mathbf{P}_j^k, \beta_{i,j}) = \left\| \mathbf{P}_i^h - \mathbf{H}_{i,j} \mathbf{P}_j^k \right\|$$

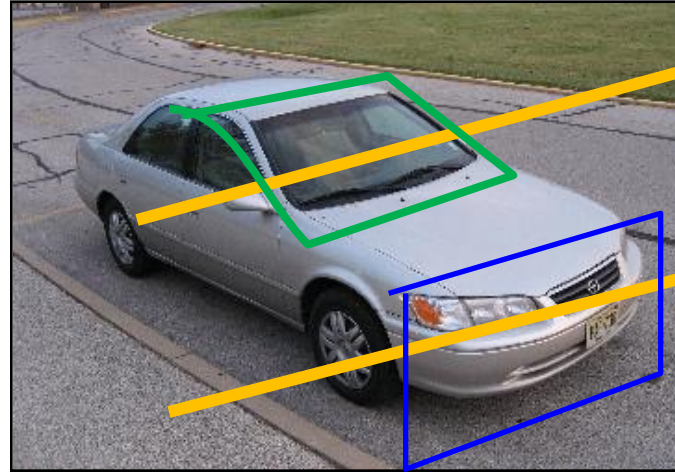
$$f(\mathbf{P}^h, \mathbf{P}^k, \gamma) < \delta$$

$$f(\mathbf{P}^h, \mathbf{P}^k, \gamma) = \left\| [\mathbf{P}_1^h, \dots, \mathbf{P}_N^h] \mathbf{F} [\mathbf{P}_1^k, \dots, \mathbf{P}_N^k]^T \right\|$$

$$\mathbf{I}^h = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M]$$



$$\mathbf{I}^k$$



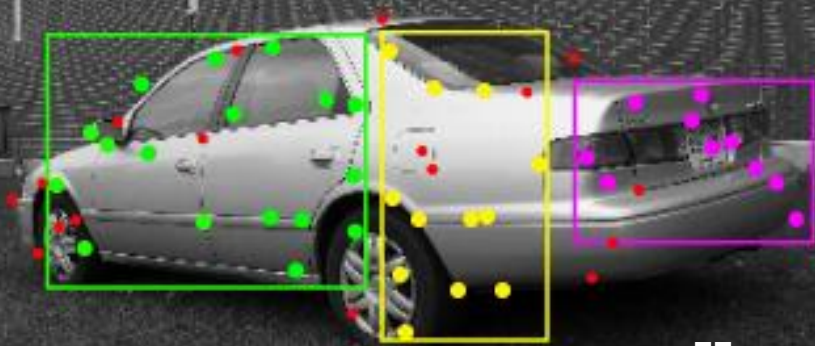
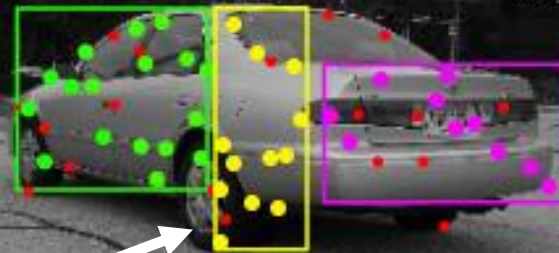
$$\pi : \mathbf{I}^h \rightarrow \{ \mathbf{P}_1^h, \mathbf{P}_2^h, \dots, \mathbf{P}_N^h, \mathbf{O}^h \}$$

$$\tau : \mathbf{I}^k \rightarrow \{ \mathbf{P}_1^k, \mathbf{P}_2^k, \dots, \mathbf{P}_N^k, \mathbf{O}^k \}$$

$$\min \left| \mathbf{O}^h \cup \mathbf{O}^k \right| + c N$$

GOAL:

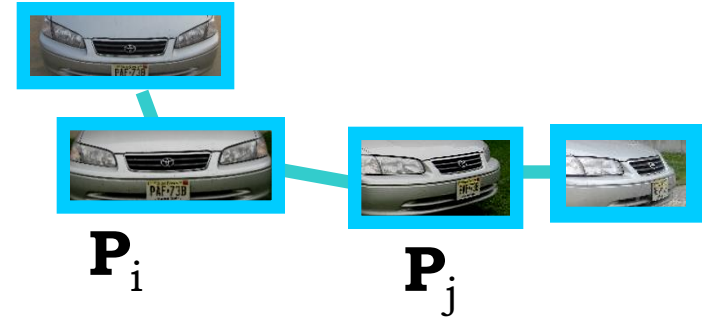
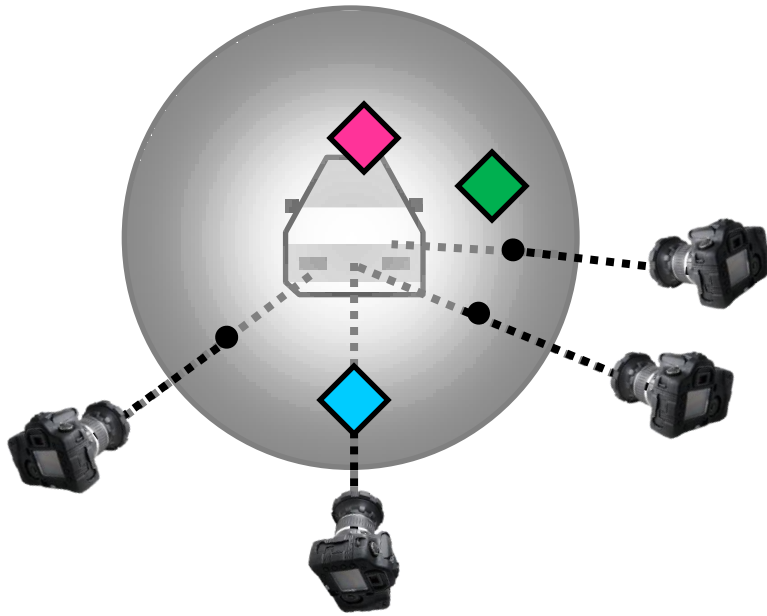
- discover partition while fitting multiple homographies
- fit global constraint
- minimize outlier set
- Use sequential RANSAC or RANSAC & J-linkage [toldo, fusiello eccv 08]

I^h  H I^k 

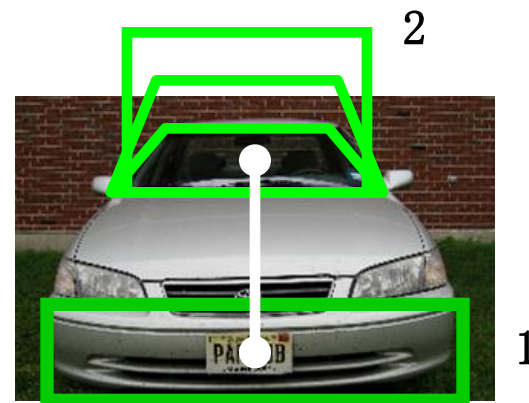
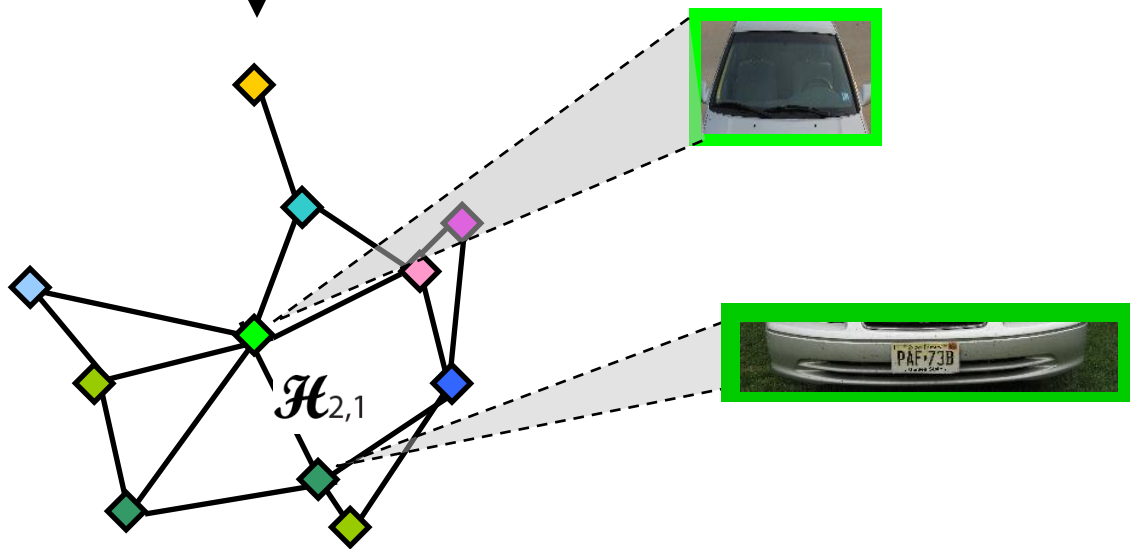
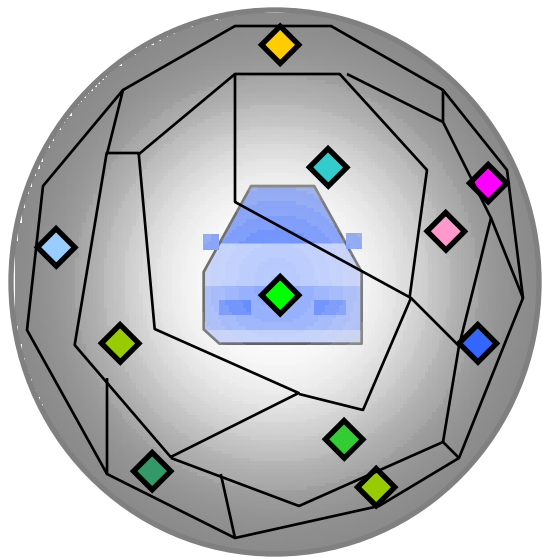
$$\pi : I^h \rightarrow \left\{ \mathbf{P}_1^h, \mathbf{P}_2^h, \mathbf{P}_3^h, \mathbf{O}^h \right\}$$

$$\tau : I^k \rightarrow \left\{ \mathbf{P}_1^k, \mathbf{P}_2^k, \mathbf{P}_3^k, \mathbf{O}^k \right\}$$

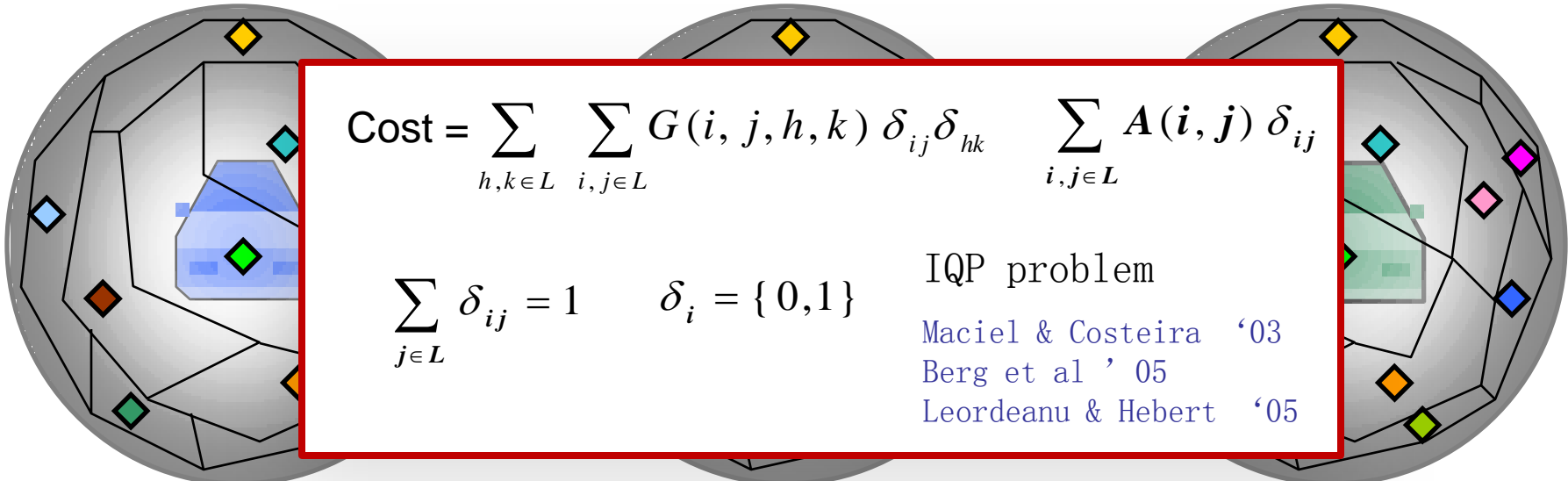
Canonical parts



Linkage structure



$$\mathcal{H}_{2,1} = \begin{pmatrix} \boxed{H_{2,1}} & \boxed{t_{2,1}} \\ 0 & 1 \end{pmatrix}$$


$$\text{Cost} = \sum_{h,k \in L} \sum_{i,j \in L} G(i,j,h,k) \delta_{ij} \delta_{hk} + \sum_{i,j \in L} A(i,j) \delta_{ij}$$

