# Matching single objects...
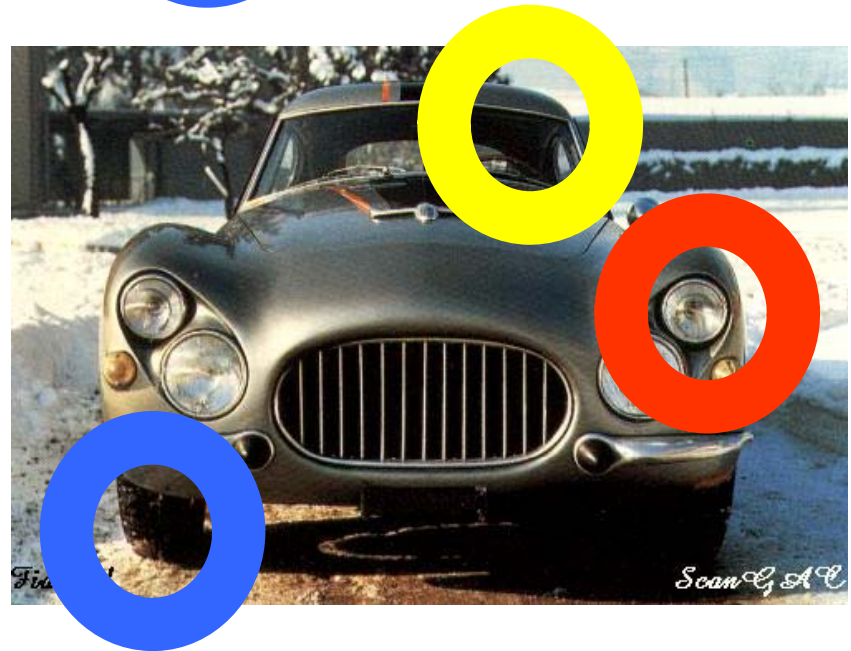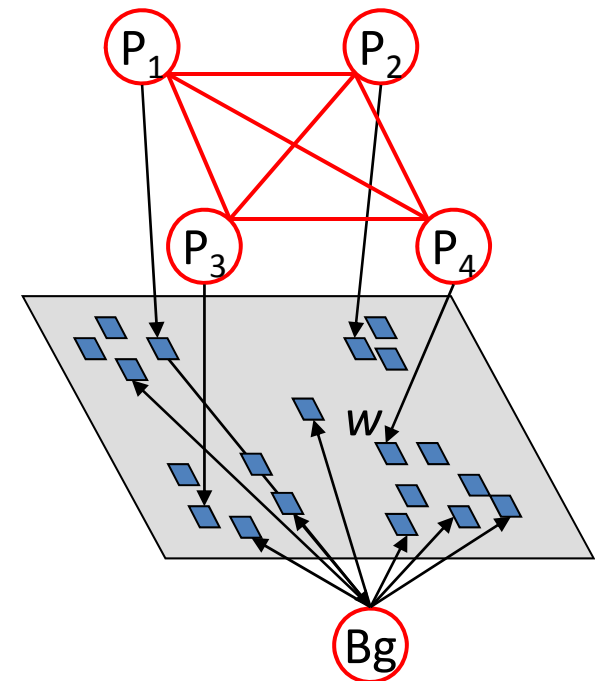
SIFT features

