

# A System for Appearance-Based Probabilistic 3D Object Recognition and Its Applications

13/12/2007, Knowledge Engineering Group  
University of Economics, Prague



Marcin Grzegorzek

Multimedia & Vision Research Group  
Queen Mary, University of London, UK



# Overview

- 1 Fundamental Concept
- 2 Statistical Modeling
- 3 Classification and Localization
- 4 Experiments and Results
- 5 System Applications
- 6 Conclusion

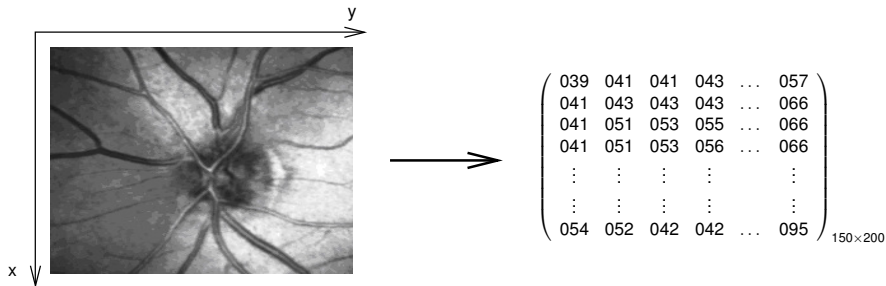


# Overview

- 1** Fundamental Concept
- 2 Statistical Modeling
- 3 Classification and Localization
- 4 Experiments and Results
- 5 System Applications
- 6 Conclusion



# Digital Representation of Gray Level Images



$$f(x, y) = \begin{pmatrix} f(0, 0) & f(0, 1) & \dots & f(0, 199) \\ f(1, 0) & f(1, 1) & \dots & f(1, 199) \\ \vdots & \vdots & & \vdots \\ f(149, 0) & f(149, 1) & \dots & f(149, 199) \end{pmatrix} ; f(x, y) \in \{0, 1, 2, \dots, 255\}$$



# Digital Representation of Color Images



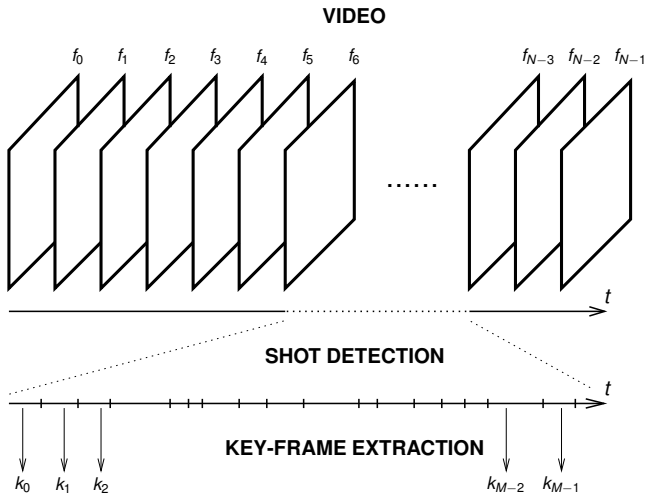
$$R = \begin{pmatrix} 178 & 182 & 182 & \dots & 185 \\ 179 & 180 & 182 & \dots & 185 \\ 179 & 181 & 182 & \dots & 184 \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ 032 & 034 & 037 & \dots & 111 \end{pmatrix}_{160 \times 240}$$

$$G = \begin{pmatrix} 067 & 069 & 069 & \dots & 066 \\ 067 & 066 & 070 & \dots & 067 \\ 068 & 067 & 069 & \dots & 067 \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ 101 & 104 & 104 & \dots & 098 \end{pmatrix}_{160 \times 240}$$

$$B = \begin{pmatrix} 189 & 188 & 188 & \dots & 188 \\ 190 & 189 & 190 & \dots & 189 \\ 190 & 188 & 190 & \dots & 189 \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ 198 & 200 & 200 & \dots & 065 \end{pmatrix}_{160 \times 240}$$



# Digital Video Representation





# Image Acquisition - Stochastic Process



$$f(120, 180) = 219$$



$$f(120, 180) = 210$$



$$f(120, 180) = 208$$



$$f(120, 180) = 204$$



$$f(120, 180) = 198$$



# Digital Image Processing

Input Image  $f = f(x, y)$



$$T\{f\} = g$$



Output Image  $g = g(x, y)$

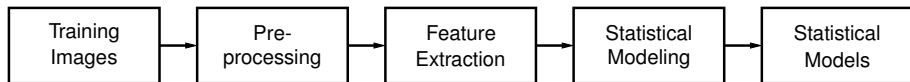




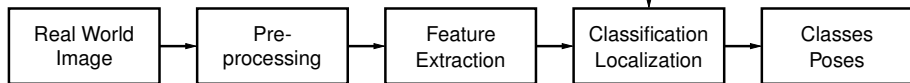


# Object Recognition Task

## TRAINING



## RECOGNITION



**Classification** – Which objects occur in the image?

**Localization** – In which poses known objects occur in the image?

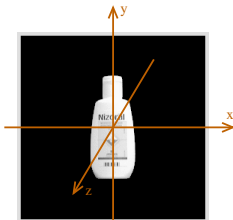
**Recognition** – Which objects and in which poses occur in the image?



# Object Pose

$$\mathbf{t}_{int} = (t_x, t_y)^T$$

$$\phi_{int} = \phi_z$$



$$\mathbf{t} = (t_x, t_y, t_z)^T = (0, 0, 100)^T$$

$$\phi = (\phi_x, \phi_y, \phi_z)^T = (0, 0, 0)^T$$



$$\mathbf{t} = (50, 25, 100)^T$$

$$\phi = (0, 0, 0)^T$$



$$\mathbf{t} = (50, 25, 100)^T$$

$$\phi = (0, 0, -30)^T$$

$$t_{ext} = t_z$$

$$\phi_{ext} = (\phi_x, \phi_y)^T$$



$$\mathbf{t} = (0, 0, 100)^T$$

$$\phi = (22.5, 0, 0)^T$$



$$\mathbf{t} = (0, 0, 100)^T$$

$$\phi = (0, 45, 0)^T$$



$$\mathbf{t} = (0, 0, 80)^T$$

$$\phi = (0, 0, 0)^T$$

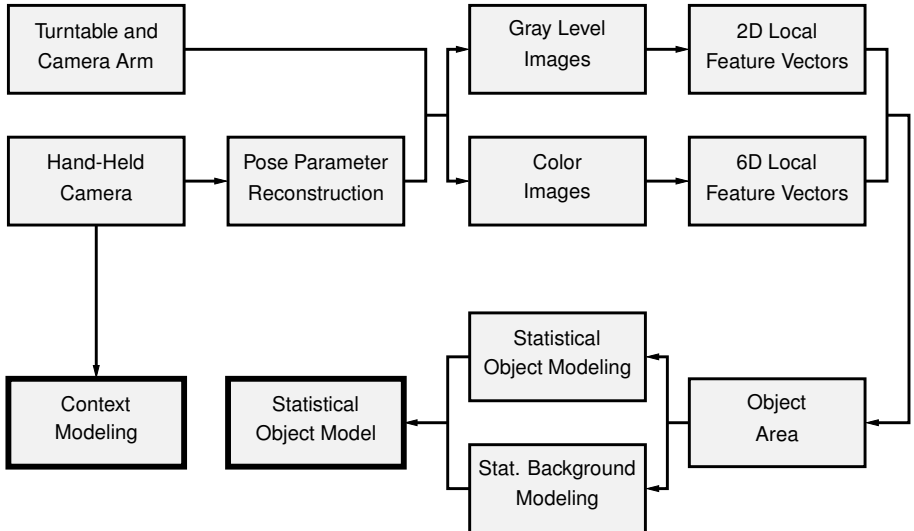


# Overview

- 1 Fundamental Concept
- 2 Statistical Modeling**
- 3 Classification and Localization
- 4 Experiments and Results
- 5 System Applications
- 6 Conclusion

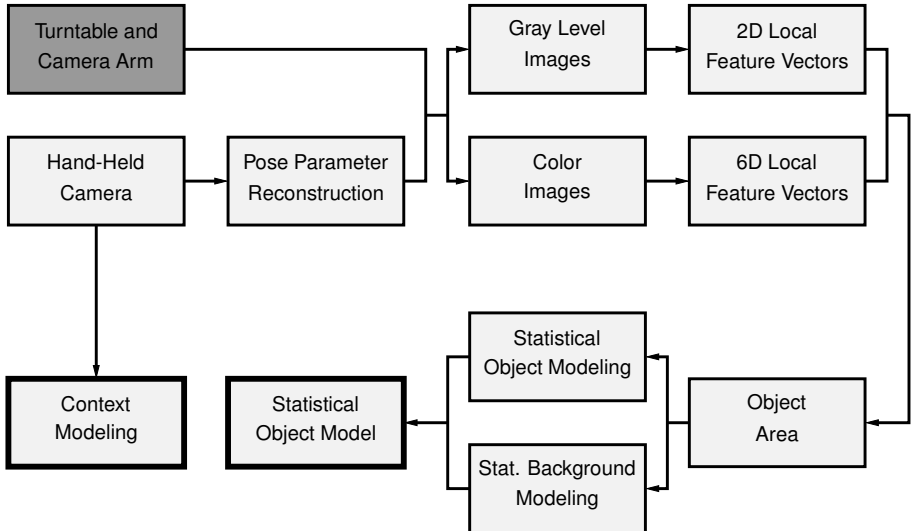


# Training Phase



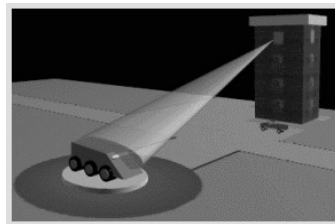
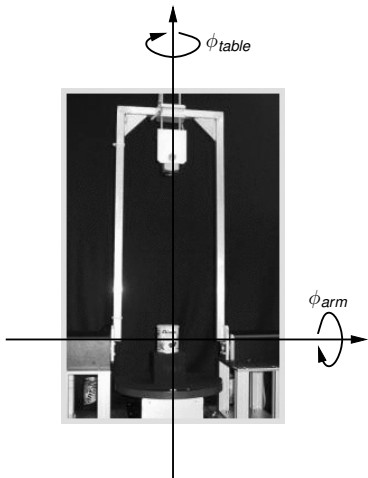