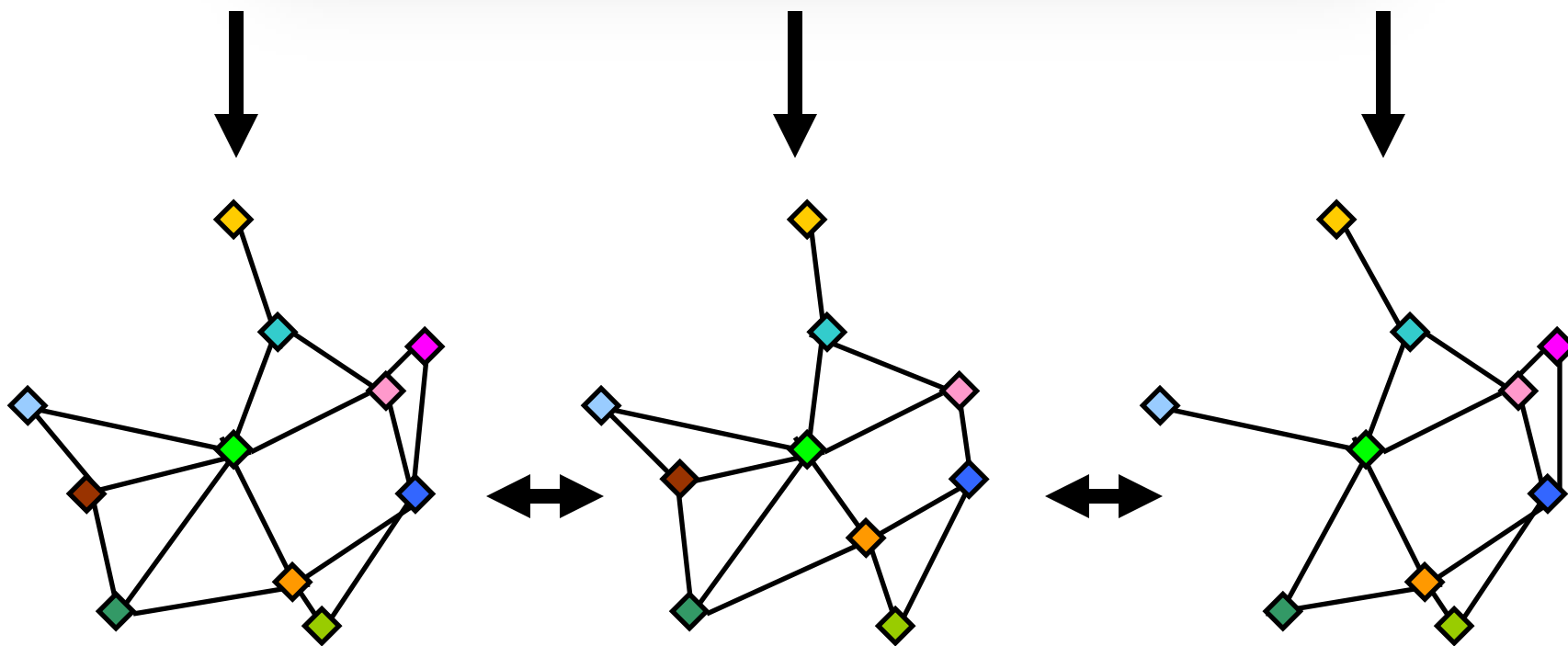
$$\sum_{j \in L} \delta_{ij} = 1 \quad \delta_i = \{0,1\}$$

IQP problem

Maciel & Costeira '03

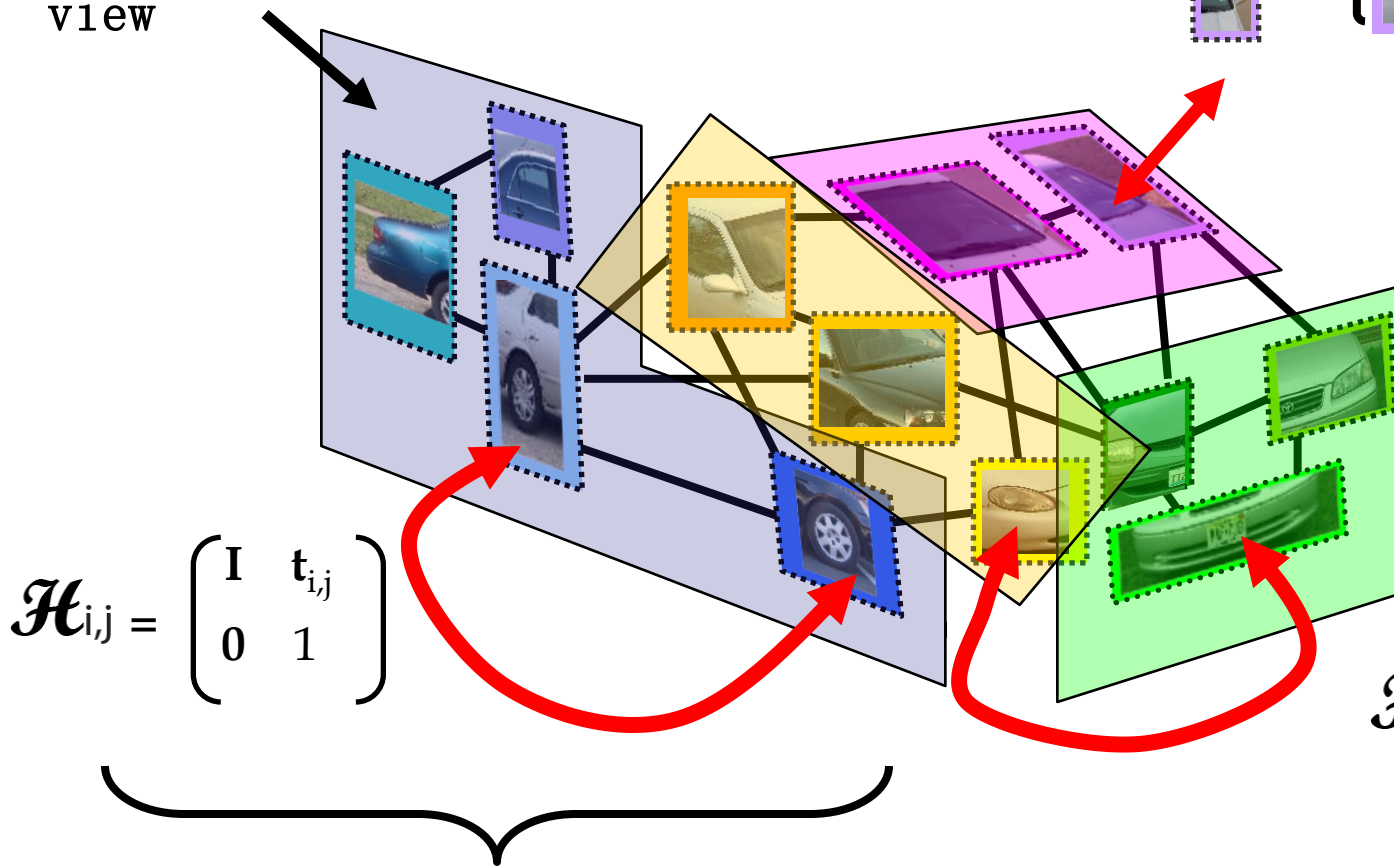
Berg et al '05

Leordeanu & Hebert '05



Category Model

Canonical view



Aspect Graphs

Koenderink & V. Doorn 76

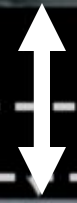
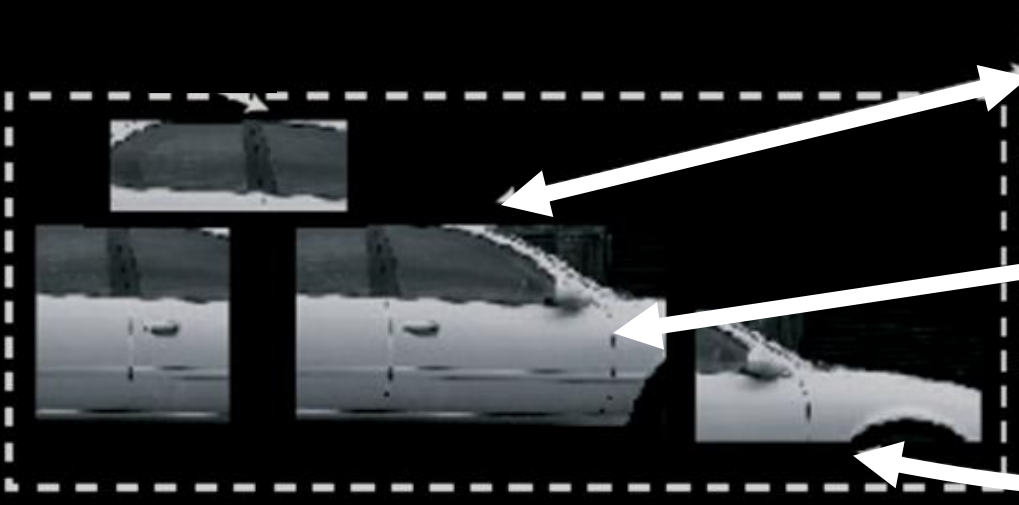
Bowyer & Dyer 90

Cyr & Kimia 04

$$\mathcal{H}_{i,j} = \begin{pmatrix} \mathbf{I} & \mathbf{t}_{i,j} \\ 0 & 1 \end{pmatrix}$$

$$\mathcal{H}_{i,j} = \begin{pmatrix} \mathbf{H}_{i,j} & \mathbf{t}_{i,j} \\ 0 & 1 \end{pmatrix}$$