Image gradients

Keypoint descriptor

Lowe, 1999
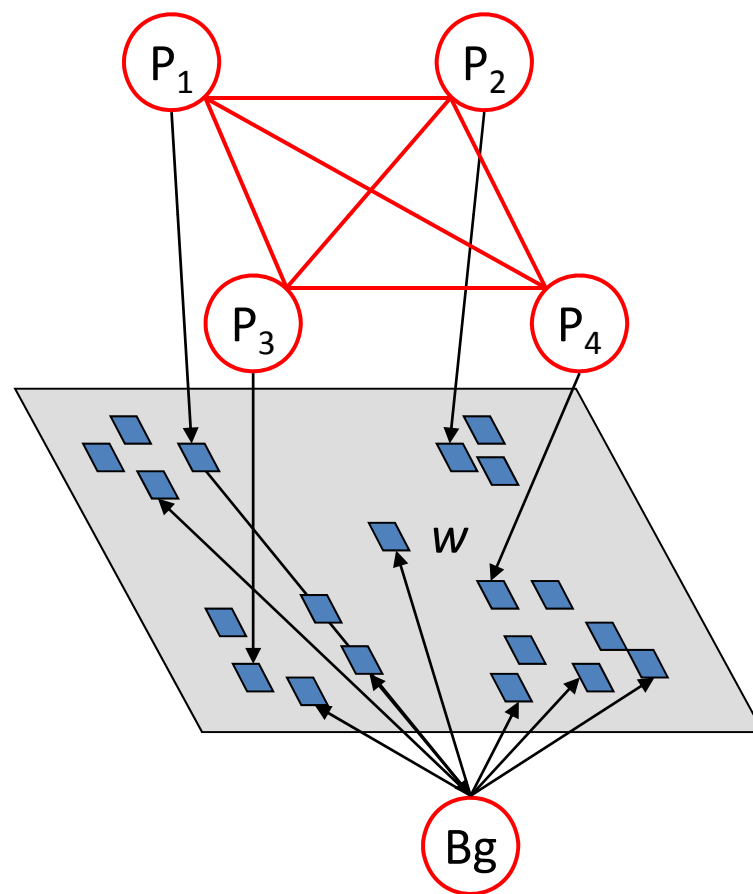
# Matching a class of objects...
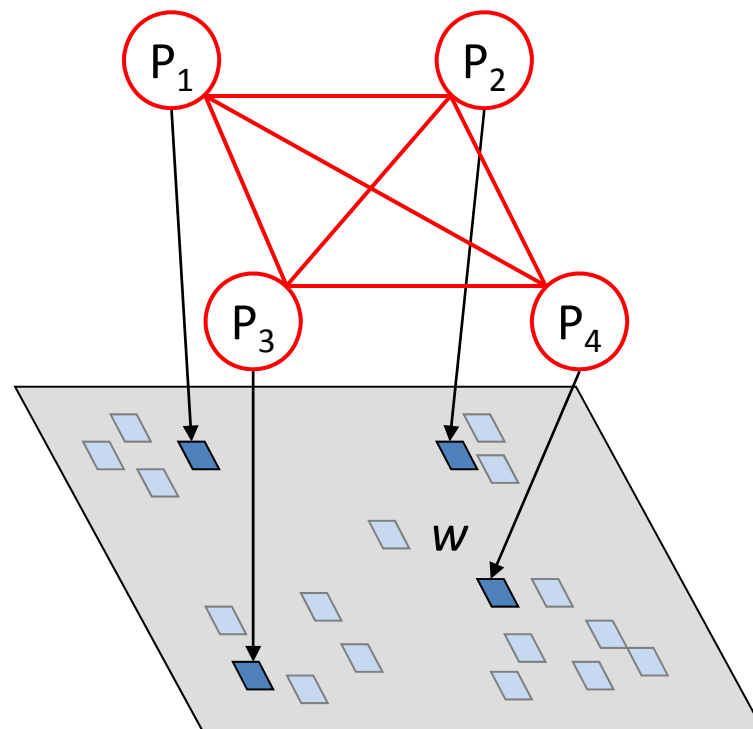
# Part-based representation: constellation model