# Turntable and Camera Arm





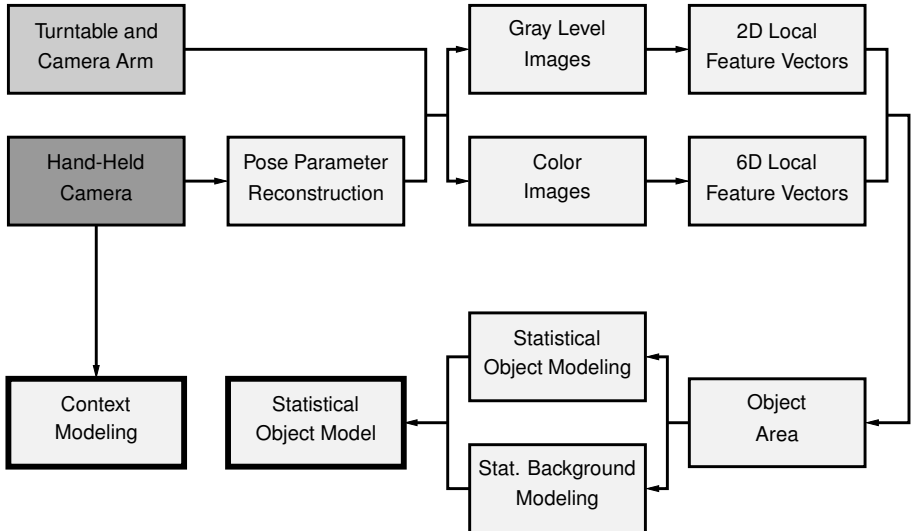
# Turntable and Camera Arm



Object poses  $(\phi_\rho, \mathbf{t}_\rho)$  for all  $N_\rho$  training images  $\mathbf{f}_{\rho=1, \dots, N_\rho}$  are known.



# Hand-Held Camera





# Hand-Held Camera

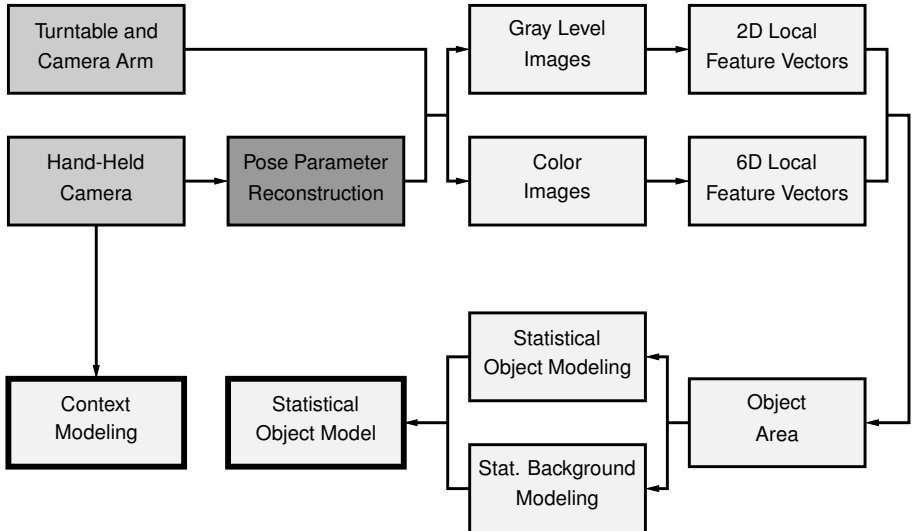


Object poses  $(\phi_\rho, \mathbf{t}_\rho)$  for the training images  $\mathbf{f}_\rho$  are unknown.



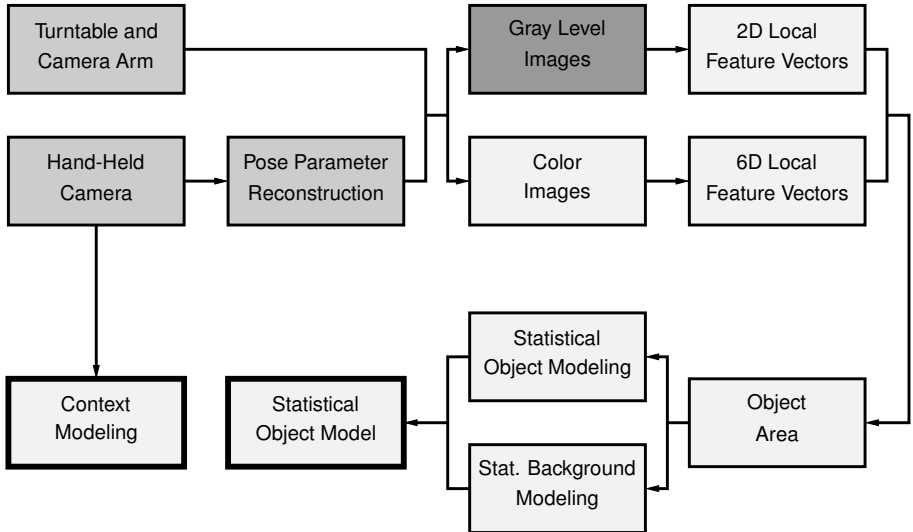


# Pose Parameter Reconstruction



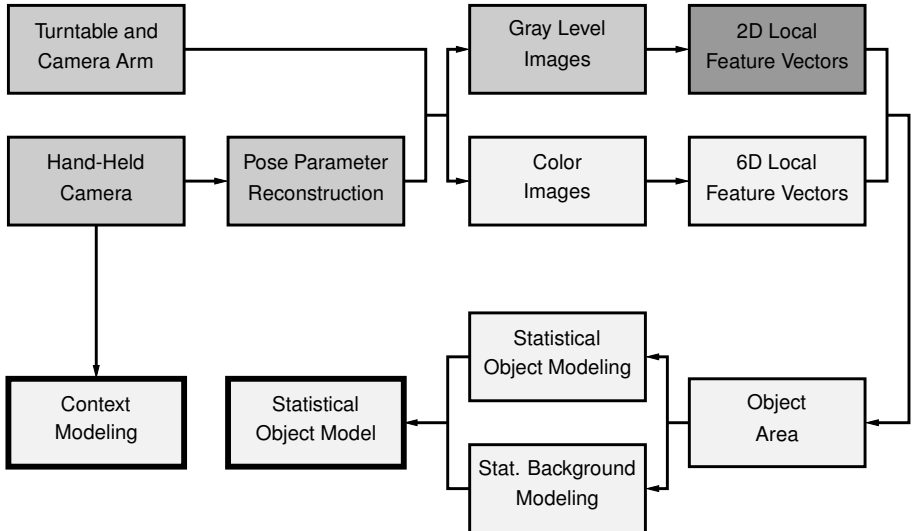


# Gray Level Images





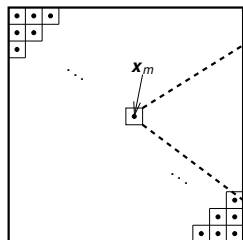
# 2D Local Feature Vectors





# 2D Feature Extraction with Wavelet Transform

$2^n \times 2^n$  Gray Level Image



$s = 0$

|             |             |             |             |
|-------------|-------------|-------------|-------------|
| $b_{s,0,0}$ | $b_{s,1,0}$ | $b_{s,2,0}$ | $b_{s,3,0}$ |
| $b_{s,0,1}$ | $b_{s,1,1}$ | $b_{s,2,1}$ | $b_{s,3,1}$ |
| $b_{s,0,2}$ | $b_{s,1,2}$ | $b_{s,2,2}$ | $b_{s,3,2}$ |
| $b_{s,0,3}$ | $b_{s,1,3}$ | $b_{s,2,3}$ | $b_{s,3,3}$ |

$s = -1$

|             |             |           |
|-------------|-------------|-----------|
| $b_{s,0,0}$ | $b_{s,1,0}$ | $d_{2,s}$ |
| $b_{s,0,1}$ | $b_{s,1,1}$ |           |
| $d_{0,s}$   |             | $d_{1,s}$ |

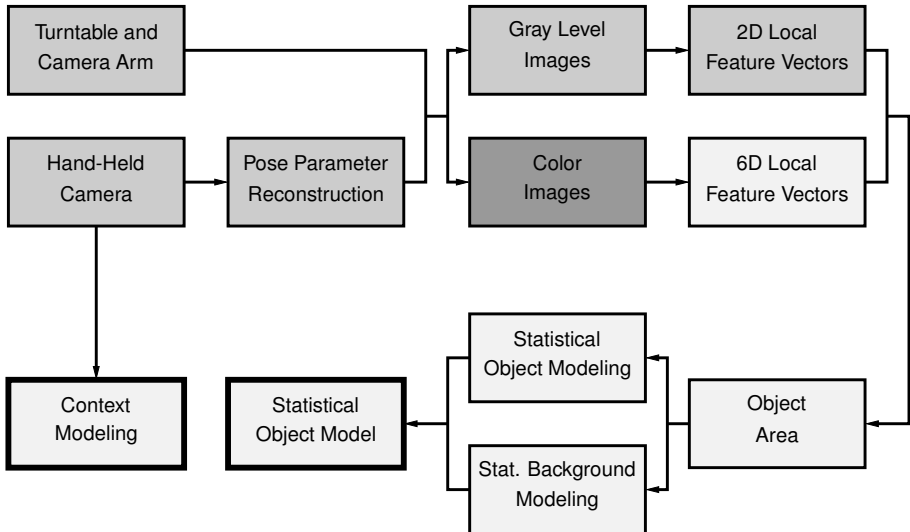
$s = \hat{s} = -2$

|             |           |             |
|-------------|-----------|-------------|
| $b_s$       | $d_{2,s}$ | $d_{2,s+1}$ |
| $d_{0,s}$   | $d_{1,s}$ |             |
| $d_{0,s+1}$ |           | $d_{1,s+1}$ |

$$\mathbf{c}_m = \mathbf{c}(\mathbf{x}_m) = \begin{pmatrix} c_{m,1} \\ c_{m,2} \end{pmatrix} = \begin{pmatrix} \ln(2^s |b_s|) \\ \ln[2^{2s} (|d_{0,s}| + |d_{1,s}| + |d_{2,s}|)] \end{pmatrix}$$