2D single view model of the object



Object Recognition

Query image



model

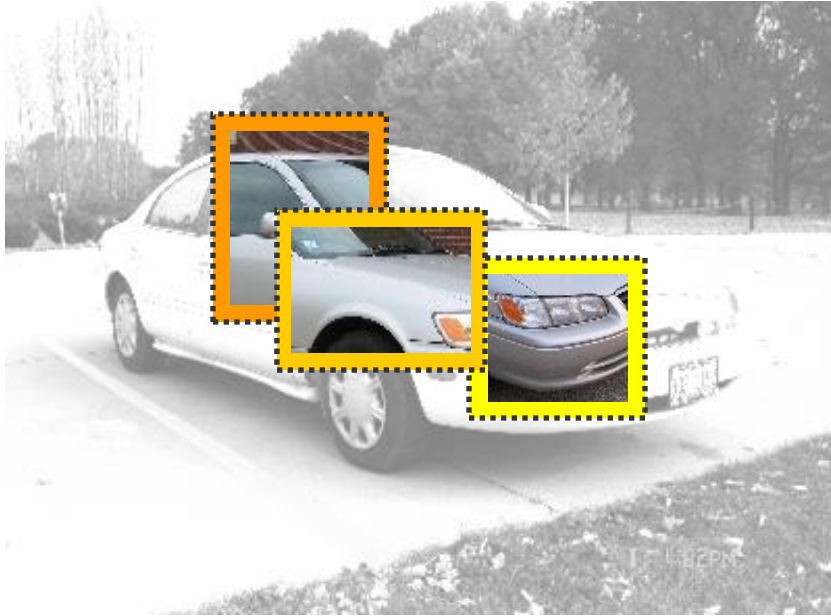


Algorithm

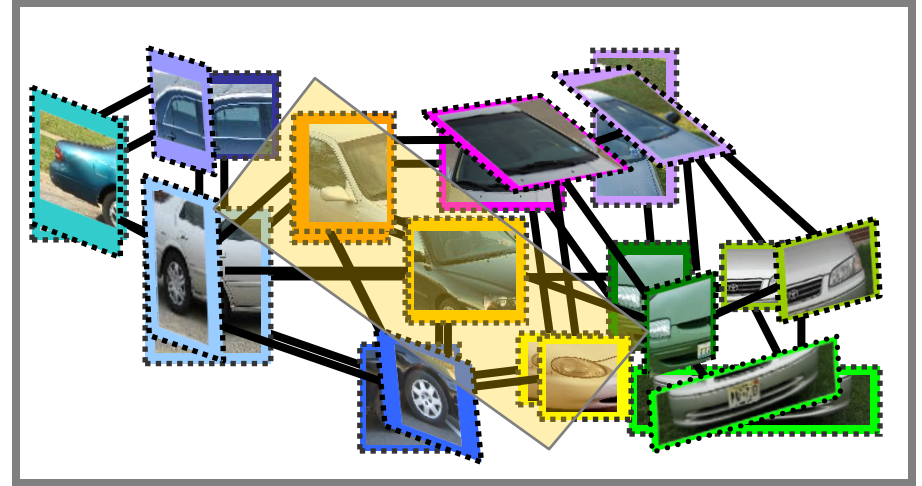
1. Find hypotheses of canonical parts consistent with a given pose

Object Recognition

Query image



model

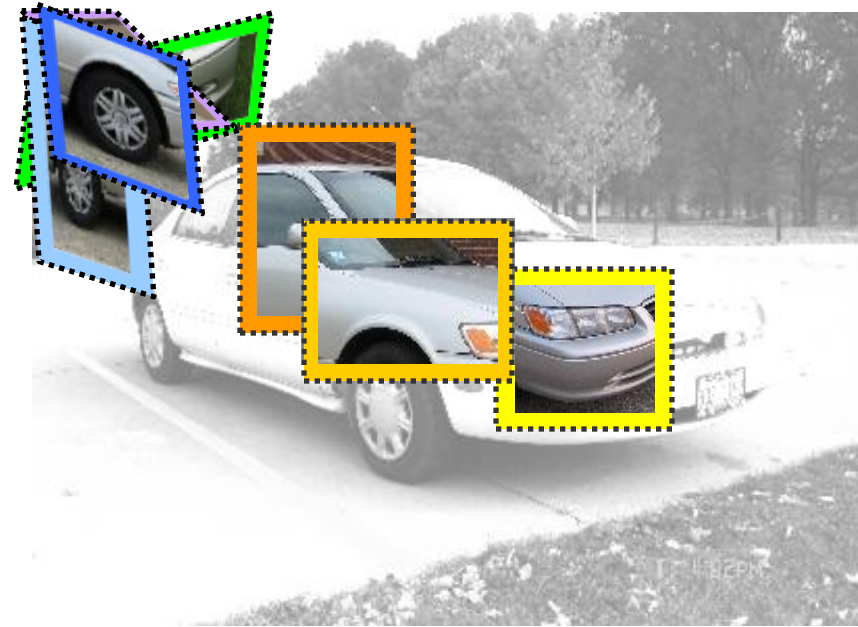


Algorithm

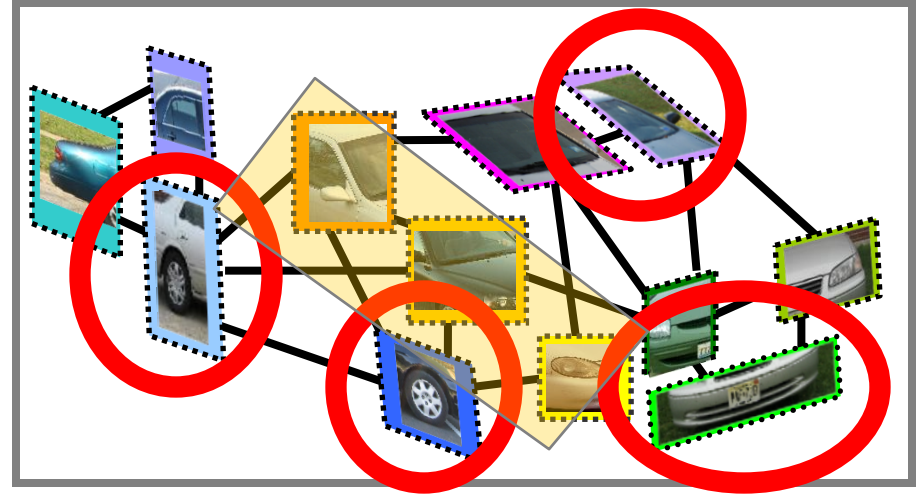
1. Find hypotheses of canonical parts consistent with a given pose
2. Infer position and pose of other canonical parts

Object Recognition

Query image



model

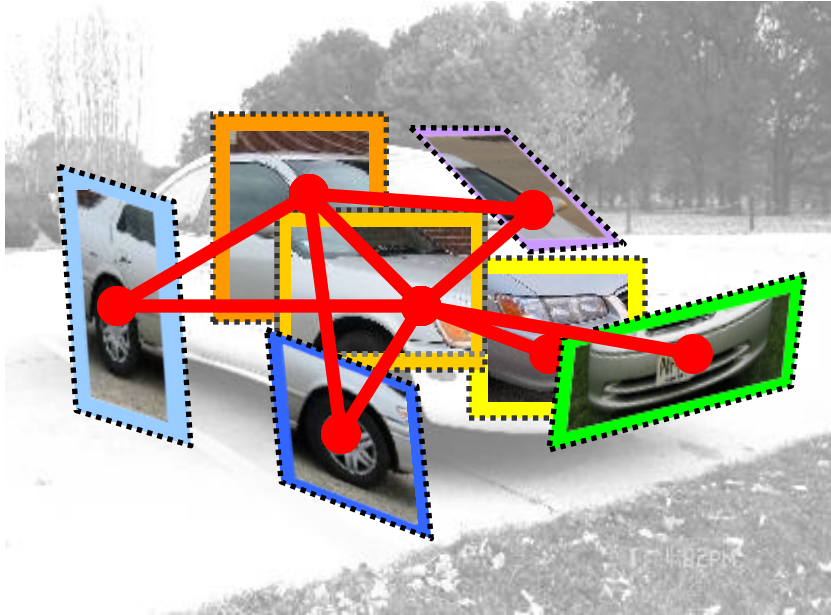


Algorithm

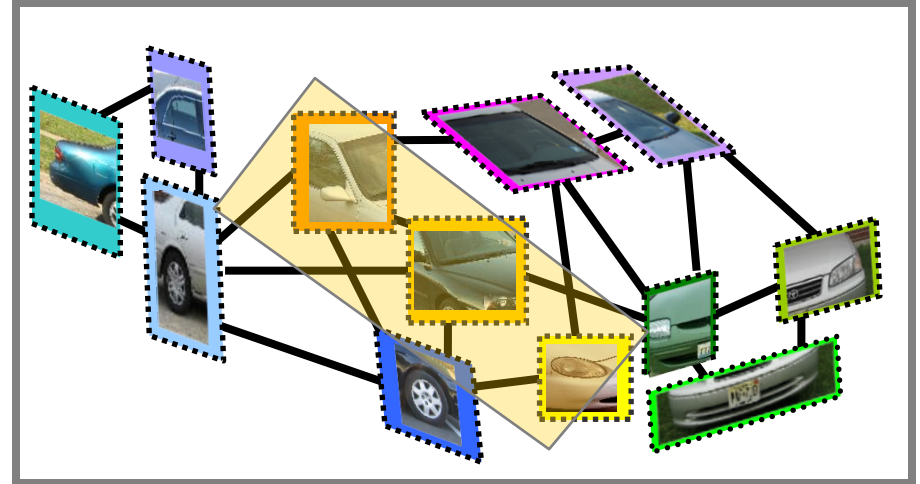
1. Find hypotheses of canonical parts consistent with a given pose
2. Infer position and pose of other canonical parts

Object Recognition

Query image



model



Algorithm

1. Find hypotheses of canonical parts consistent with a given pose
2. Infer position and pose of other canonical parts
3. Optimize over \mathbf{E} , \mathbf{G} and \mathbf{s} to find best combination of hypothesis
→ error

3D object class dataset

bicycle



car



cellphone



iron



shoe



stapler



toaster

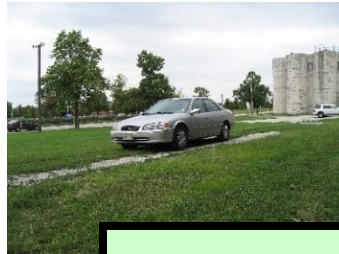


Poses

72

⋮

1



⋮



- 8 azimuth angles
- 3 zenith
- 3 distances

~ 7000 images!