- Fischler & Elschlager 1973

- Yuille '91
- Brunelli & Poggio '93
- Lades, v.d. Malsburg et al. '93
- Cootes, Lanitis, Taylor et al. '95
- Amit & Geman '95, '99
- Perona et al. '95, '96, '98, '00, '03, '04, '05
- Felzenszwalb & Huttenlocher '00, '04
- Crandall & Huttenlocher '05, '06
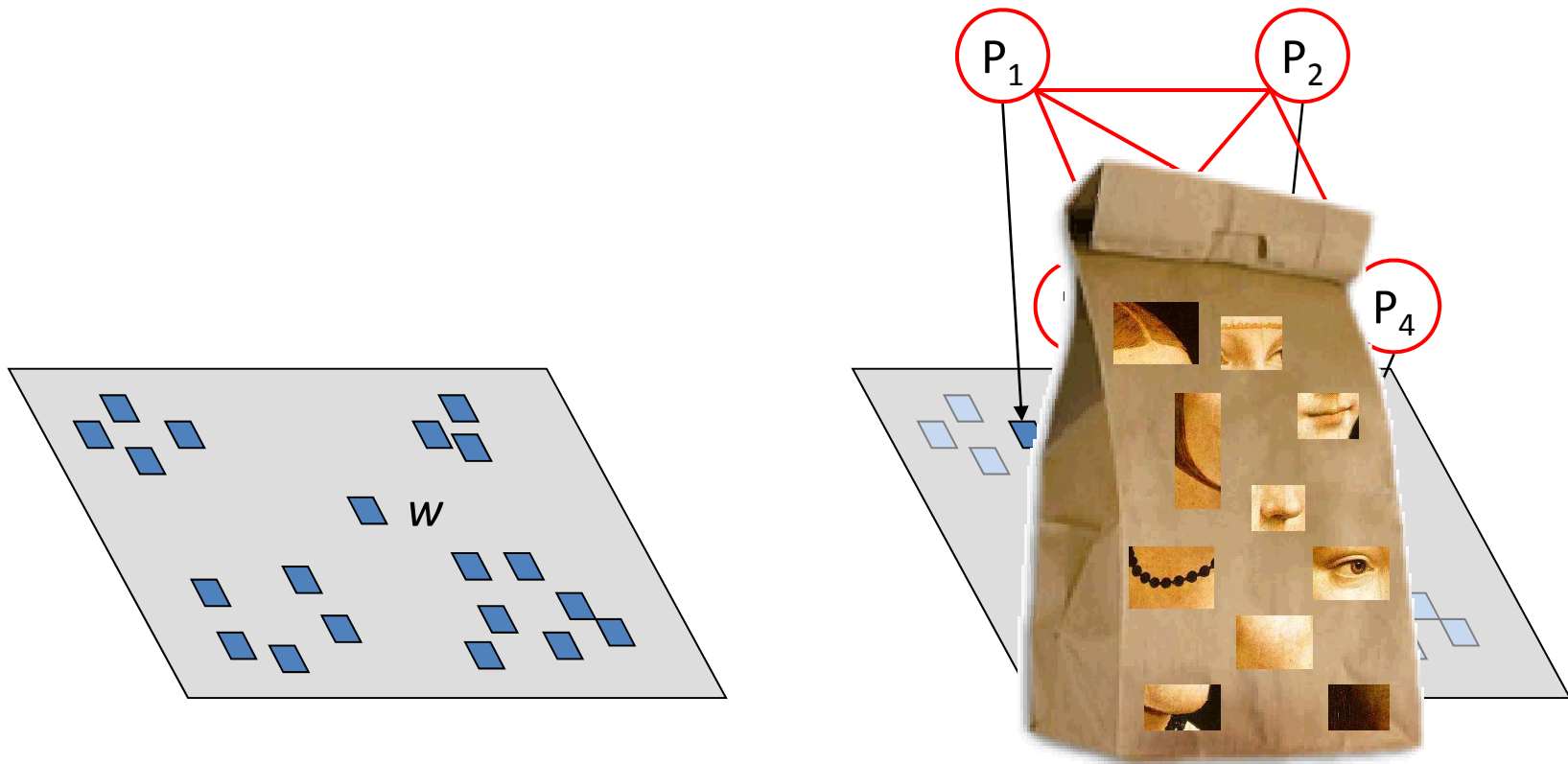- Leibe & Schiele '03, '04
- Many papers since '00

