

LOCUS: Learning Object Classes with Unsupervised Segmentation

16-721: Advanced Perception
Nik Melchior
April 10, 2006

Outline

- Learning Flexible Sprites in Video Layers (Jojic & Frey 2001)
- LOCUS

Motivation



Object category recognition and segmentation in cluttered scenes

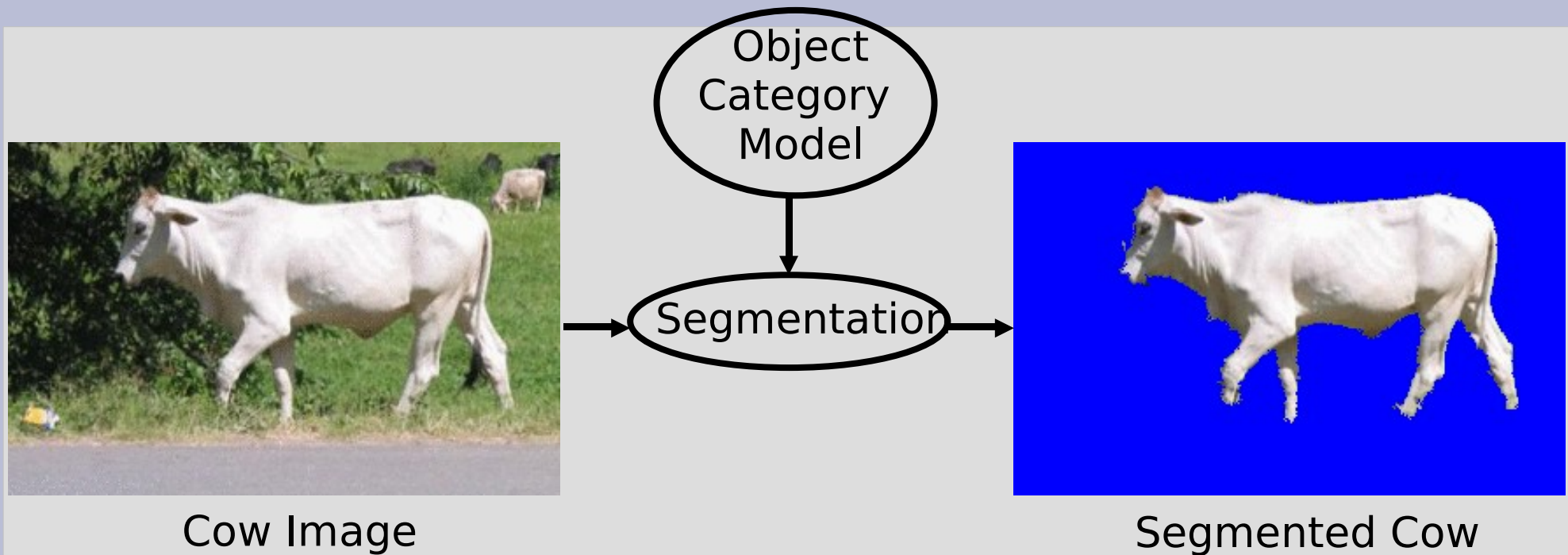
Challenges

- ✓ Spatial variability
- ✓ Texture variability
- ✓ Occlusion
- ✓ Pose
- ✓ Lighting