1

2

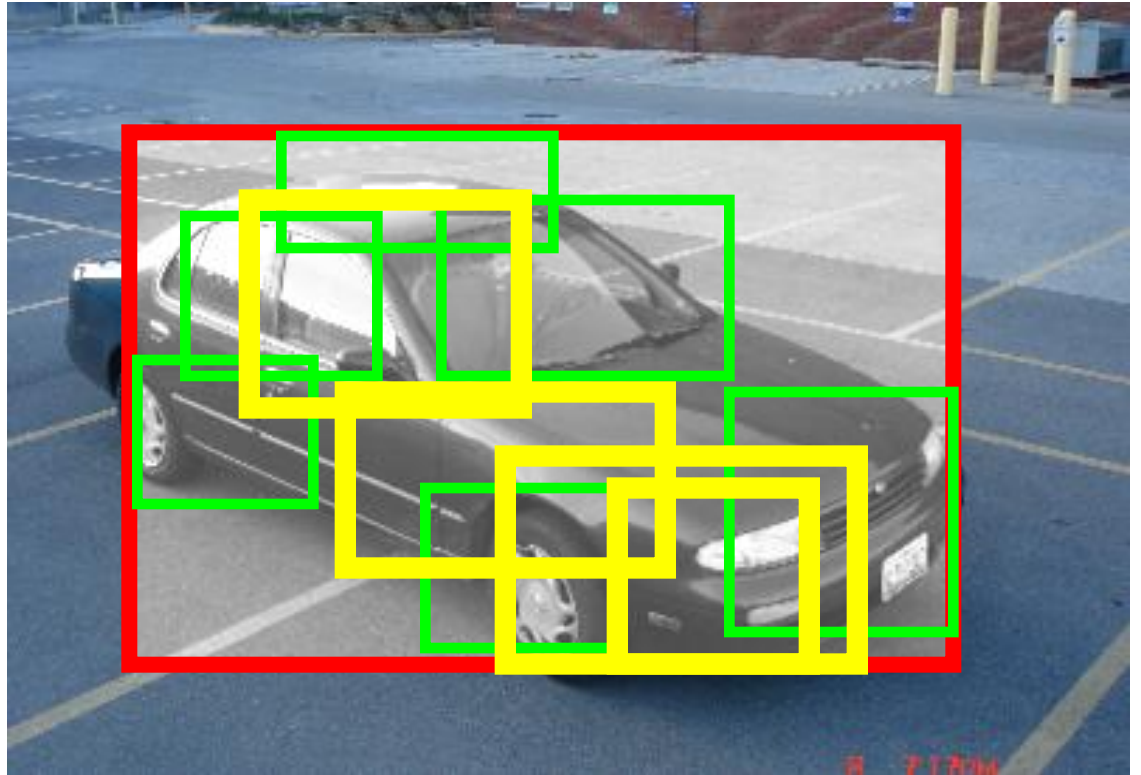
⋮

10

Instances

Examples

Category: car
Azimuth = 45°
Zenith = 30°
Distance = close



Examples

Category: mouse

Azimuth = 315°

Zenith = 0°

Distance = medium



Classification accuracy

Average Perf. = **75.7%**

cellphone	.76	.03	.03	.02	.10	.03	.03	
bike	.02	.81	.07	.02	.03	.02	.03	
iron			.77	.09	.06	.04	.04	
mouse	.04	.04		.87	.02	.02	.02	
shoe	.04	.06	.04		.62	.12	.12	
stapler		.11	.04	.04		.77	.04	
toaster	.08	.06	.03		.06		.75	
car	.04	.04		.12	.04	.07	.70	
	c	b	i	m	s	s	t	c

Average accuracy in classifying 1 out 8 categories
Random=12%

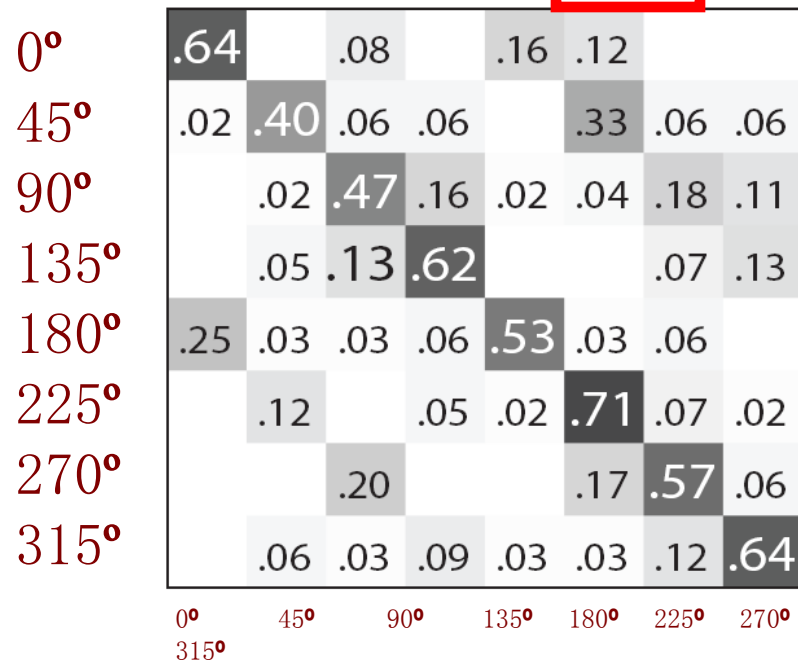
Failure example

Category: car
Azimuth = 225°
Zenith = 30°
Distance = close



Pose estimation accuracy

Pose estimation (1 catit.)
av. of 8 categ (57.2%)



Summary

	Single view	Mixture / Multi-view	Sav. et al, 07
View point invariant	X	✓ <small>Category</small>	✓
No supervision	✓	X <small>View point</small>	X → ✓ <small>category : all views all instances available</small>
# Categories	~300	2	8

- lack of a coherent methodology for learning parameters
- need of multi-view observations of the same object instance
- no generative model for robust learning and recognition

Part-based Multi-view Models

Savarese, Fei-Fei, ICCV 07

Savarese, Fei-Fei, ECCV 08

Sun, Su, Savarese, Fei-Fei, CVPR 09

Su, Sun, Fei-Fei, Savarese, ICCV 09



Min Sun

University of Michigan, USA



Hao Su

Beihang University, China



Fei-Fei Li.

Stanford U, USA

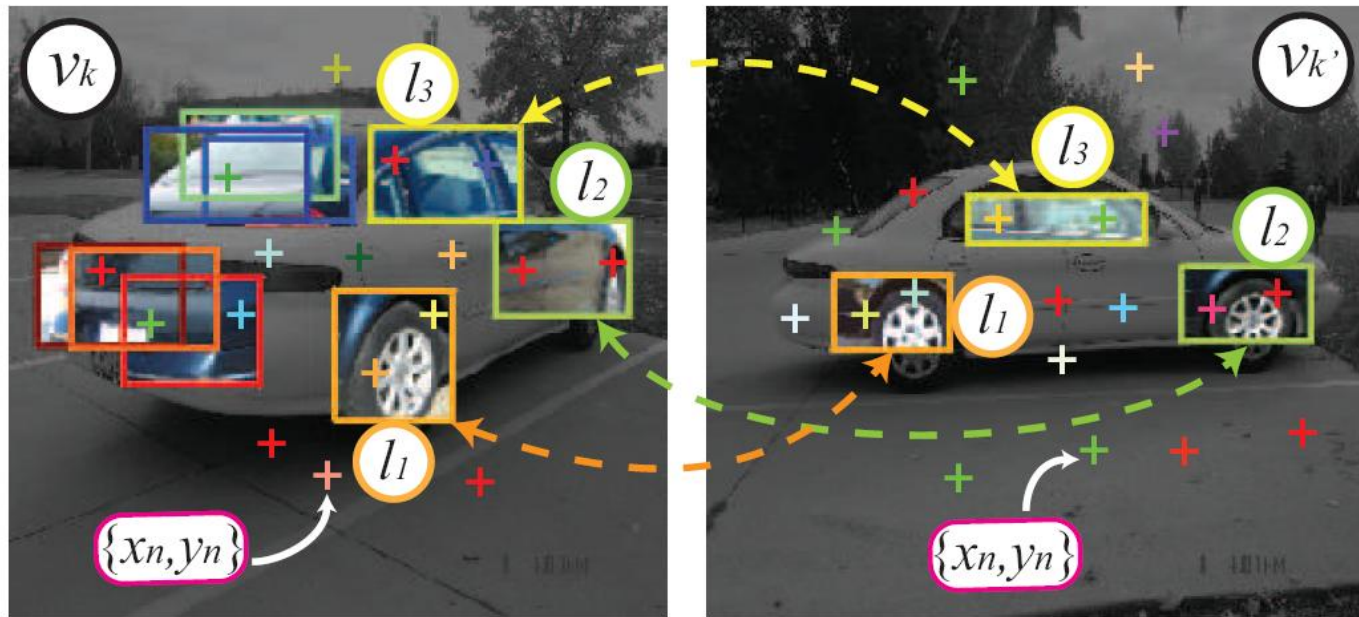
Part-based Multi-view Models

Savarese, Fei-Fei, ICCV 07

Savarese, Fei-Fei, ECCV 08

Sun, Su, Savarese, Fei-Fei, CVPR 09

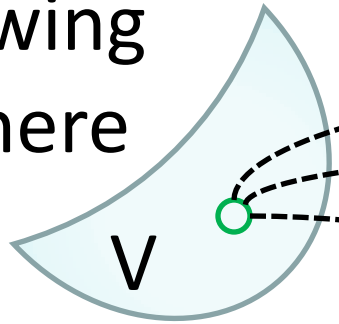
Su, Sun, Fei-Fei, Savarese, ICCV 09



- Probabilistic generative part-based model
- Dense multi-view representation on the viewing sphere

Part-based generative model

Viewing
Sphere



α = Part Proportion Prior

η = Part Appearance

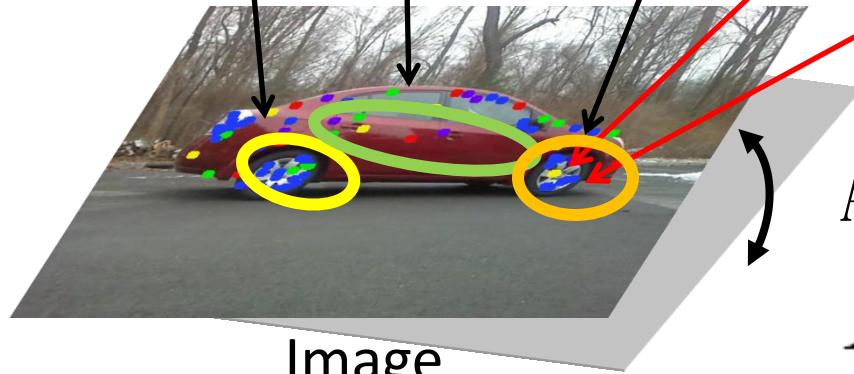
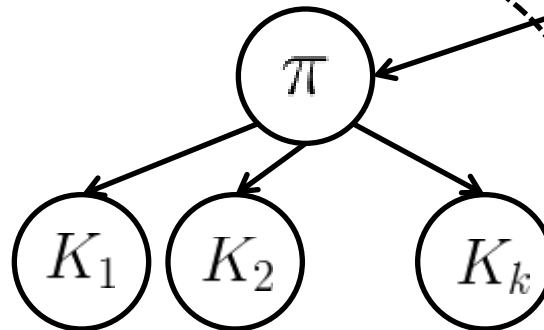
θ = Part Location/shape