# Part-based representation: constellation model
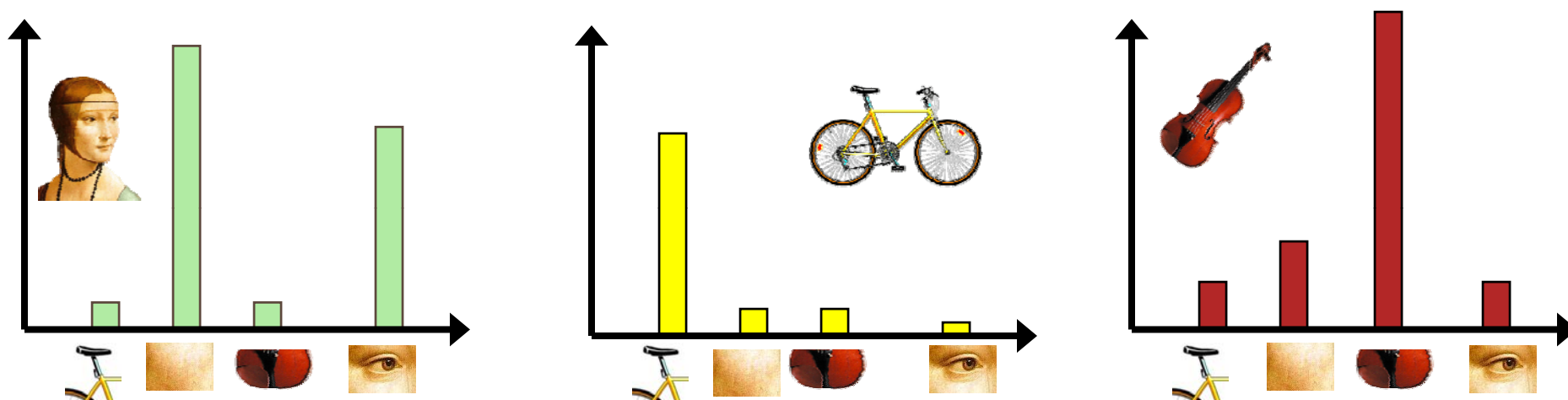
# Part-based representation: constellation model
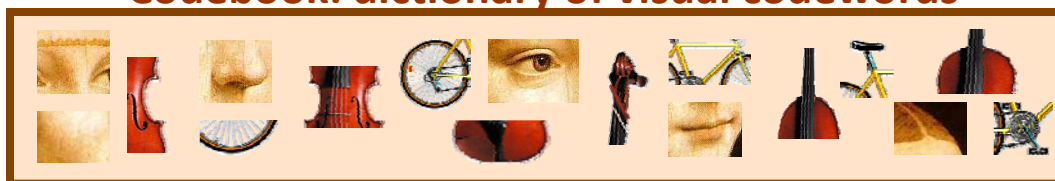
# "Bag of words" representation