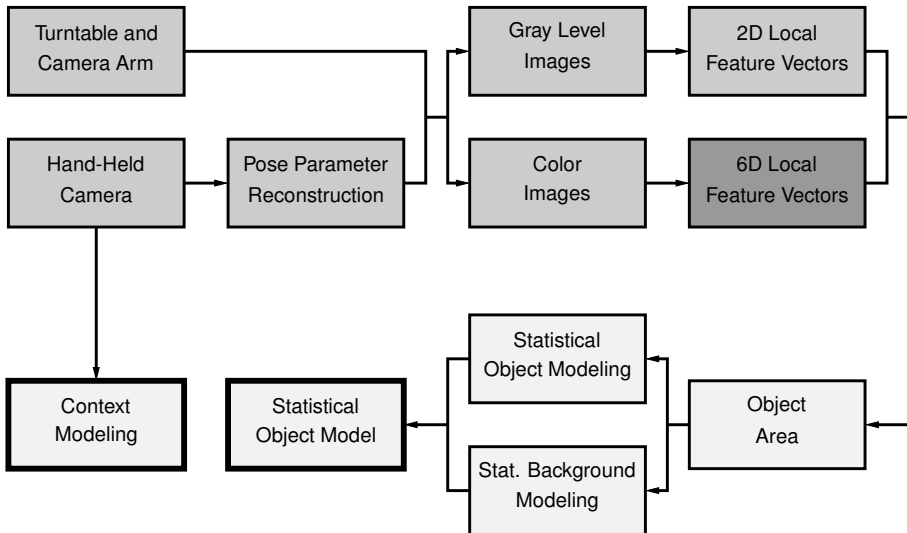


# Color Images



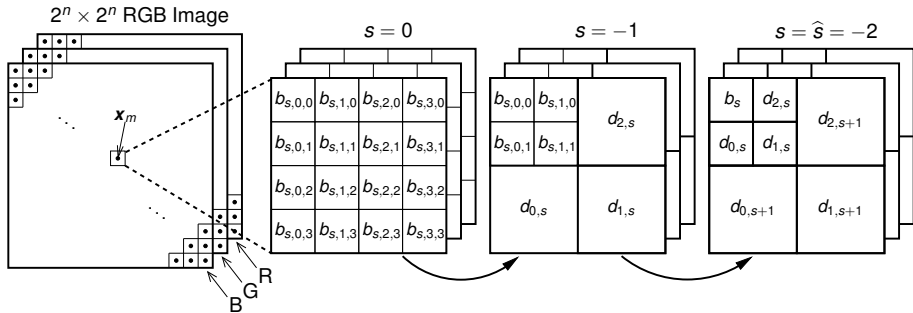


# 6D Local Feature Vectors





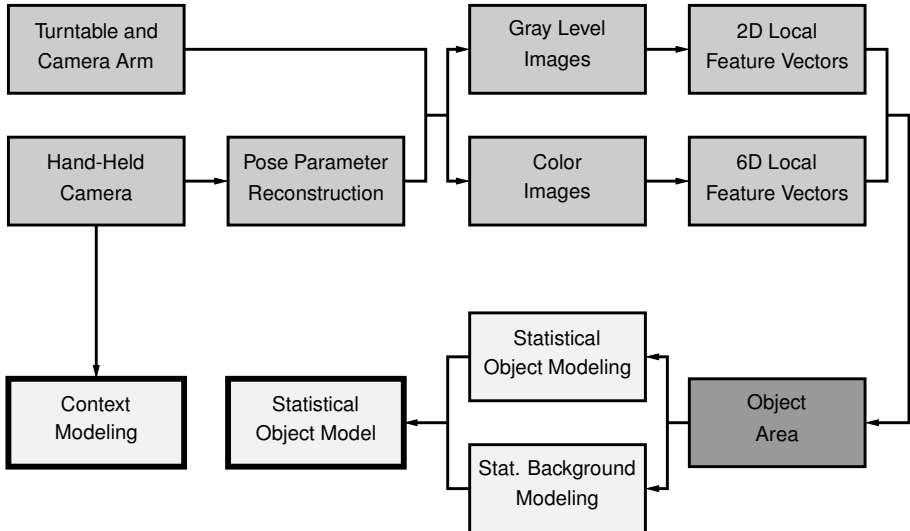
# 6D Feature Extraction with Wavelet Transform



$$\mathbf{c}_m = \mathbf{c}(\mathbf{x}_m) = \begin{pmatrix} C_{m,r,1} \\ C_{m,r,2} \\ C_{m,g,1} \\ C_{m,g,2} \\ C_{m,b,1} \\ C_{m,b,2} \end{pmatrix}$$



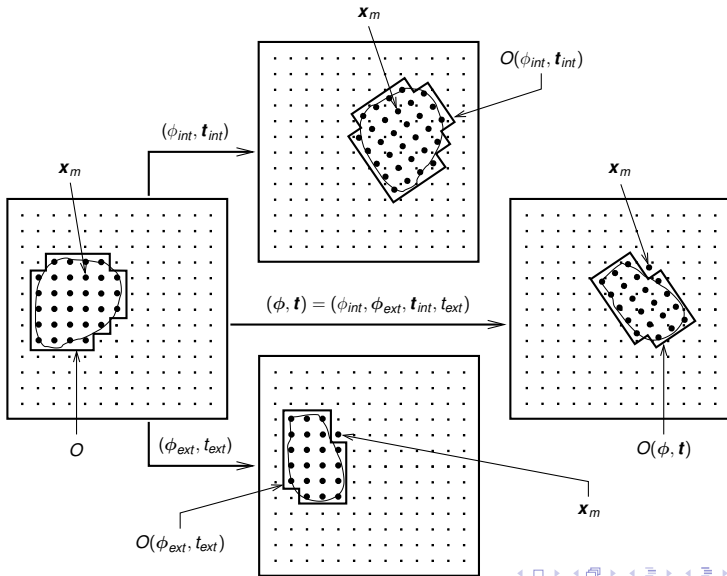
# Object Area





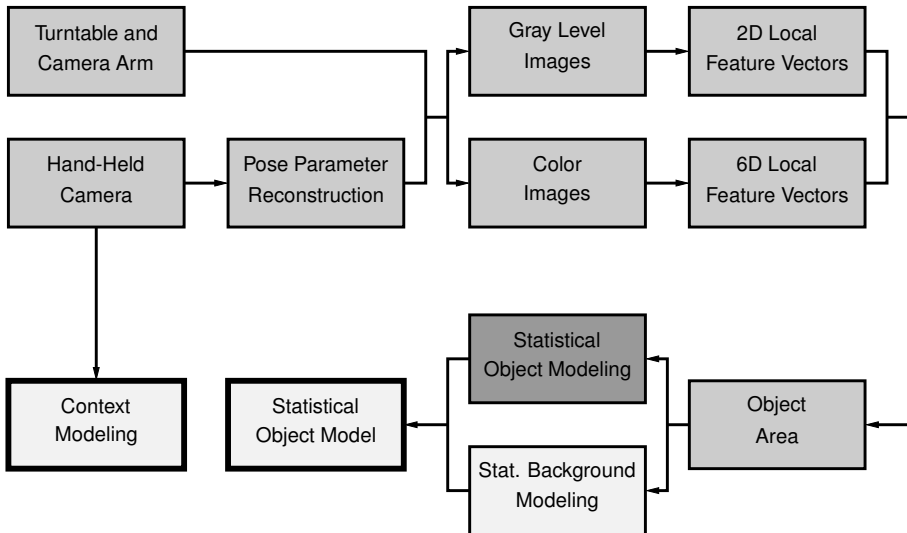


# Object Area $O = O(\phi, \mathbf{t})$





# Statistical Object Modeling





# Object Features as Normal Random Variables

- The elements  $c_{m,q}$  of  $\mathbf{c}_m \in C_O$  are considered as normal random variables

$$p(c_{m,q} | \mu_{m,q}, \sigma_{m,q}, \phi, \mathbf{t}) = \frac{1}{\sigma_{m,q} \sqrt{2\pi}} \exp \left( \frac{(c_{m,q} - \mu_{m,q})^2}{-2\sigma_{m,q}^2} \right)$$