$\pi \sim Dir(\alpha)$

$Y_n \sim Mult(\eta)$

$X_n \sim N(\theta)$

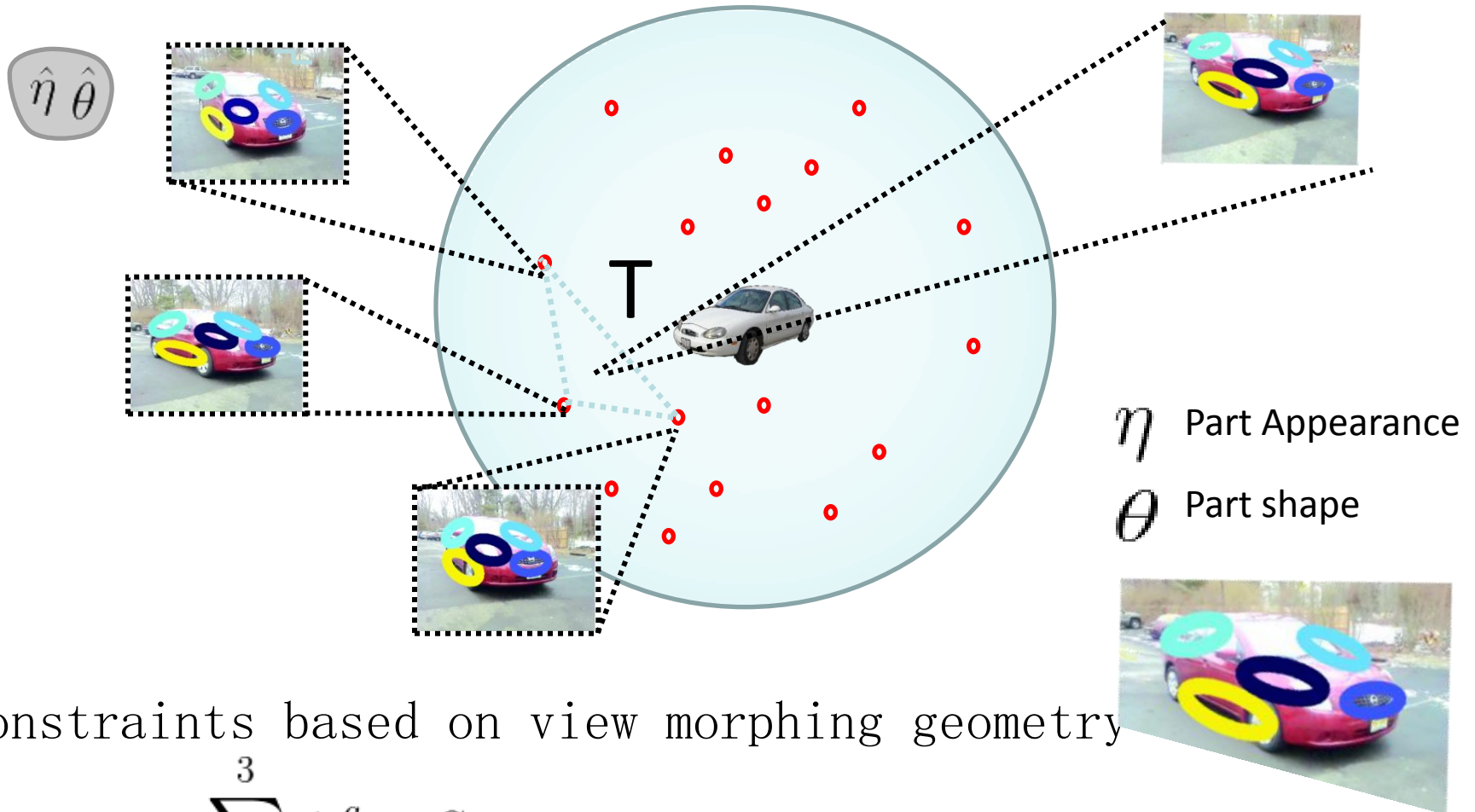
Y=Codeword, X=Location



A

$$X_n \leftarrow A \cdot X$$

Dense representation on view-sphere

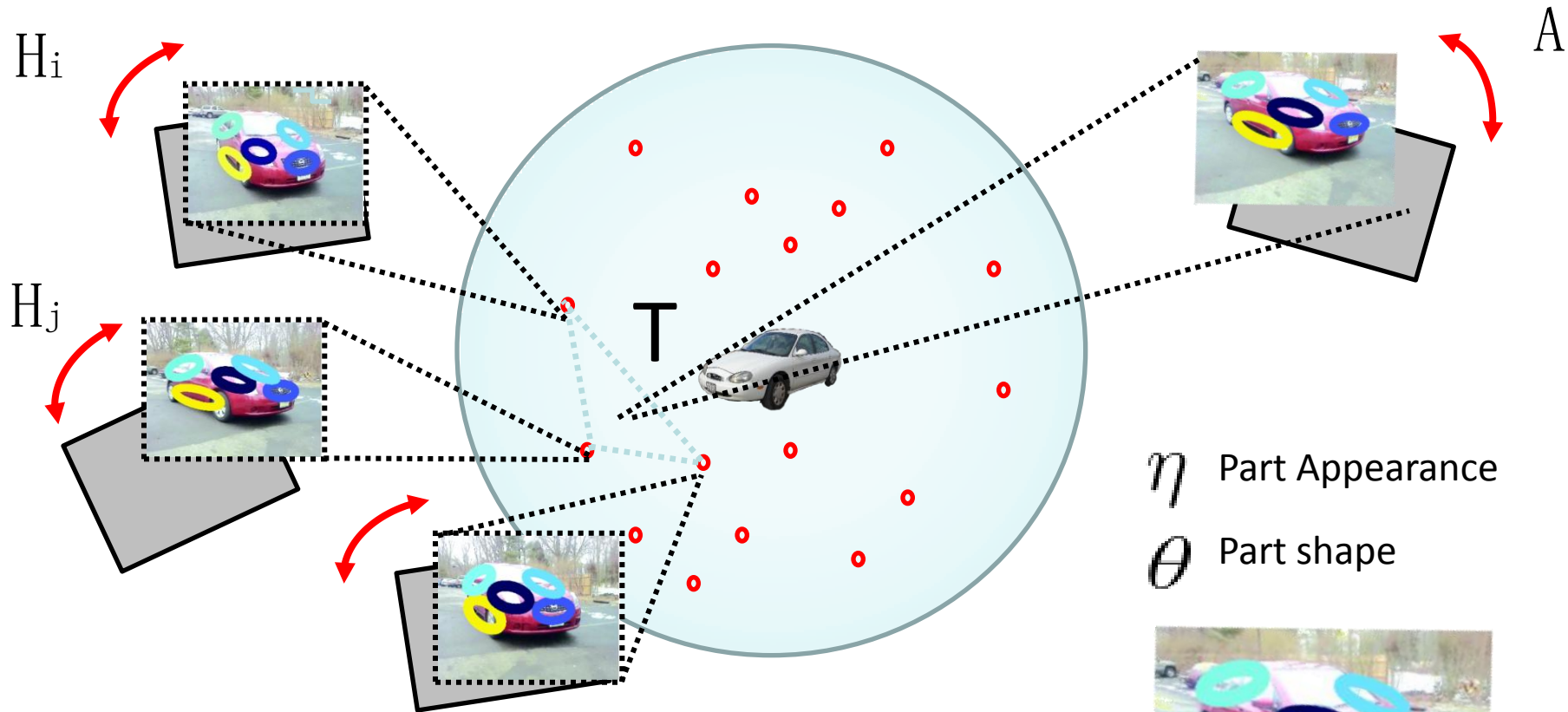


Constraints based on view morphing geometry

$$m = \sum_{g=1}^3 \hat{m}_{Tk}^g \cdot S_g$$

$$\theta = (m, \Sigma)$$

Dense representation on view-sphere



- Pre-warping transformations H
- Post-warping transformation A

η Part Appearance

θ Part shape



$$\theta = (m, \Sigma)$$

Dense representation on view-sphere

Joint probability of the model:

$$P(X, Y, T, S, K, \pi) = P(T|\phi)P(\pi|\alpha_T)P(S|\beta) \\ \prod_n^N \{P(x_n|\theta_{TK_n}(S), A)P(y_n|\eta_{TK_n}(S))P(K_n|\pi)\}$$

Observable variables: X, Y, T, S

Latent variables: K, π + relevant priors

Learning

Variational EM

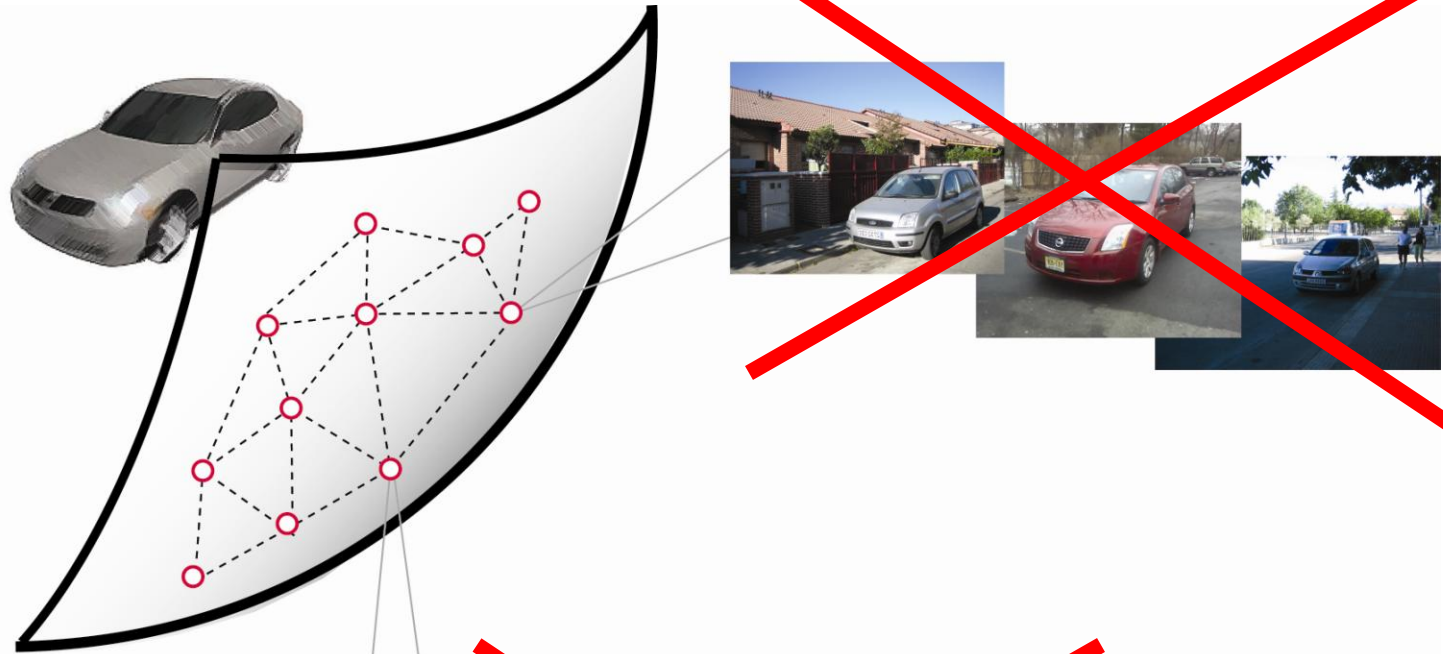
E step:

- update hidden part assignments π
part proportion K

M step:

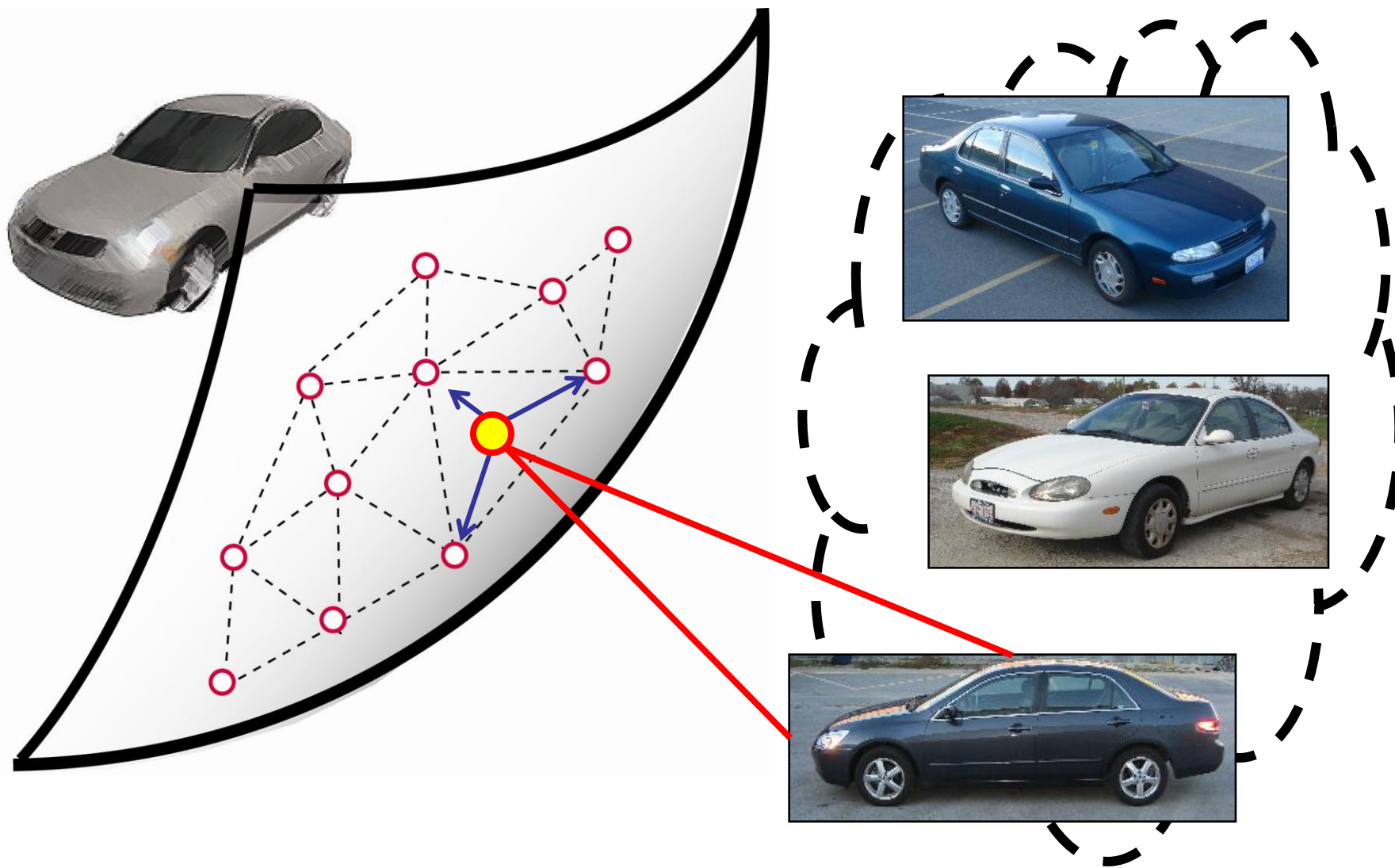
- update part proportion prior α
appearance η
Location/shape θ
- weakly supervised
- incremental (training image)