Csurka et al. 2004; earlier work in texture: Leung & Malik, 1999

# "Bag of words" representation



**Codebook: dictionary of visual codewords**

Csurka et al. 2004; earlier work in texture: Leung & Malik, 1999
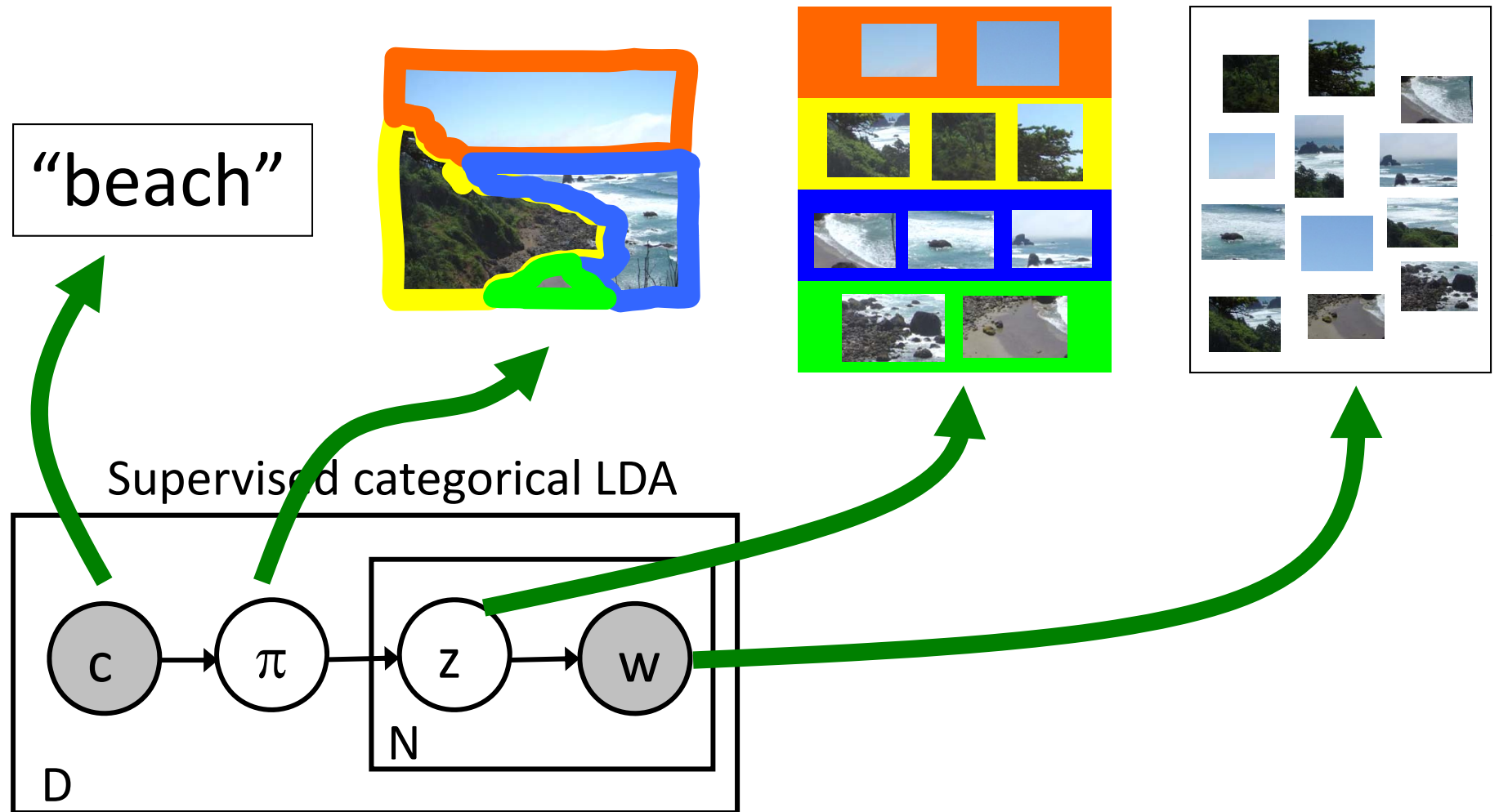
# Analogy to textual documents

Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach the brain from our eyes. For a long time it was thought that the retinal image was transmitted point by point to visual centers in the brain; the cerebral cortex was a movie screen, so to speak, upon which the image in the eye was projected. Through the discoveries of Hubel and Wiesel we now know that behind the origin of the visual perception in the brain there is a considerably more complicated course of events. By following the visual impulses along their path to the various cell layers of the optical cortex, Hubel and Wiesel have been able to demonstrate that the *message about the image falling on the retina undergoes a step-wise analysis in a system of nerve cells stored in columns. In this system each cell has its specific function and is responsible for a specific detail in the pattern of the retinal image.*