Motivation

- Given an image and object category, to segment the object

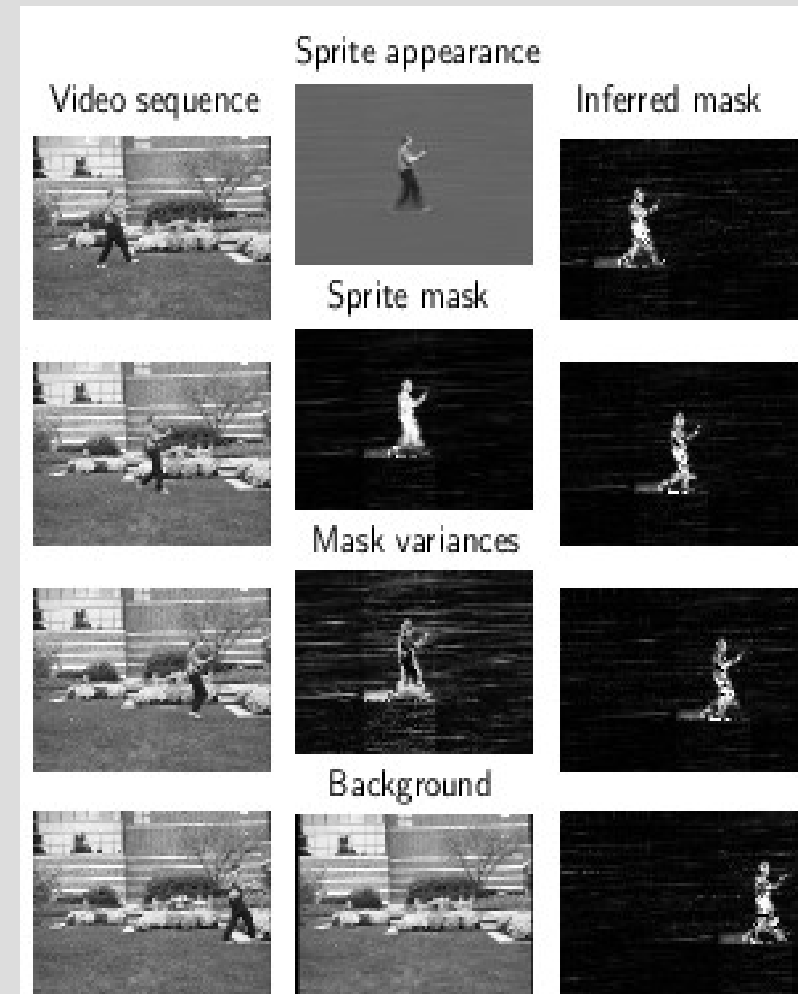


Segmentation should (ideally) be

- shaped like the object e.g. cow-like
- obtained efficiently in an unsupervised manner
- able to handle self-occlusion

Flexible Sprites

- Provides a stable segmentation of whole moving objects in video
- Each flexible sprite is a separate layer containing the transformation of a rigid sprite



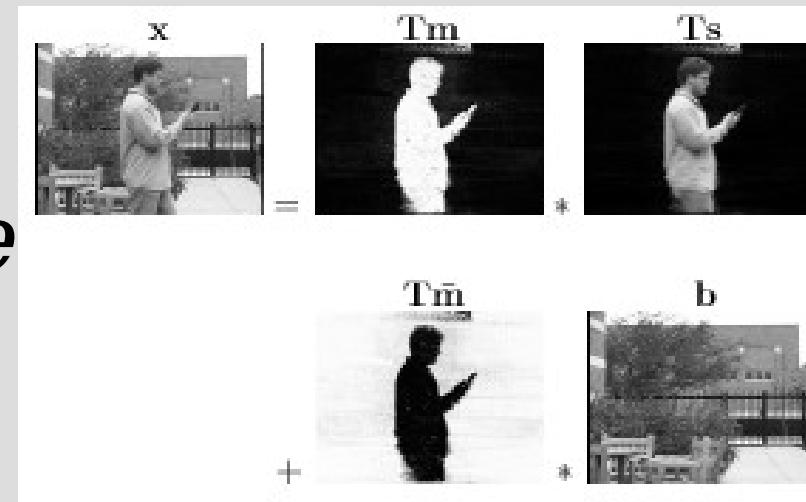
Big Picture

- Generative model
- Seeks the best explanation for the pixels of all images, constrained by number of layers
- Operates directly on pixels rather than low-level features
- Rigorous formulation similar to DDMCMC