- Assuming the statistical independence of the elements  $c_{m,q}$

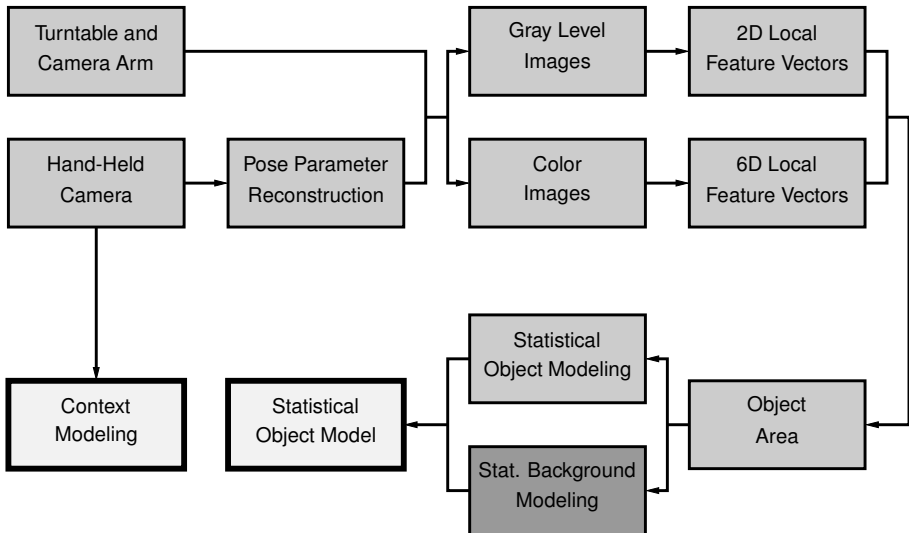
$$p(\mathbf{c}_m | \mu_m, \sigma_m, \phi, \mathbf{t}) = \prod_{q=1}^{N_q} p(c_{m,q} | \mu_{m,q}, \sigma_{m,q}, \phi, \mathbf{t})$$

$N_q = 2$  for gray level modeling

$N_q = 6$  for color modeling



# Statistical Background Modeling





# Background Features as Uniform Random Variables

- The elements  $c_{m,q}$  of  $\mathbf{c}_m \notin C_O$  are considered as uniform random variables

$$p(c_{m,q}) = \frac{1}{\max(c_{m,q}) - \min(c_{m,q})}$$

- Assuming the statistical independence of the elements  $c_{m,q}$

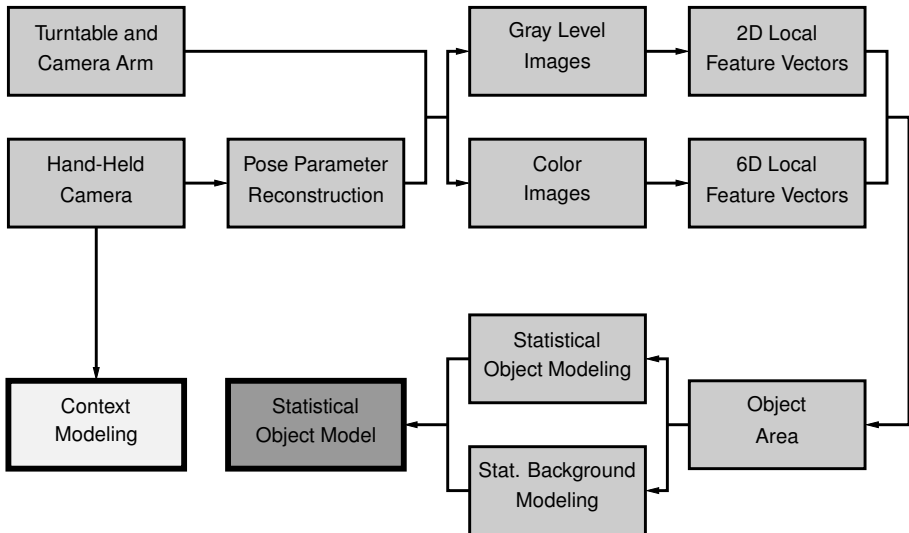
$$p(\mathbf{c}_m) = \prod_{q=1}^{N_q} \frac{1}{\max(c_{m,q}) - \min(c_{m,q})} = p_b$$

$N_q = 2$  for gray level modeling

$N_q = 6$  for color modeling



# Statistical Object Model





# Statistical Object Model - Summary

For all  $\Omega_{\kappa=1, \dots, N_{\Omega}}$  the system creates statistical models  $\mathcal{M}_{\kappa=1, \dots, N_{\Omega}}$

$$\mathcal{M}_{\kappa} = \mathcal{M}_{\kappa}(\phi, \mathbf{t})$$

containing:

- Object area and set of object feature vectors

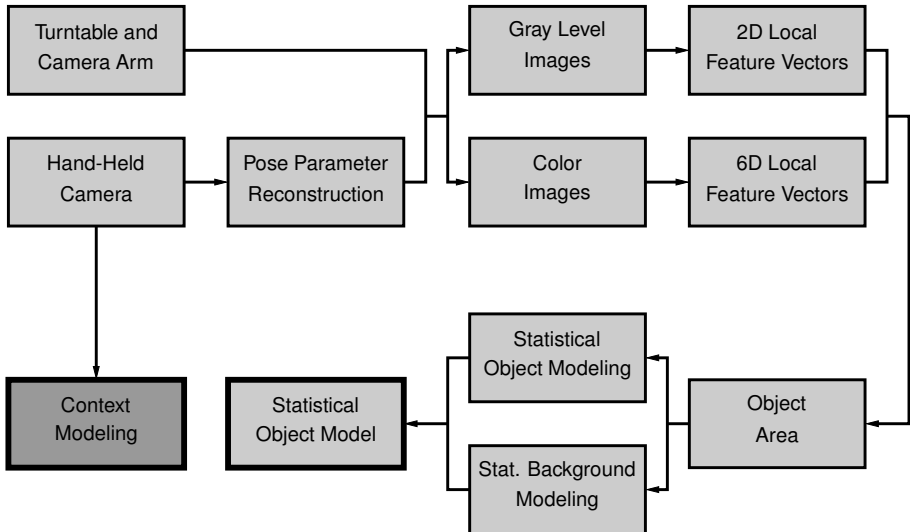
$$O_{\kappa} = O_{\kappa}(\phi, \mathbf{t}) \quad C_{O_{\kappa}} = C_{O_{\kappa}}(\phi, \mathbf{t})$$

- Object density and background density

$$p(\mathbf{c}_m | \boldsymbol{\mu}_m, \boldsymbol{\sigma}_m, \phi, \mathbf{t}) \quad p_b$$



# Context Modeling







# Motivation for Context Modeling



- Without context modeling it is assumed that

$$p(\Omega_1) = \dots = p(\Omega_\kappa) = \dots = p(\Omega_{N_\Omega})$$

- Considering context dependencies the a-priori probabilities  $p(\Omega_\kappa)$  cannot be assumed to be equal and they have to be trained

$$p(\Omega_1) \neq \dots \neq p(\Omega_\kappa) \neq \dots \neq p(\Omega_{N_\Omega})$$



# Training of A-Priori Probabilities $p(\Omega_\kappa)$

- Set of contexts is introduced

$$\Upsilon = \{\Upsilon_1, \Upsilon_2, \dots, \Upsilon_\ell, \dots, \Upsilon_{N_\Upsilon}\}$$

- $N_\ell$  images for each context  $\Upsilon_\ell$  are taken
- $N_{\ell, \kappa}$  denotes how often  $\Omega_\kappa$  occurs in  $\Upsilon_\ell$
- The a-priori probability for  $\Omega_\kappa$  in  $\Upsilon_\ell$  is defined by

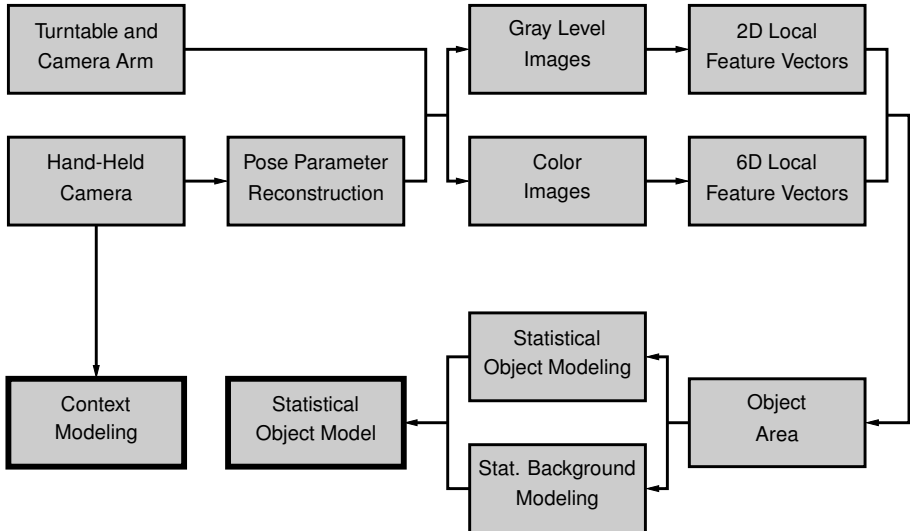
$$p_\ell(\Omega_\kappa) = \eta_\ell N_{\ell, \kappa}$$

- By  $\eta_\ell$  a normalization factor is denoted so that

$$\eta_\ell [p_\ell(\Omega_1) + p_\ell(\Omega_2) + \dots + p_\ell(\Omega_\kappa) + \dots + p_\ell(\Omega_{N_\Omega})] = 1$$



# Training Phase Completed



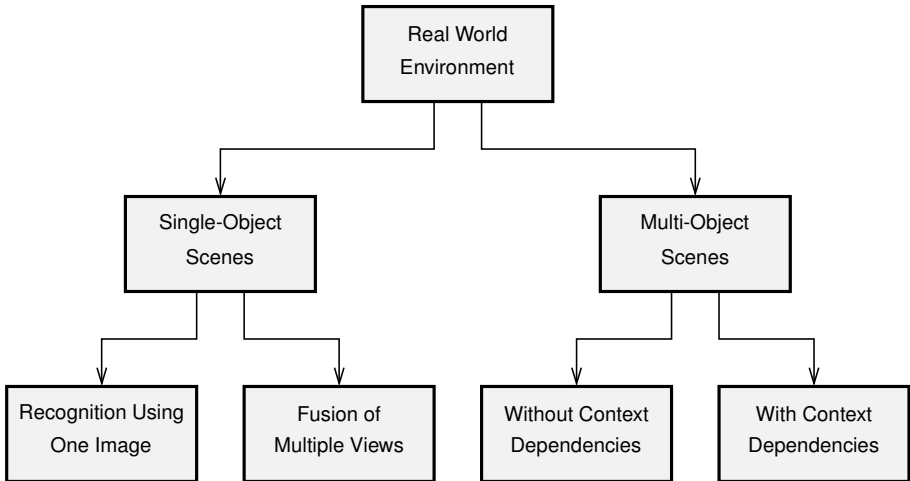


# Overview

- 1 Fundamental Concept
- 2 Statistical Modeling
- 3 Classification and Localization**
- 4 Experiments and Results
- 5 System Applications
- 6 Conclusion