Weakly supervised

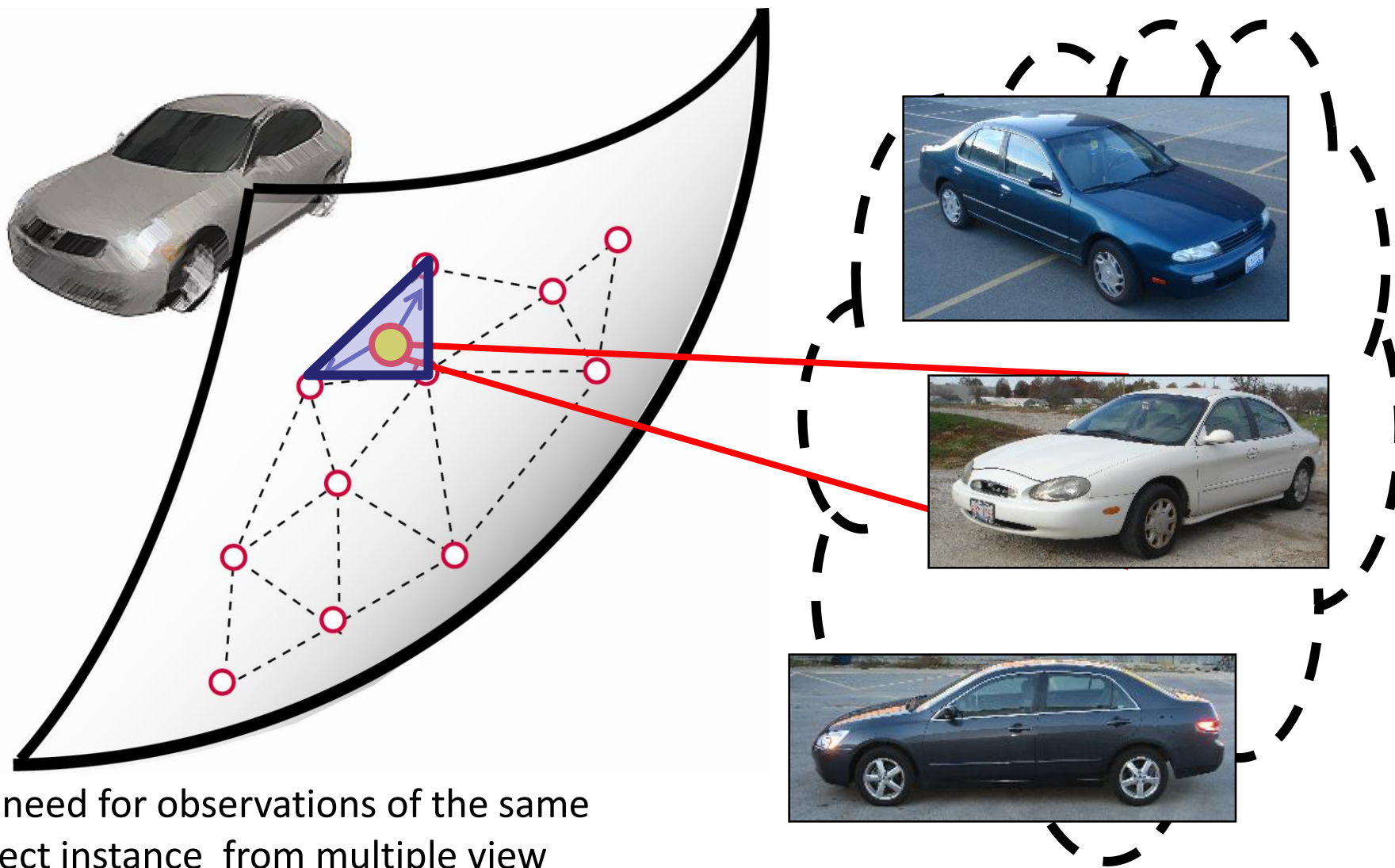


- Class label
- No pose label

Incremental learning

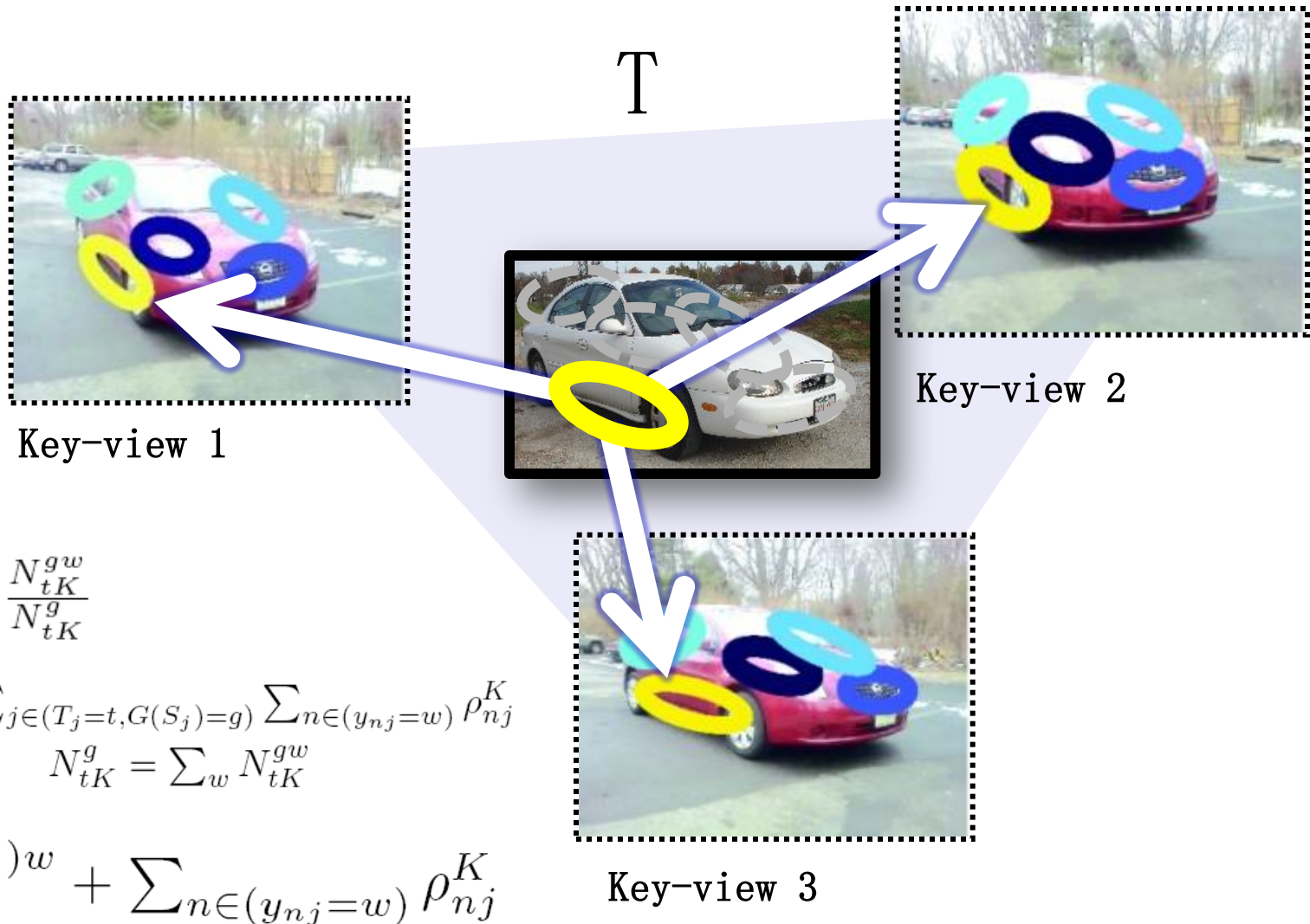


Incremental learning



No need for observations of the same object instance from multiple view

Incremental learning



$$\hat{\eta}_{tK}^{gw} = \frac{N_{tK}^{gw}}{N_{tK}^g}$$

$$N_{tK}^{gw} = \sum_{j \in (T_j=t, G(S_j)=g)} \sum_{n \in (y_{nj}=w)} \rho_{nj}^K$$

$$N_{tK}^g = \sum_w N_{tK}^{gw}$$

$$N_{T_j K}^{G(S_j)w} + \sum_{n \in (y_{nj}=w)} \rho_{nj}^K$$

key ingredients for weakly supervised learning

Initialization

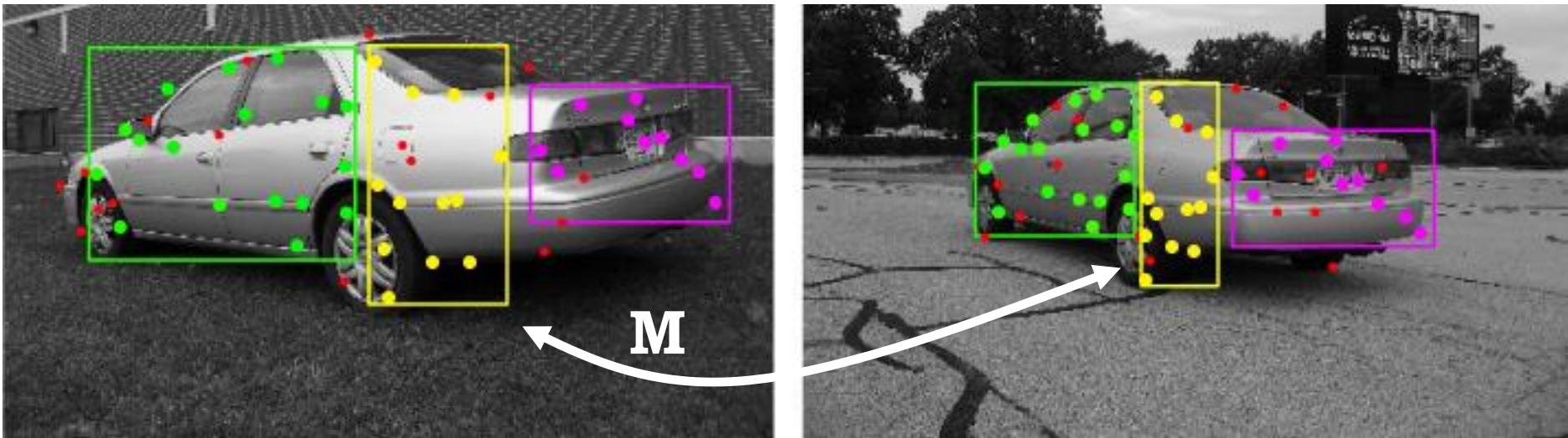
Constraints

- Across Triangle
- Within Triangle

Initialization



Initialization



$$\pi : \mathbf{I}^h \rightarrow \{ \mathbf{P}_1^h, \mathbf{P}_2^h, \mathbf{P}_3^h, \mathbf{O}^h \}$$
$$\tau : \mathbf{I}^k \rightarrow \{ \mathbf{P}_1^k, \mathbf{P}_2^k, \mathbf{P}_3^k, \mathbf{O}^k \}$$