Flexible Sprites

- User specifies number of layers
- Each layer contains:
 - Sprite appearance 's'
 - Real-valued mask 'm'
 - Discretized simple transformation 'T'
“permutation matrix that rearranges the pixels in s”

$$x = Tm * Ts + T\bar{m} * b + noise$$



Appearance, mask, and transform

Recursive formula for image appearance

Zero-mean Gaussian noise

$$p(x|\{s_l, m_l, T_l\}) = N\left(x; \sum_L \left(\prod_L T_i \bar{m}_i \right) * T_l m_l * T_l s_l, \beta\right)$$

- Assuming independence of each class, appearance, mask, and transformation

Transformation prior (uniform)

Appearance

Mask

Class prior (to be learned)

$$p(x, \{s_l, m_l, T_l\}) = N\left(x; \sum_L \left(\prod_L T_i m_i \right) * T_l s_l, \beta\right) * \prod_L \left(N(s_l; \mu_{c_l}, \phi_{c_l}) N(m_l; \eta_{c_l}, \psi_{c_l}) \pi_{c_l} \rho_{T_l} \right)$$

Formulation

- Use EM to maximize the posterior

$$p(\{c_l, T_l, s_l, m_l\} | X) = \frac{p(x, \{c_l, T_l, s_l, m_l\})}{\sum \int p(x, \{c_l, T_l, s_l, m_l\})} \\ \approx \prod_L q(c_l, T_l) q(s_l) q(m_l)$$

$$D = \sum \int \left(\prod_L q(c_l, T_l) q(s_l) q(m_l) \right) = \ln \frac{p(\{c_l, T_l, s_l, m_l\} | x)}{\prod_L q(c_l, T_l) q(s_l) q(m_l)}$$

Formulation

$$F = D + \ln p(\mathbf{x})$$
$$= \sum \int \left(\prod_L q(c_l, T_l) q(s_l) q(m_l) \right) = \ln \frac{p(\mathbf{x}, \{c_l, T_l, s_l, m_l\})}{\prod_L q(c_l, T_l) q(s_l) q(m_l)}$$

- Maximize F
 - log of numerator is a sum of Mahalanobis distances
 - $q \ln 1/q$ is the entropy of q

E step

- For each single image:
 - the mean should stay close to the overall mean, but large where image does not match background
 - variance should be high where the transformed image is similar to the background

$$\nu \leftarrow \left(\psi^{-1} + \sum_{\mathbf{T}} \xi_{\mathbf{T}} \mathbf{T}^{-1} (\beta^{-1} * (\mathbf{T}\boldsymbol{\mu} - \mathbf{b})) \right)^{-1}$$

and

$$\gamma \leftarrow \nu * \left(\psi^{-1} \boldsymbol{\eta} + \sum_{\mathbf{T}} \xi_{\mathbf{T}} \mathbf{T}^{-1} (\beta^{-1} * (\mathbf{T}\boldsymbol{\mu} - \mathbf{b}) * (\mathbf{x} - \mathbf{b})) \right),$$

M step

$$\psi \leftarrow \frac{1}{N} \sum_n \left(\nu^{(n)} + (\gamma^{(n)} - \eta)^2 \right),$$