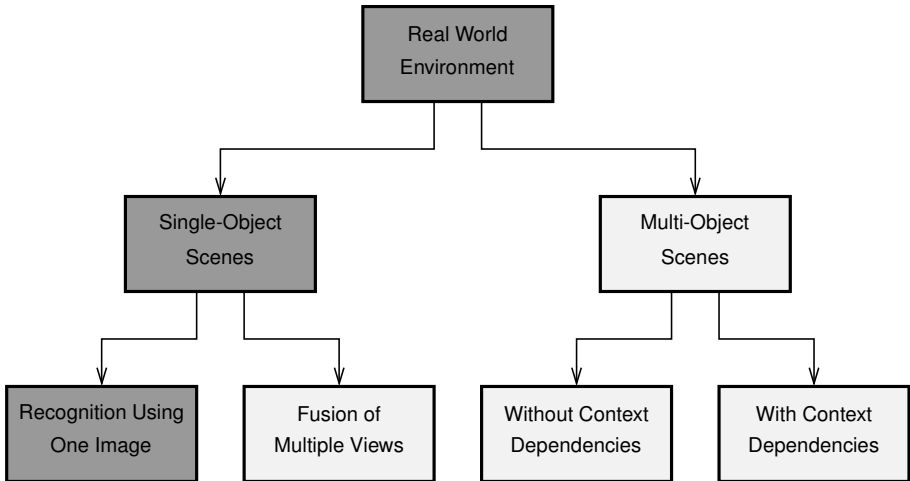


# Recognition Phase





# Single-Object, One Image





# Maximum Likelihood Estimation

## Assumptions

- Observation  $\mathbf{f} \rightarrow C$
- A-priori probabilities for all  $\Omega_{\kappa=1, \dots, N_{\Omega}}$  are equal

$$p(\Omega_1) = \dots = p(\Omega_{\kappa}) = \dots = p(\Omega_{N_{\Omega}})$$

- All  $(\phi, \mathbf{t})$  are also equiprobable

## Classification

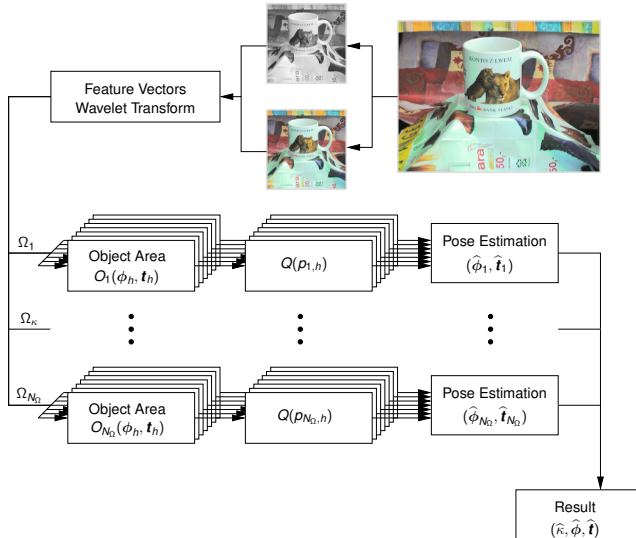
$$\hat{\kappa} = \operatorname{argmax}_{\kappa} p(\Omega_{\kappa} | C) = \operatorname{argmax}_{\kappa} \frac{p(\Omega_{\kappa})p(C|\Omega_{\kappa})}{p(C)} = \operatorname{argmax}_{\kappa} p(C|\Omega_{\kappa}) = \operatorname{argmax}_{\kappa} p(C|\mathcal{M}_{\kappa})$$

## Recognition

$$(\hat{\kappa}, \hat{\phi}, \hat{\mathbf{t}}) = \operatorname{argmax}_{(\kappa, \phi, \mathbf{t})} p(C|\mathcal{M}_{\kappa}(\phi, \mathbf{t}))$$



# Classification and Localization Algorithm







# Object Density Value

## Problem

- How to compute  $p_{\kappa,h}$  for given  $\mathbf{f}$ ,  $\mathcal{M}_{\kappa}$ , and  $(\phi_h, \mathbf{t}_h)$ ?

## Solution

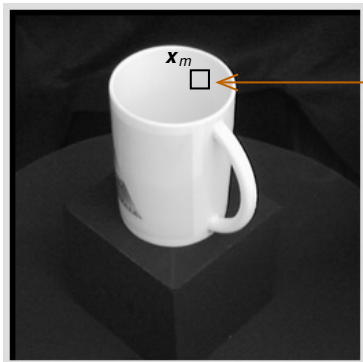
- Compute feature vectors in  $\mathbf{f}$ !
- Determine  $O_{\kappa}(\phi_h, \mathbf{t}_h)$  using  $\mathcal{M}_{\kappa}$ !
- Hence,  $C_{O_{\kappa}} = \{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_M\}$  is known.
- Determine  $(p_{\mathbf{c}_1}, p_{\mathbf{c}_2}, \dots, p_{\mathbf{c}_M})$  using  $\mathcal{M}_{\kappa}$ !
- Compute  $p_{\kappa,h}$  as follows

$$p_{\kappa,h} = \prod_{i=0}^M \max \{p_{\mathbf{c}_i}, p_b\} \quad !$$



# Motivation for Background Density

Training Image of  $\Omega_\kappa$  in  $(\phi_h, \mathbf{t}_h)$



$$\boldsymbol{\mu}_{\kappa,m} \neq \mathbf{c}_m$$

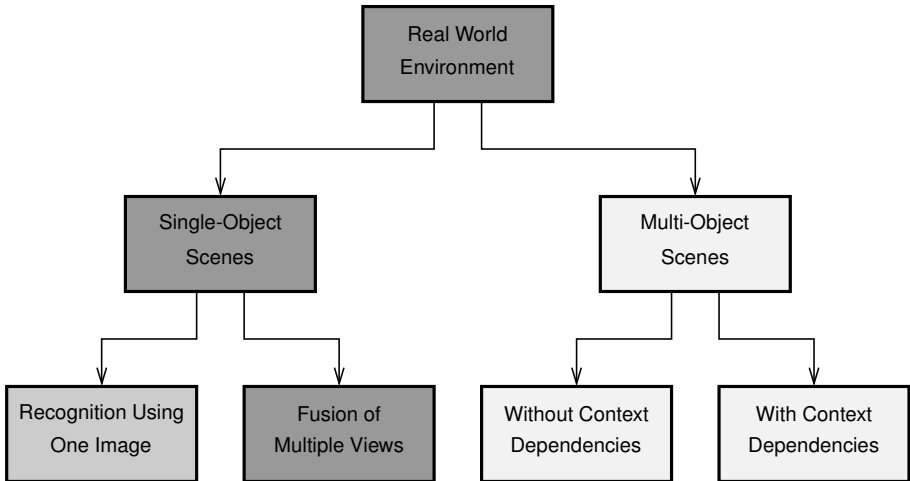
Test Image of  $\Omega_\kappa$  in  $(\phi_h, \mathbf{t}_h)$



$$p_{\mathbf{c}_m} = p(\mathbf{c}_m | \boldsymbol{\mu}_{\kappa,m}, \boldsymbol{\sigma}_{\kappa,m}, \phi_h, \mathbf{t}_h) \approx 0$$