Sequential ransac
J-linkage

key ingredients for weakly supervised learning

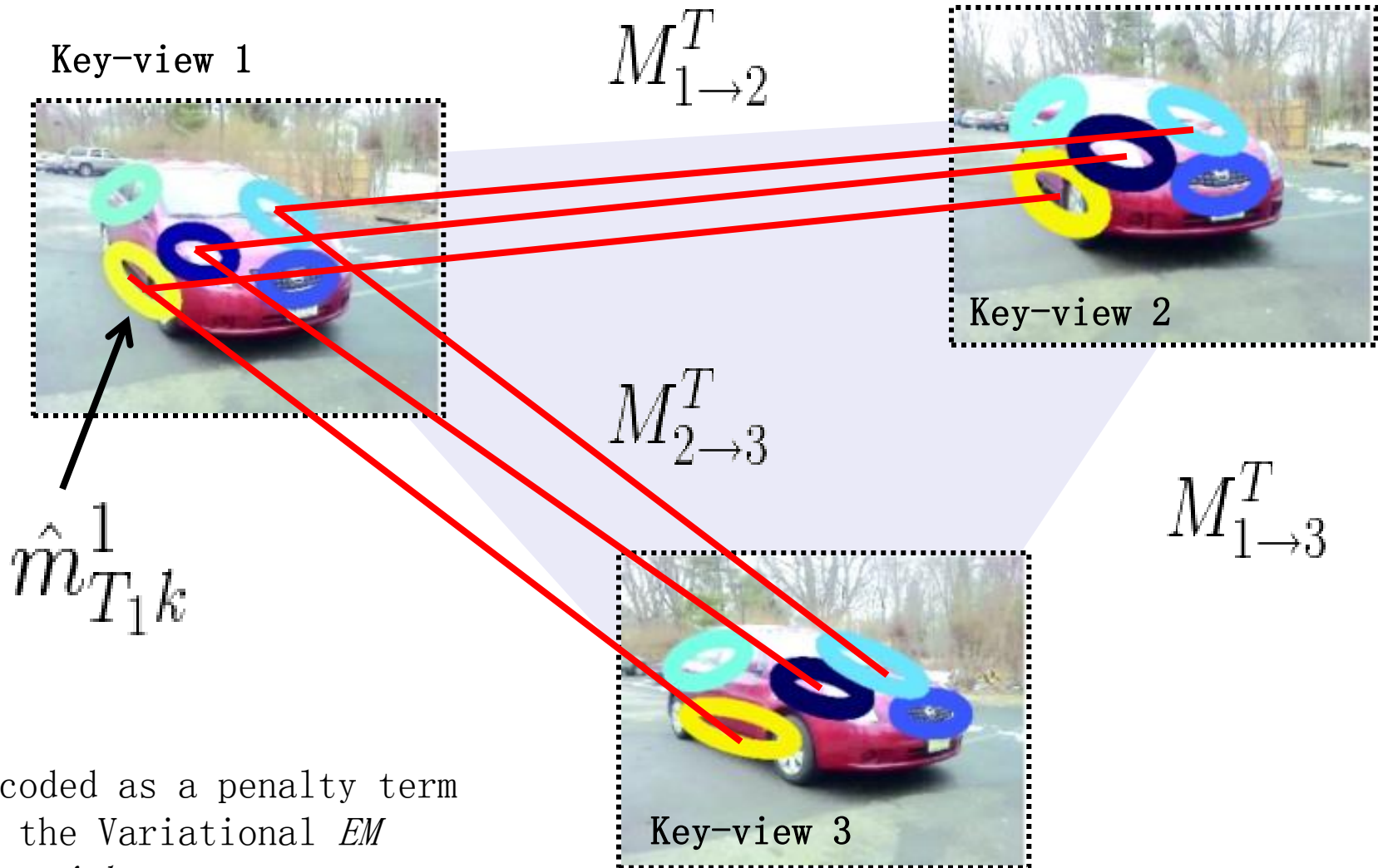
Initialization

Constraints

- Across Triangle
- Within Triangle

Constraints

Within Triangle Constraints: $M_{i \rightarrow j}^T \cdot \hat{m}_{Tk}^i \approx \hat{m}_{Tk}^j$



Encoded as a penalty term
in the Variational *EM*
algorithm

3D object class dataset

bicycle



car



iron



mouse



shoe



stapler

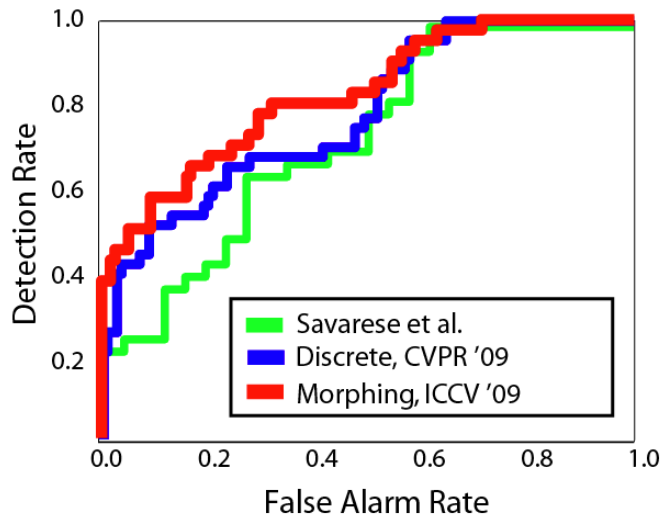


toaster

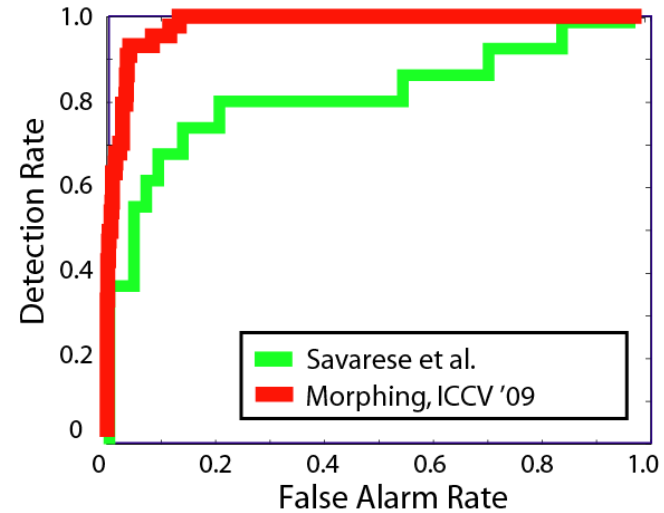


Detection

Car

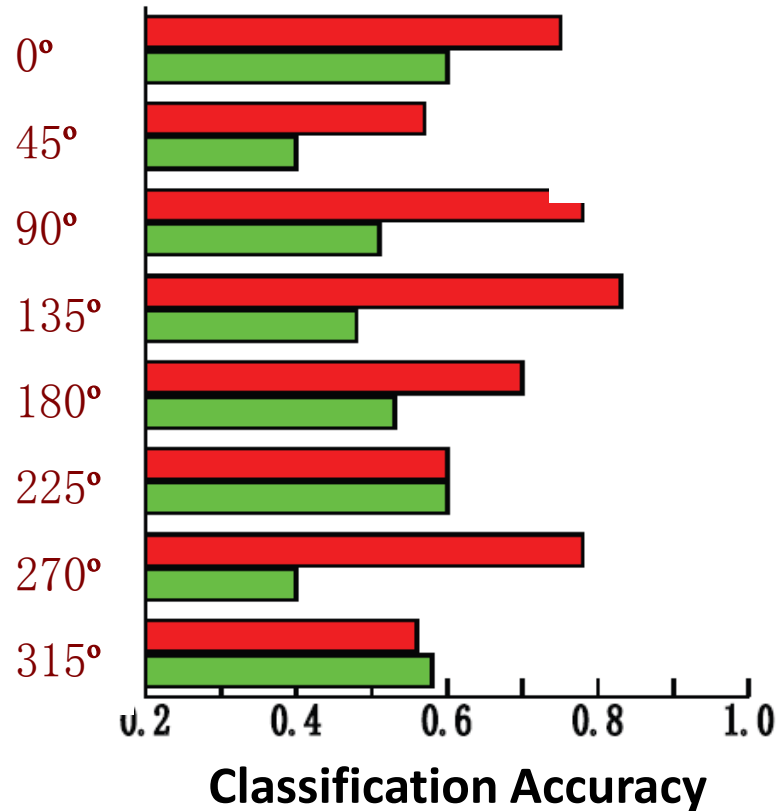


Bicycle



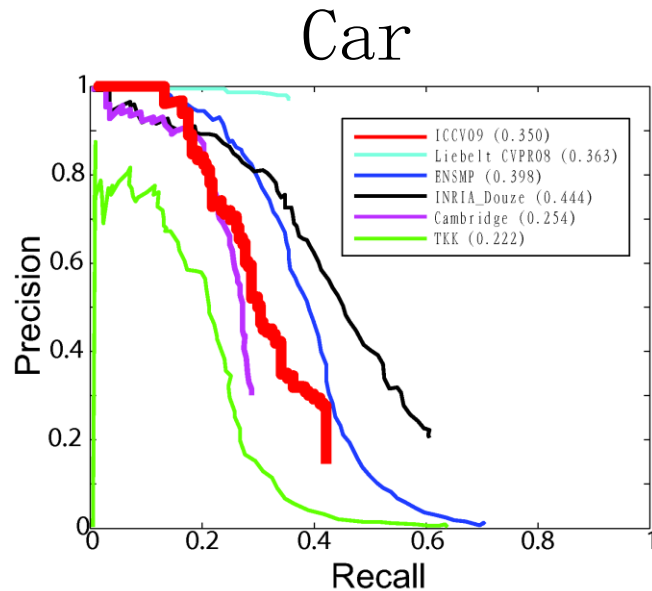
- Morphing model, 2009
- Savarese, & Fei-Fei ICCV '07