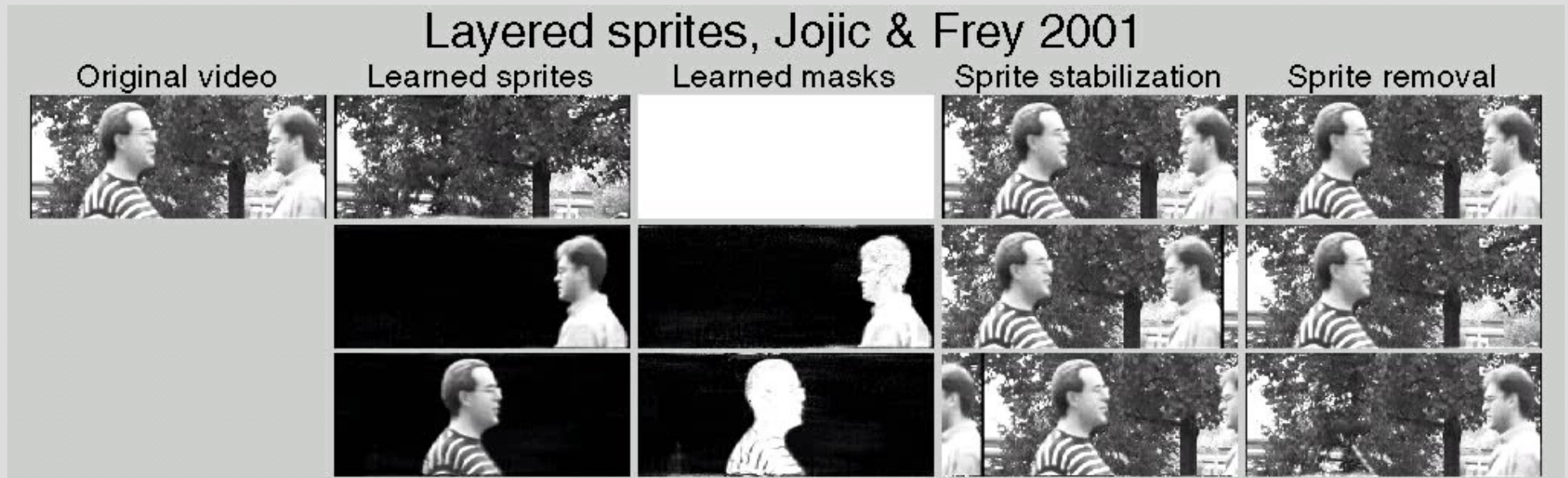
$$\eta \leftarrow \frac{1}{N} \sum_n \gamma^{(n)},$$

$$\begin{aligned} \boldsymbol{\mu} \leftarrow & \left(\sum_n \sum_{\mathbf{T}} \xi_{\mathbf{T}}^{(n)} (\gamma^{(n)2} + \nu^{(n)}) \right)^{-1} * \left(\sum_n \sum_{\mathbf{T}} \xi_{\mathbf{T}}^{(n)} \right. \\ & \left. \cdot (\gamma^{(n)} * (\mathbf{T}^{-1} \mathbf{x}^{(n)} - \bar{\gamma}^{(n)} * \mathbf{T}^{-1} \mathbf{b}) + \nu^{(n)} * \mathbf{T}^{-1} \mathbf{b}) \right). \end{aligned}$$

$$\begin{aligned} \mathbf{b} \leftarrow & \left(\sum_n \sum_{\mathbf{T}} \xi_{\mathbf{T}}^{(n)} \mathbf{T} (\bar{\gamma}^{(n)} + \nu^{(n)}) \right)^{-1} * \left(\sum_n \sum_{\mathbf{T}} \xi_{\mathbf{T}}^{(n)} \right. \\ & \left. \cdot (\mathbf{T} \bar{\gamma}^{(n)} * (\mathbf{x}^{(n)} - \mathbf{T} \gamma^{(n)} * \boldsymbol{\mu}) + \mathbf{T} \nu^{(n)} * \boldsymbol{\mu}) \right). \end{aligned}$$