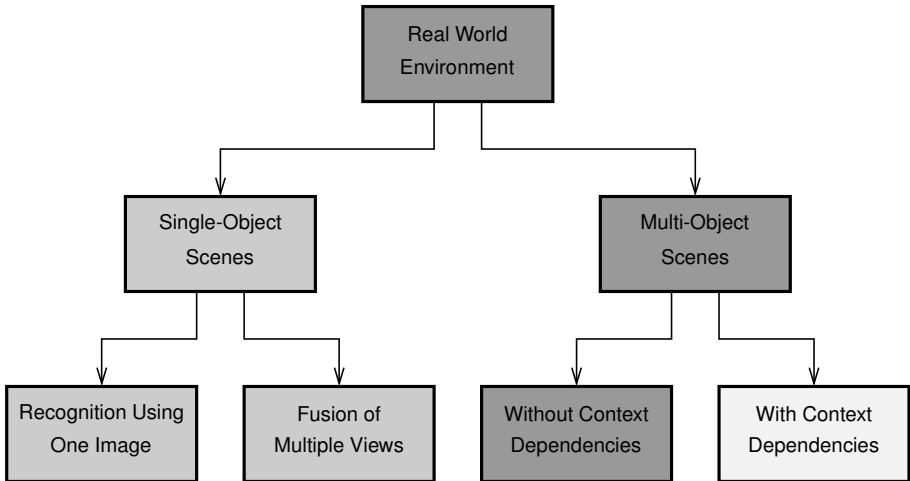


# Single-Object, Multiple Views



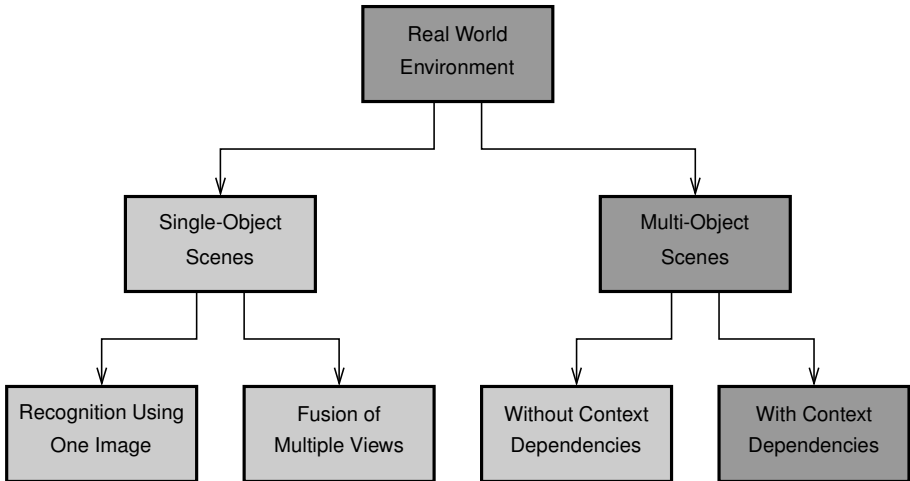


# Multi-Object Scenes without Context





# Multi-Object Scenes with Context





# Recognition Task for Multi-Object Scenes

## Given

- multi-object image  $f$

## Expected Result

|               |  |
|---------------|--|
| first object  | $(\kappa_1, \hat{\phi}_{\kappa_1}, \hat{\mathbf{t}}_{\kappa_1})$                         |
| second object | $(\kappa_2, \hat{\phi}_{\kappa_2}, \hat{\mathbf{t}}_{\kappa_2})$                         |
|               | $\vdots$   |
| last object   | $(\kappa_{\hat{i}}, \hat{\phi}_{\kappa_{\hat{i}}}, \hat{\mathbf{t}}_{\kappa_{\hat{i}}})$ |

## Unknown

- object classes and poses  $(\kappa_i, \hat{\phi}_{\kappa_i}, \hat{\mathbf{t}}_{\kappa_i})$
- number of objects  $\hat{i}$
- context  $\Upsilon_{\hat{i}}$



# Multi-Object Scenes with Context

## First Object

- recognition with ML estimation  $(\kappa_1, \hat{\phi}_{\kappa_1}, \hat{\mathbf{t}}_{\kappa_1})$
- context  $\Upsilon_{\hat{\iota}}$  determined with

$$\hat{\iota} = \underset{\iota}{\operatorname{argmax}} p_{\iota}(\Omega_{\kappa_1})$$

## Ranking of Remaining Objects

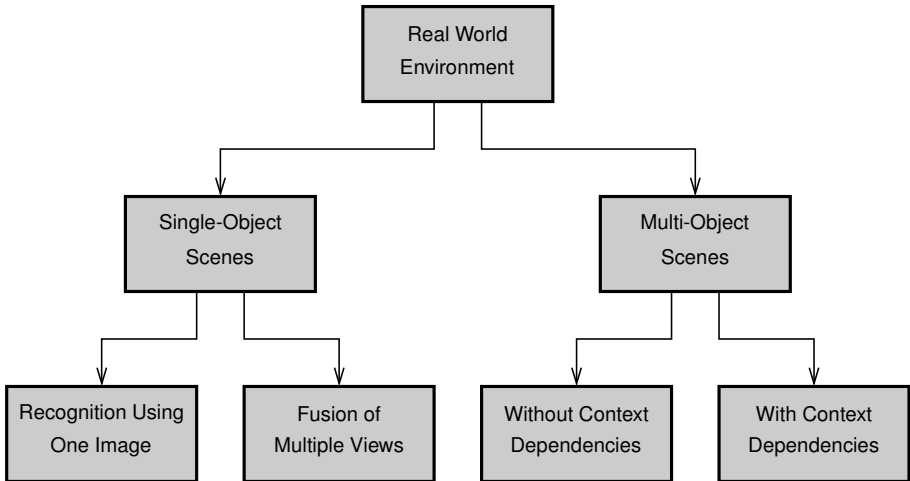
- object densities weighted with the trained a-priori probabilities

## Last Object

- determined by the highest distance in the ranking



# Recognition Phase Completed





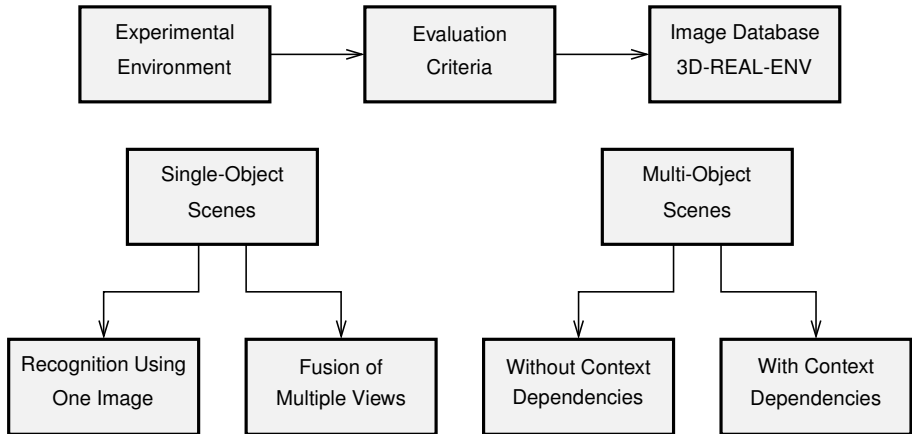


# Overview

- 1 Fundamental Concept
- 2 Statistical Modeling
- 3 Classification and Localization
- 4 Experiments and Results**
- 5 System Applications
- 6 Conclusion