Viewpoint Classification: Car

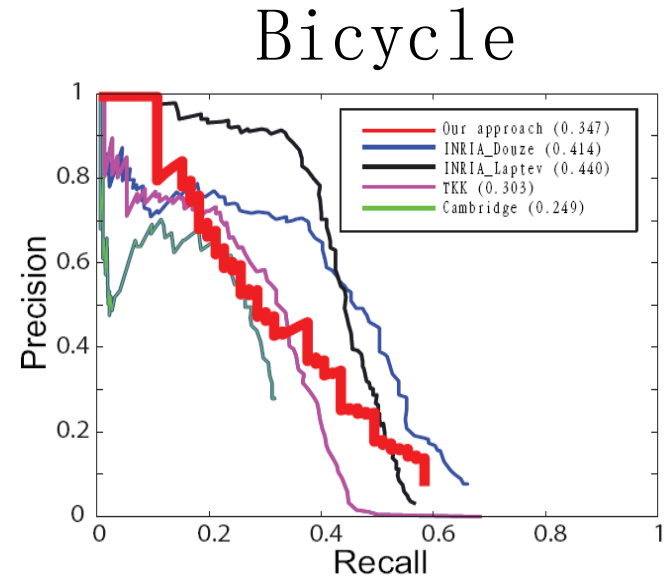


- Morphing model, 2009
- Savarese, & Fei-Fei ICCV '07

Detection: Pascal 2006 dataset

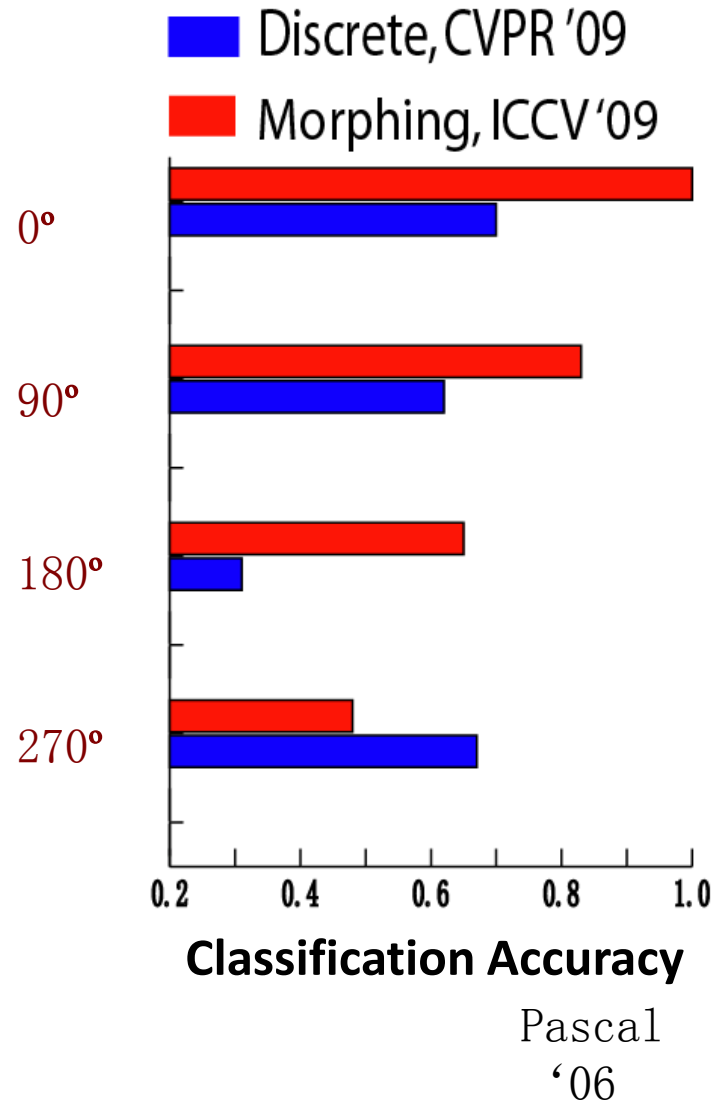


ICCV09 Morphing Model
0.35(average p)

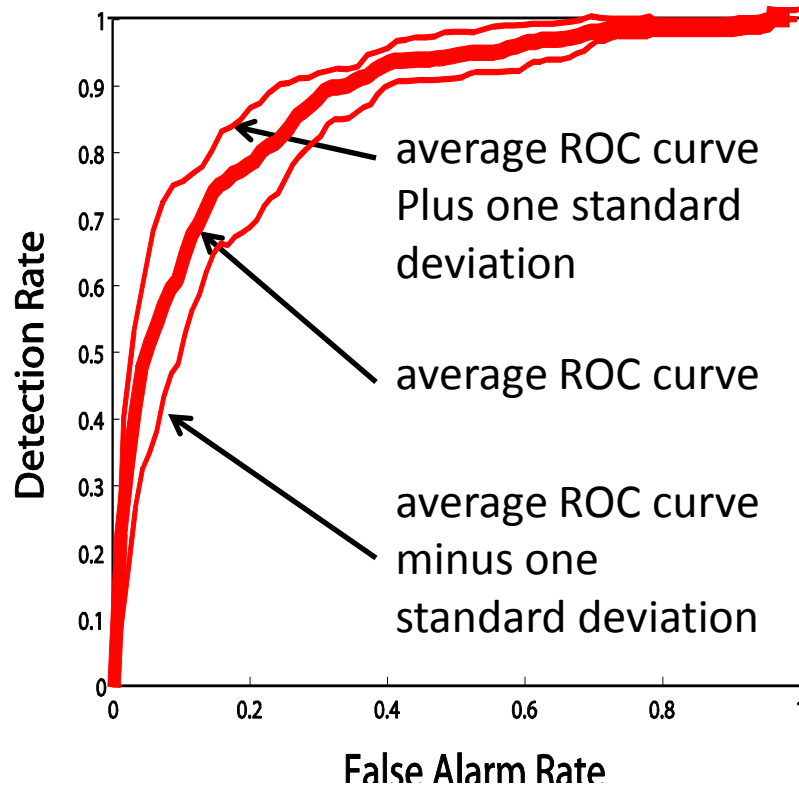


ICCV09 Morphing Model
0.347(average p)

Viewpoint Classification: Car- Pascal 2006 dataset



Household Item Dataset: Detection



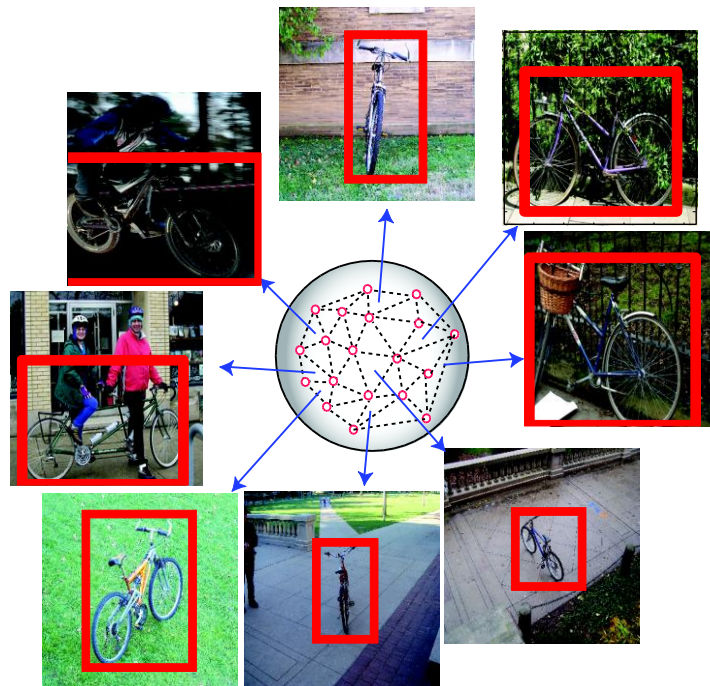
Object Class	AUC
Sewing Machine	98.1
Microscope	87.9
Travel Iron	88.1
Swivel Chair	91.2
Calculator	97.2
Flashlight	87.1
Teapot	86.4
Watch	84.9
All	90.1

Household Item Dataset: 8-Viewpoint Classification

Object Class	Accuracy
Sewing Machine	71.4
Microscope	63.9
Travel Iron	73.5
Swivel Chair	58.6
Calculator	69.2
Flashlight	68.4
Teapot	60.0
Watch	61.9
All	70.2

Typical Examples

Bicycle



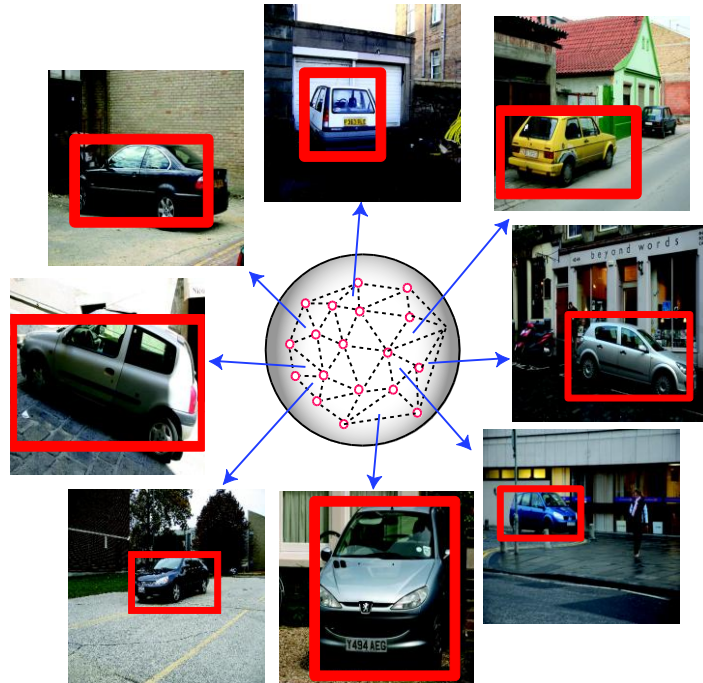
Binocular Microscope



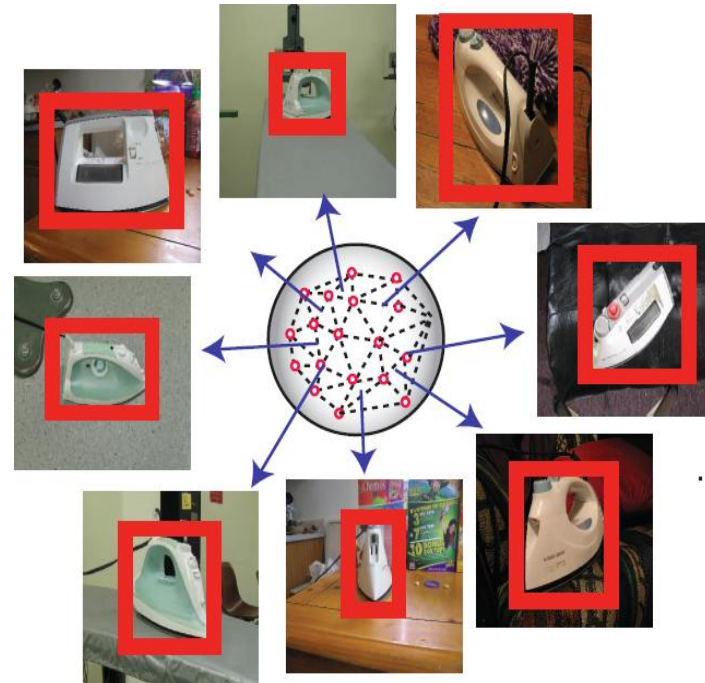
Blue arrows indicate the viewpoint for the detected object (in red bounding box).

Typical Examples

Car



Travel Iron



Blue arrows indicate the viewpoint for the detected object (in red bounding box).

Novel view object synthesis from a single image

For the first time!

car calculator flashlight bike

For natural or artificial scenes, see [hoeim 07](#); [saxena](#)

Conclusions

	Single view	Mixture / Multi-view	Sav. et al, 07	Morphing model
View point invariant	X	✓	✓	✓
No supervision	✓	X Category View point	X → ✓ category ; all views all instances available	X → ✓ category
# Categories	~300	2	8	16
Share information across views	X	✓	✓	✓
View synthesis	X	X	X	X → ✓
Pose estimation	X → ✓	X	X → ✓	✓

Thank you!