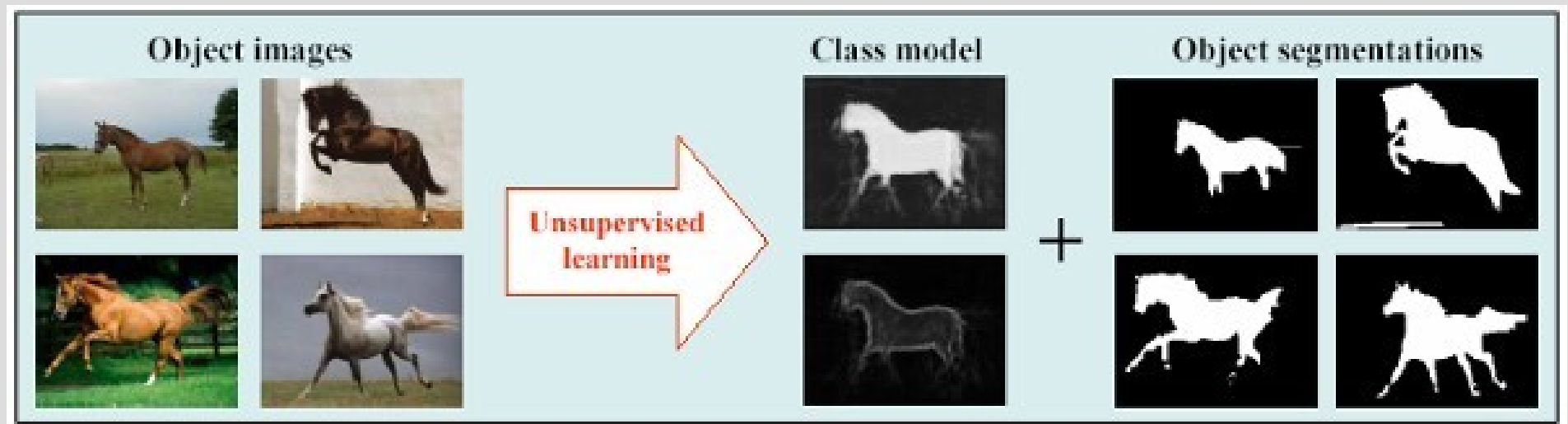
Results

- 1 fps at 320x240



LOCUS