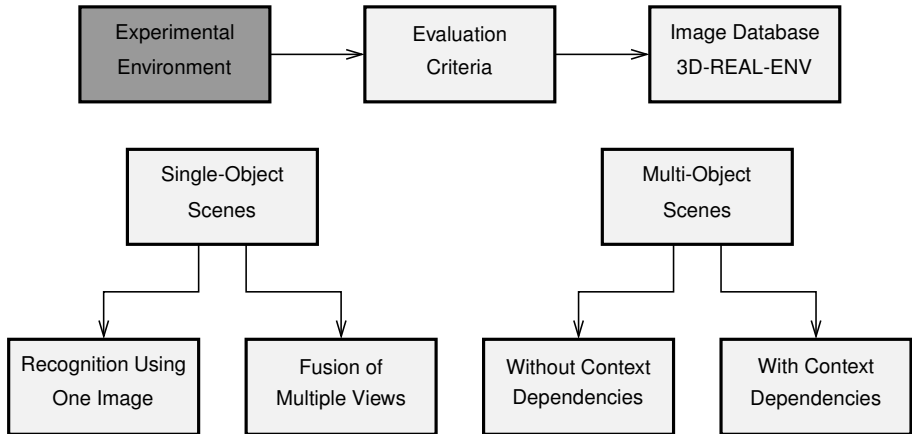


# Experiments and Results



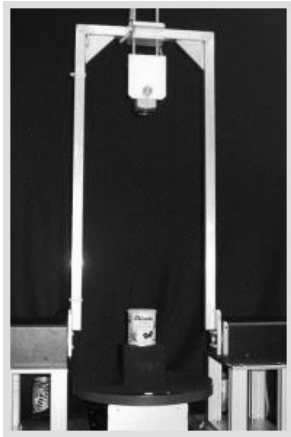


# Experimental Environment



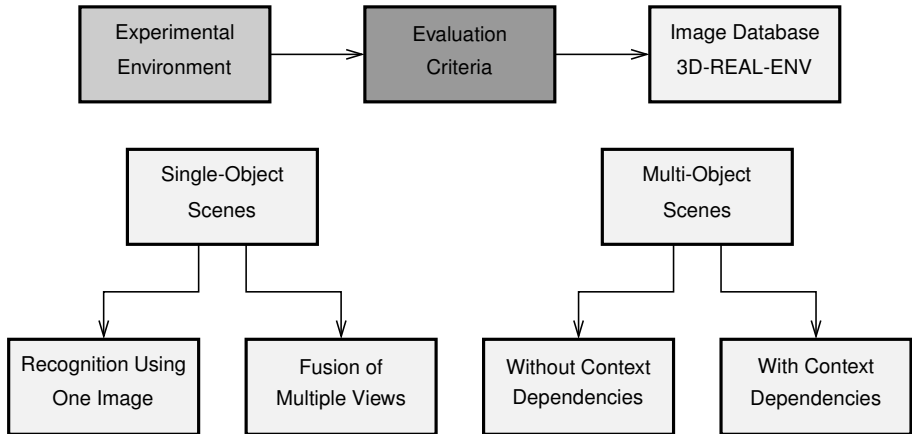


# Experimental Environment





# Evaluation Criteria





# Evaluation Criteria

- **Classification** is correct or not
- **Localization** is correct if

$$\Delta t_x \leq 10P$$

$$\Delta \phi_x \leq 15^\circ$$

$$\Delta t_y \leq 10P$$

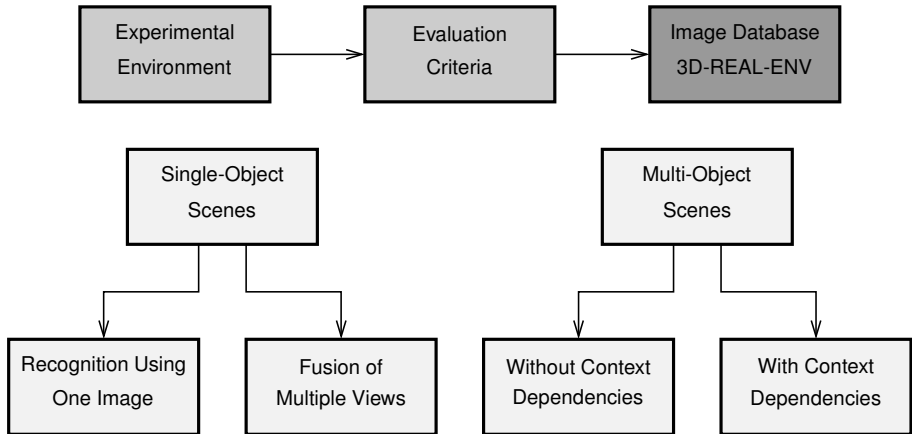
$$\Delta \phi_y \leq 15^\circ$$

$$\Delta t_z \leq 10\%$$

$$\Delta \phi_z \leq 10^\circ$$

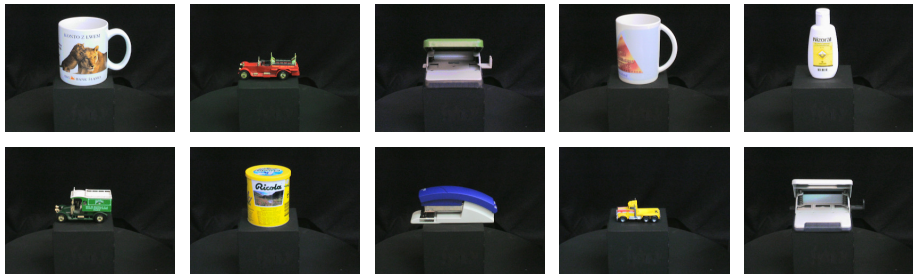


# Image Database 3D-REAL-ENV





# Training Images



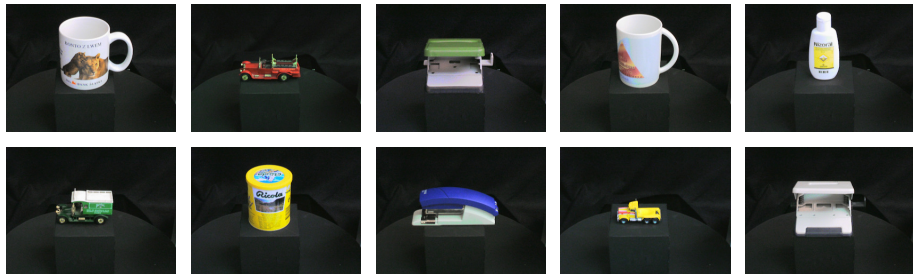
1680 Training Viewpoints, 33600 Images

$$\phi_{x,\rho} = (0.0^\circ, 4.5^\circ, 9.0^\circ, \dots, 85.5^\circ, 90.0^\circ)$$

$$\phi_{y,\rho} = (0.0^\circ, 4.5^\circ, 9.0^\circ, \dots, 351.0^\circ, 355.5^\circ)$$



# Test Images HomBack



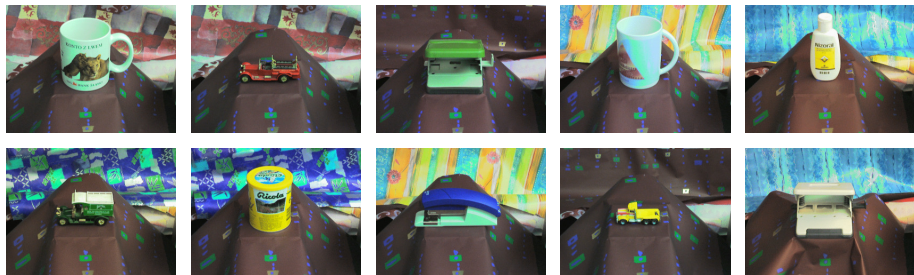
288 Test Viewpoints, 2880 Images

$$\phi_{x,\tau} = (0.00^\circ, 11.25^\circ, 22.50^\circ, \dots, 78.75^\circ, 90.00^\circ)$$

$$\phi_{y,\tau} = (0.00^\circ, 11.25^\circ, 22.50^\circ, \dots, 337.50^\circ, 348.25^\circ)$$



# Test Images LessHetBack



288 Test Viewpoints, 2880 Images

$$\phi_{x,\tau} = (0.00^\circ, 11.25^\circ, 22.50^\circ, \dots, 78.75^\circ, 90.00^\circ)$$

$$\phi_{y,\tau} = (0.00^\circ, 11.25^\circ, 22.50^\circ, \dots, 337.50^\circ, 348.25^\circ)$$



# Test Images MoreHetBack



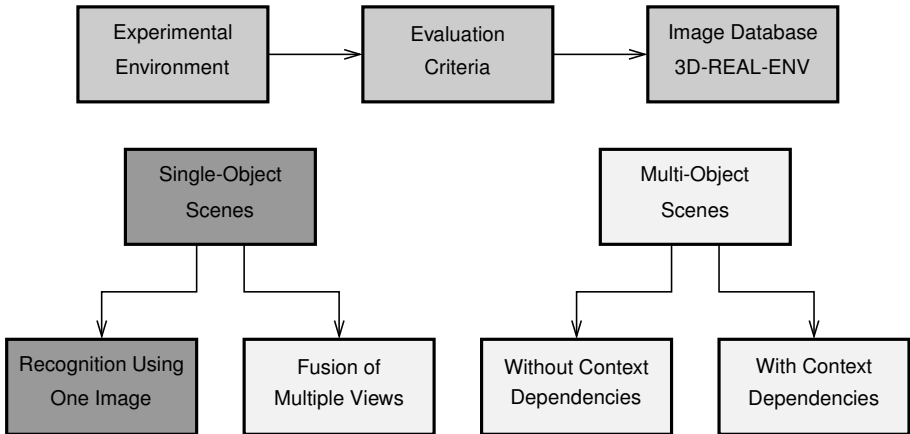
288 Test Viewpoints, 2880 Images

$$\phi_{x,\tau} = (0.00^\circ, 11.25^\circ, 22.50^\circ, \dots, 78.75^\circ, 90.00^\circ)$$

$$\phi_{y,\tau} = (0.00^\circ, 11.25^\circ, 22.50^\circ, \dots, 337.50^\circ, 348.75^\circ)$$



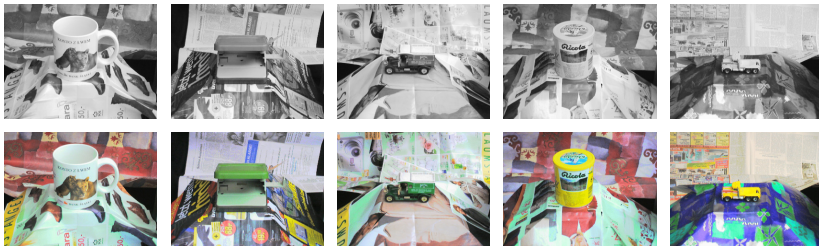
# Single-Object, One Image





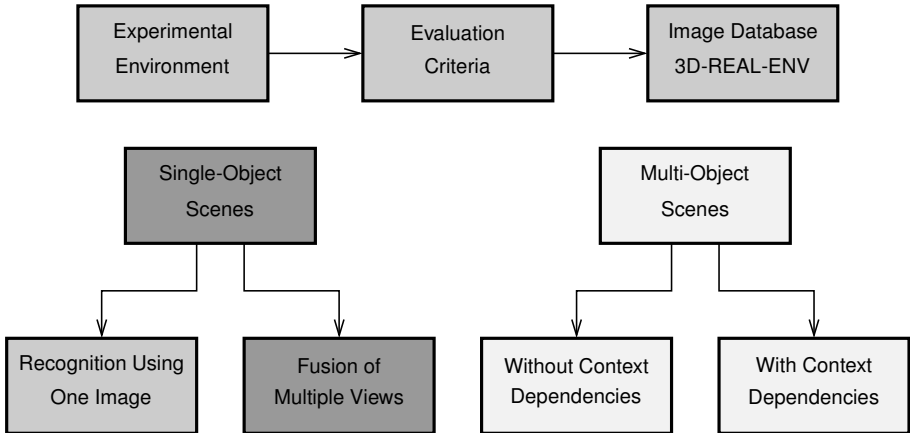
# Classification and Localization Rates

| Distance of Training Views 4.5° | Classification |                 |                 | Localization |                 |                 |
|---------------------------------|----------------|-----------------|-----------------|--------------|-----------------|-----------------|
|                                 | Hom. Back.     | Less Het. Back. | More Het. Back. | Hom. Back.   | Less Het. Back. | More Het. Back. |
| Gray Level                      | 100%           | 92.2%           | 54.1%           | 99.1%        | 80.9%           | 69.0%           |
| Color                           | 100%           | 88.0%           | <b>82.3%</b>    | 98.5%        | 77.8%           | <b>73.6%</b>    |