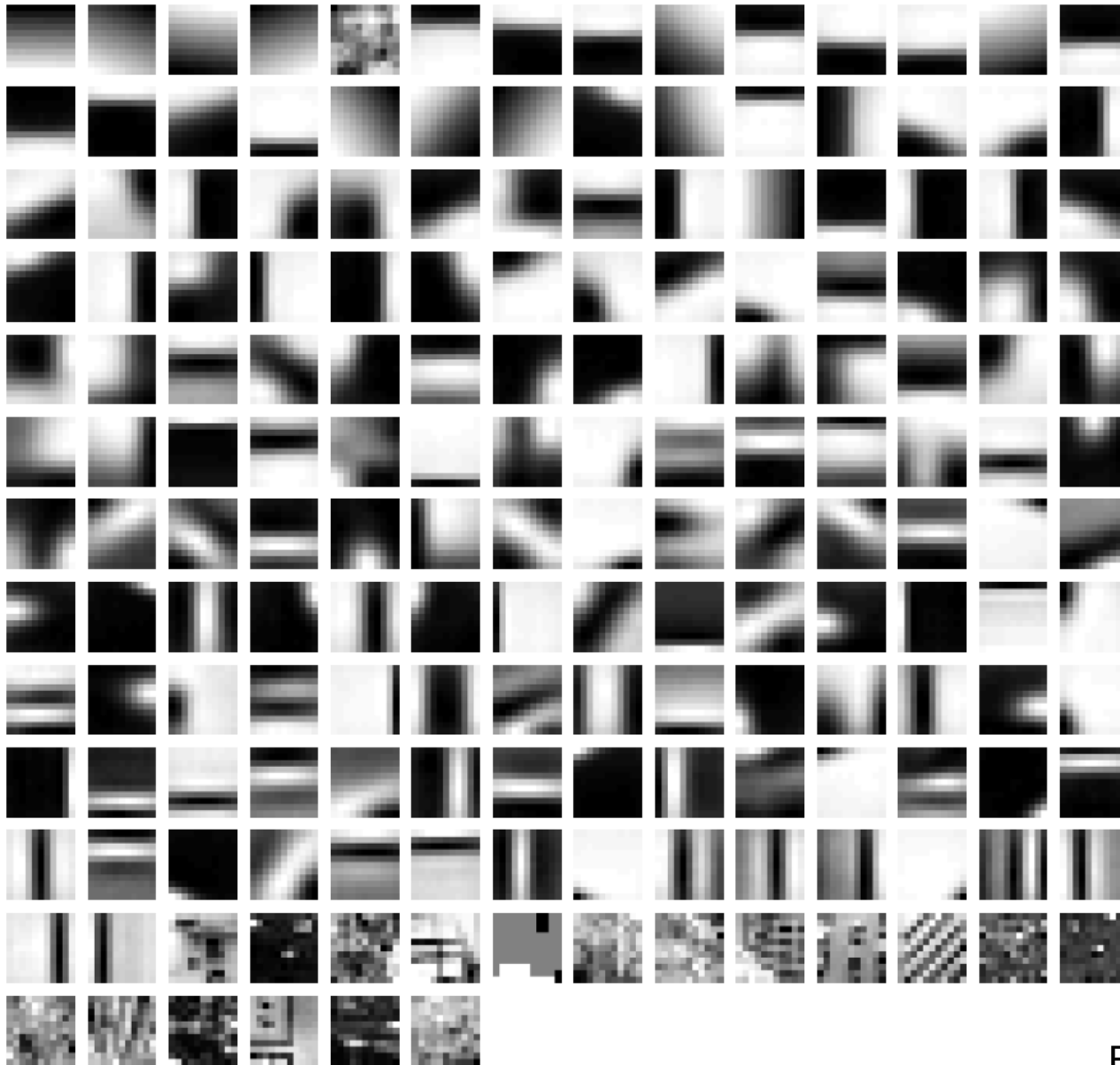**sensory, brain, visual, perception, retinal, cerebral cortex, eye, cell, optical nerve, image Hubel, Wiesel**

China is forecasting a trade surplus of $90bn (£51bn) to $100bn this year, a threefold increase on 2004's $32bn. The Commerce Ministry said the surplus would be created by a predicted 30% jump in exports to $750bn, compared with a 18% rise in imports to $660bn. The figures are likely to further annoy the US, which has long argued that China's exports are unfairly helped by a deliberately undervalued yuan. Beijing agrees the surplus is too high, but says the yuan is only one factor. Bank of China governor Zhou Xiaochuan said the country also needed to do more to boost domestic demand so more goods stay within the country. China increased the value of the yuan against the dollar by 2.1% in July and permitted it to trade within a narrow band, but the US wants the yuan to be allowed to trade freely. However, Beijing has made it clear that it will take its time and tread carefully before allowing the yuan to rise further in value.