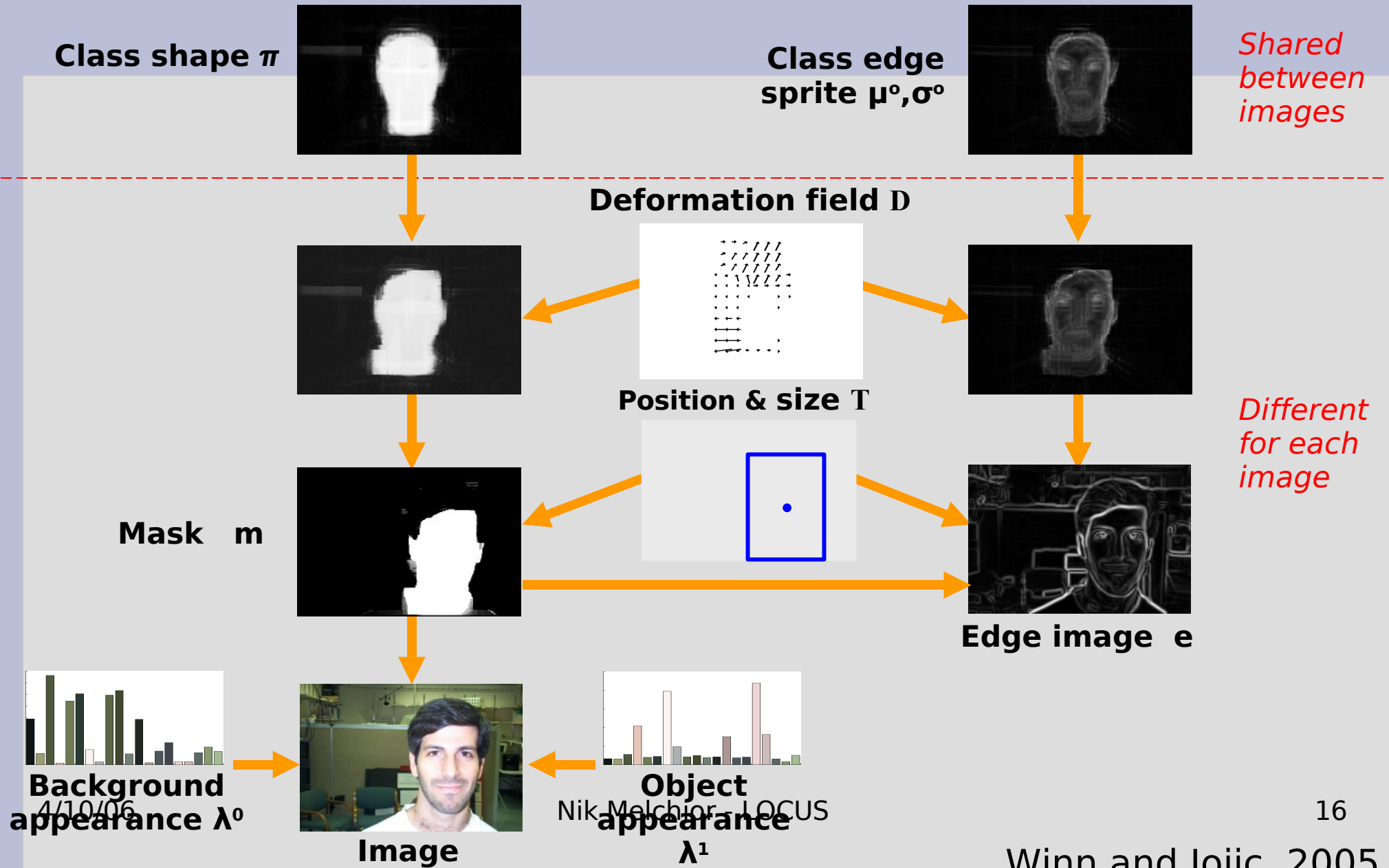
- Learning Object Classes with Unsupervised Segmentation