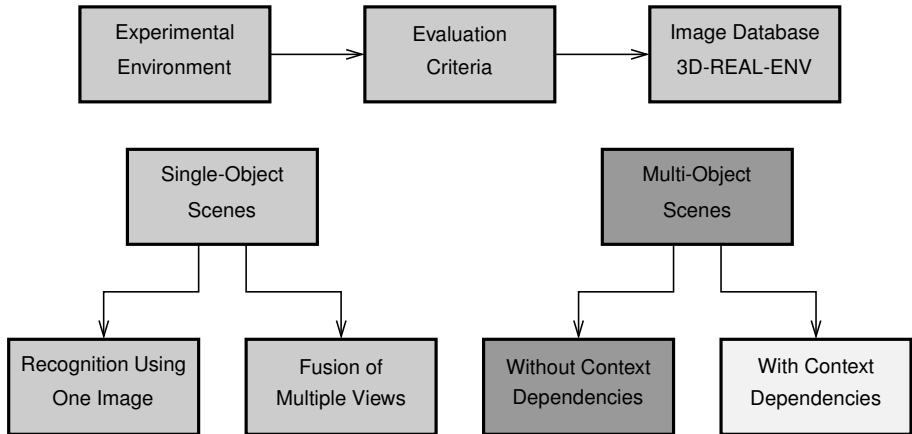


# Single-Object, Multiple Views



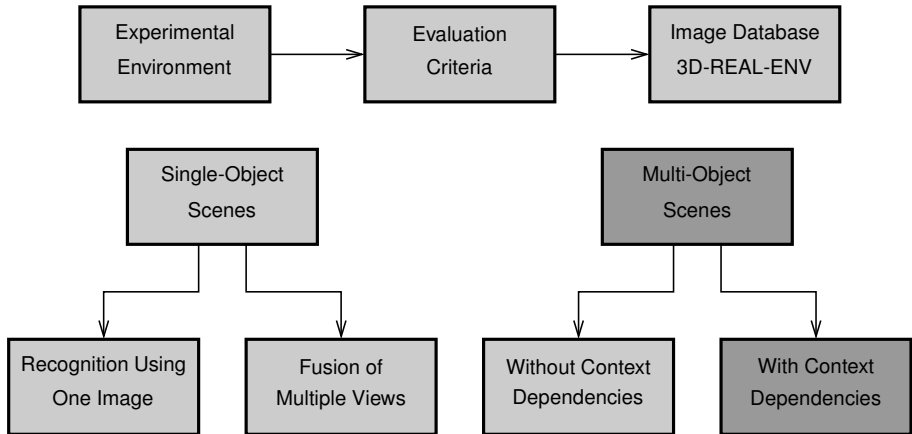


# Multi-Object, Without Context





# Multi-Object, With Context







# Test Images and Recognition Rates

1080 HomBack



1080 LessHetBack



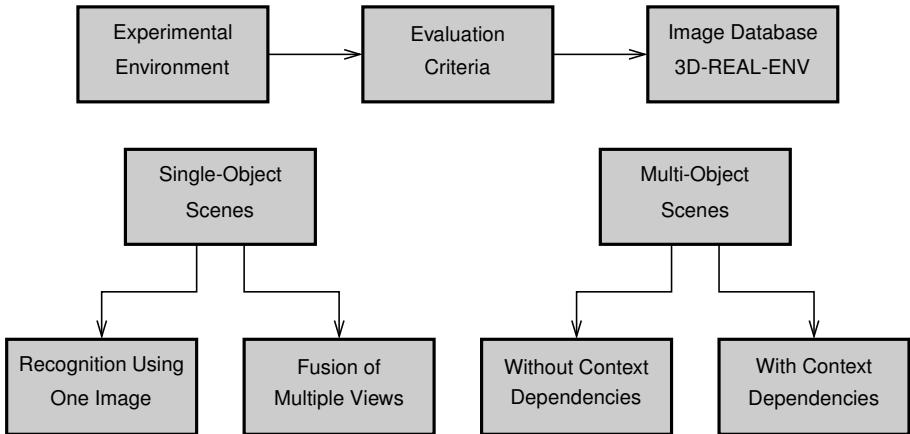
1080 MoreHetBack



| 3D-REAL-ENV<br>Image Database | Without Context Modeling |         |         | With Context Modeling |         |              |
|-------------------------------|--------------------------|---------|---------|-----------------------|---------|--------------|
|                               | Hom                      | LessHet | MoreHet | Hom                   | LessHet | MoreHet      |
| <b>Classification</b>         | 100%                     | 91.9%   | 62.9%   | 100%                  | 97.0%   | <b>87.5%</b> |
| <b>Localization</b>           | 99.7%                    | 81.7%   | 58.1%   | 99.7%                 | 81.7%   | 58.1%        |



# Evaluation Completed



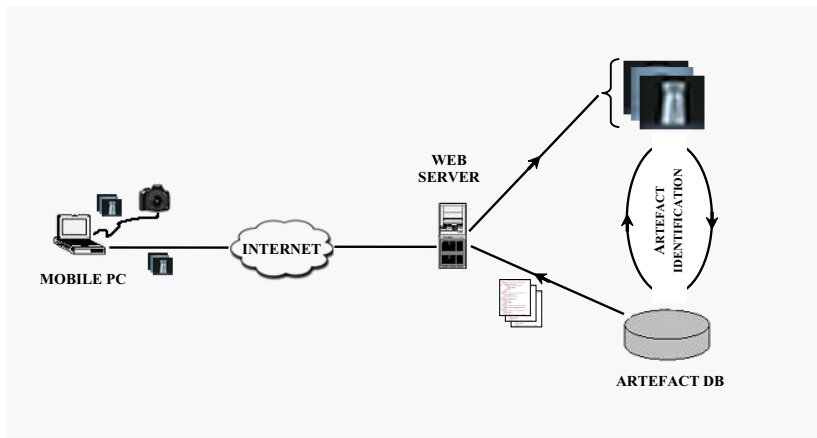


# Overview

- 1 Fundamental Concept
- 2 Statistical Modeling
- 3 Classification and Localization
- 4 Experiments and Results
- 5 System Applications**
- 6 Conclusion

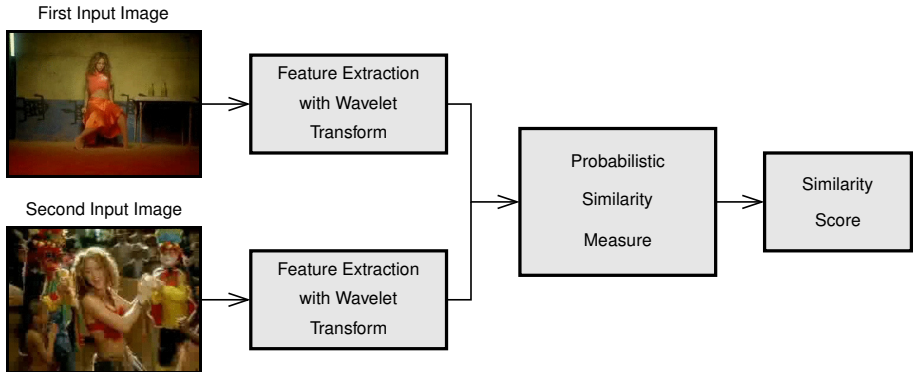


# Museum Artifact Classification





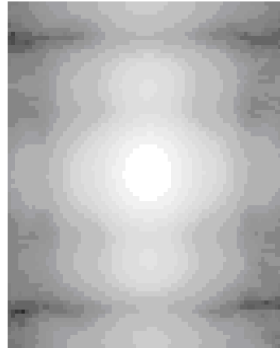
# Image Similarity Measure



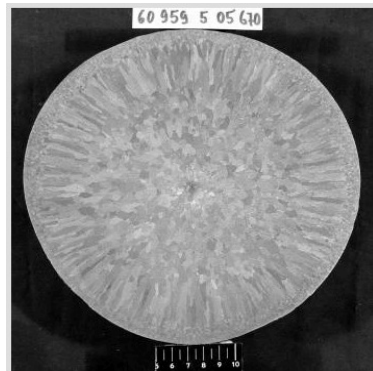
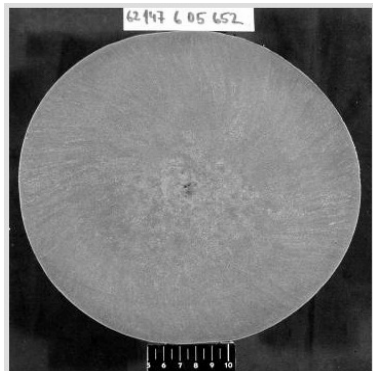
## Applications

- Similarity Measure for Antenna Patterns
- Quality Scoring of Metallography Images
- Content Management and Retrieval System

# Antenna Patterns



# Metallography Images





# Video Management and Retrieval System

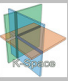
K-Space Content Management and Retrieval System - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorites

Address <http://138.37.33.138:9759/> Go Links

## K-Space Document Navigation



**Content Management System**
















Create content ...  
Input content  
Process content  
Current tasks

Content Retrieval ...  
Navigate  
Search

Auxiliary tasks ...  
User Profile  
Help


Home

Logout guest

| View  | Document title & abstract   | Transcode  | Shot edition  | Semantic annotation   | Visual search   |
|---|---|--|---|---|---|
|  | <b>shakira hips dont lie</b><br>"Hips Don't Lie" is a Grammy Award-nominated Latin pop song performed by Colombian singer Shakira and Haitian rapper Wyclef Jean. The music video was direct...                                 |  |  |  |  |
|  | <b>smokey robbins nothing compares to you</b><br>"Nothing Compares To You" is a song written around 1964 or 1965 by Prince and the New Power Generation, a funk band created as an outlet to release more of his music. Five... |  |  |  |  |
|  | <b>britney spears baby one more time</b><br>Shot at Venice High School in California, the scenario begins with Spears in a particularly boring class right before the end of the day. Her assistant Paul...                     |  |  |  |  |

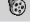









  

<http://138.37.33.138:9759/> - View Document "molocho sing it back..."



Done Internet

|   |  |  |   |   |   |
|---|--|--|---|---|---|
|  | <b>News journal4</b><br>Schroeder speaking with Romano Prodi in Berlin about the European Union Stability pact. They are standing in front of Brandenburg Gate in Berlin commenting t... |  |  |  |  |
|  | <b>News journal5</b><br>Ulrich Barths Deutsche Welle journalists speaks about mobile phones market. Business news. Ericsson  |  |  |  |  |





# Overview

- 1 Fundamental Concept
- 2 Statistical Modeling
- 3 Classification and Localization
- 4 Experiments and Results
- 5 System Applications
- 6 Conclusion**



# Final Remarks & Future Work

## Final Remarks

- Color Modeling (Classification Rate: 54.1% → 82.3%)
- Context Modeling (Classification Rate: 62.9% → 87.5%)
- 3D-REAL-ENV Image Database is Becoming Popular
- System is Applicable for Real Life Tasks

## Future Work

- Museum Artifact Classification
- Quality Scoring of Metallography Images
- Image Classification for Video Retrieval
- Appearance-Based and Shape-Based Object Recognition



# Summary

- 1. Fundamental Concept**
- 2. Statistical Modeling**
- 3. Classification and Localization**
- 4. Experiments and Results**
- 5. System Applications**
- 6. Conclusion**