Comparison to Flexible Sprites

- Classes instead of single sprites
 - more complex appearance model
 - addition of edge model
- 2 layers: object and model
- weak prior of palettes (appearance)
 - one for object, one for background

LOCUS model



LOCUS model

- Deformation field \mathbf{D}
 - Markov Random Field

$$P(\mathbf{D}) = \frac{1}{Z} \exp \sum_{(i,j) \in \bar{E}} -\alpha |d_i - d_j|^2$$

- Transform \mathbf{T}
 - scaling and translation
 - rotation or full affine transform possible

- Mask \mathbf{m}

$$P_{\text{mrf}}(\mathbf{m} | \gamma, \beta) = \frac{1}{Z} \exp \sum_{(i,j) \in \bar{E}} -\delta(m_i \neq m_j) \gamma e^{-\beta \|z_i - z_j\|^2}$$

- Edge model \mathbf{e}

- Appearance model

- Gaussian mixture model (histogram of colors or textures)

$$P(\mathbf{z}_i | m_i, \boldsymbol{\theta}) = \sum_{k=1}^K \lambda_k^{m_i} \mathcal{N}(\mathbf{z}_i | \eta_k, \Sigma_k).$$

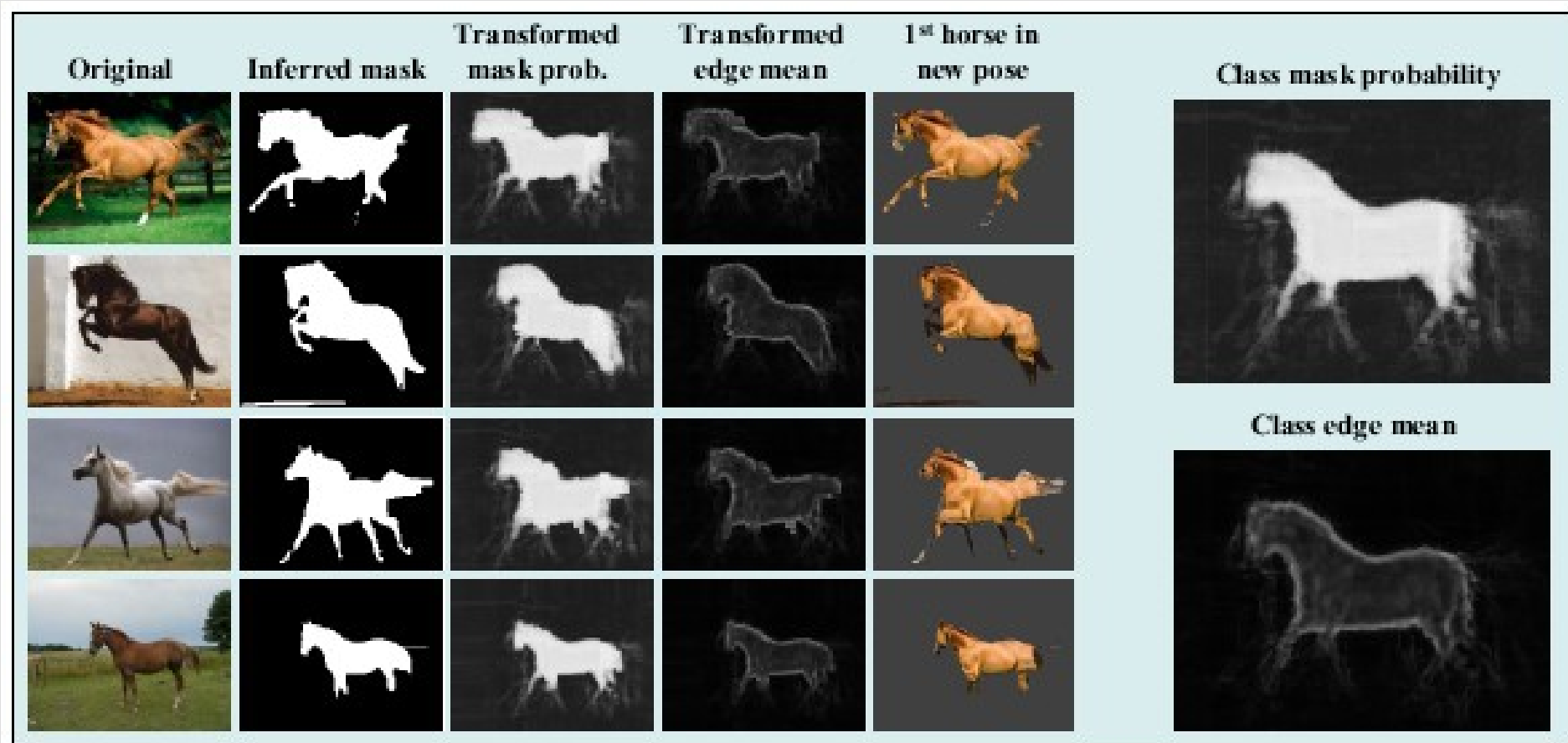
EM formulation

$$P(\pi, \mu^o, \sigma^o, \mu^b, \sigma^b | I) \approx \prod_i Q(\pi_i) Q(\mu_i^o, \sigma_i^o) Q(\mu_i^b, \sigma_i^b)$$

- Posterior similar to flexible sprites
- Iteratively optimize over 9 latent variables

$$\{\lambda, m, T, D, \pi, \mu^o, \sigma^o, \mu^b, \sigma^b\}$$

Results



Results

- Comparison to Borenstein fragments
- Percentages indicate correctly labeled pixels across all images
- Last row removes class shape and edge information

	Segmentation accuracy	
	Horses	Cars (side)
LOCUS (color)	93.1%	91.4%
LOCUS (texture)	93.0%	94.0%
unannotated training images		
Borenstein et al	93.6%	-
hand-segmented training images		
LOCUS - no class model	88.6%	82.1%

Results