**China, trade, surplus, commerce, exports, imports, US, yuan, bank, domestic, foreign, increase, trade, value**
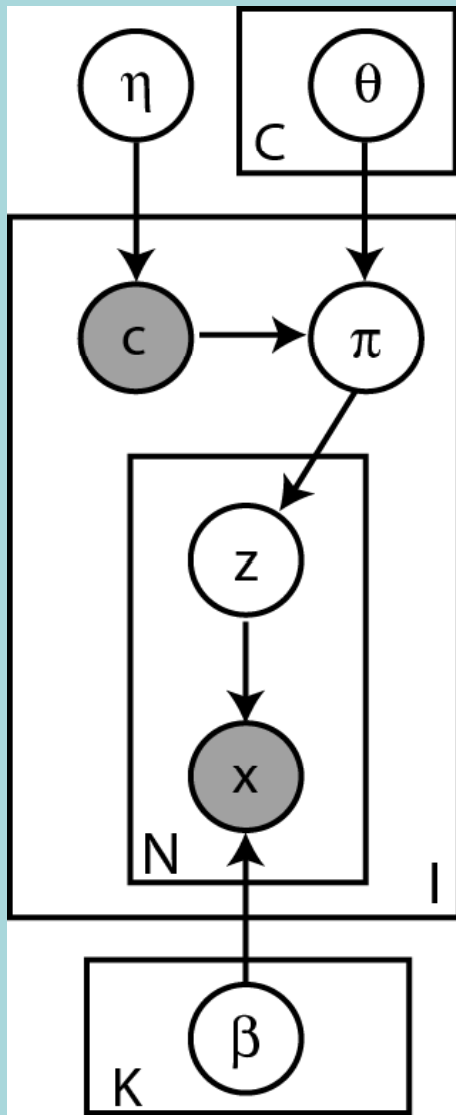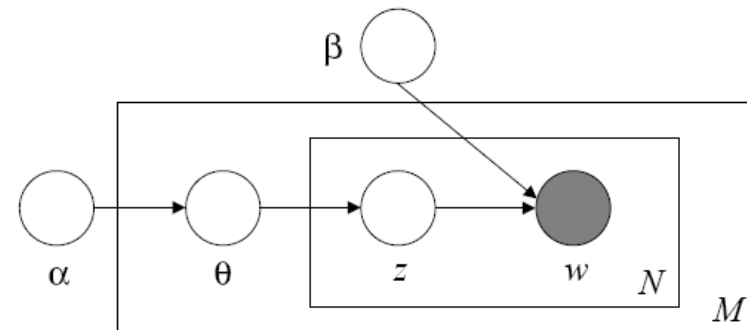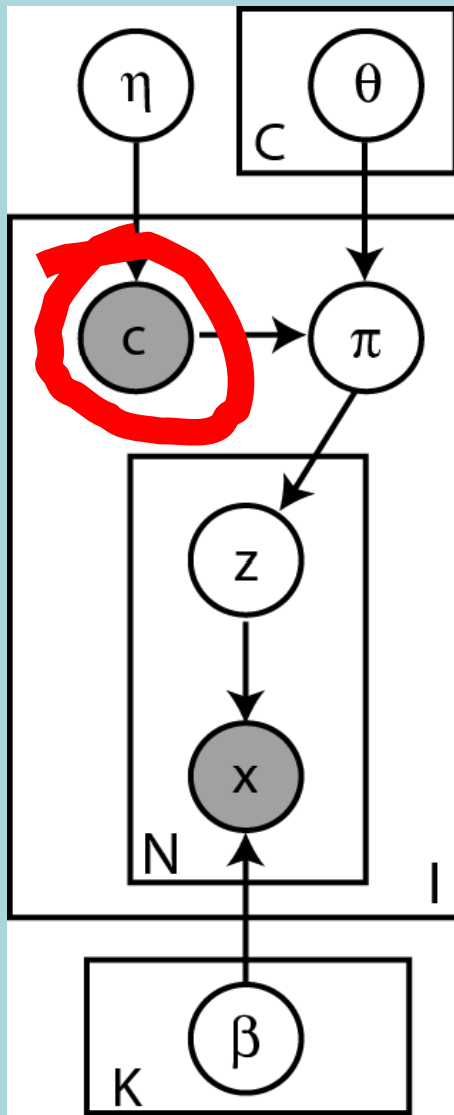
# Natural scene categorization



"beach"

Supervised categorical LDA

Fei-Fei et al. CVPR 2005

# codebook

# A Generative Model



*LDA: Blei, Ng, & Jordan. 2003*

**A Generative Model**

**scene category**

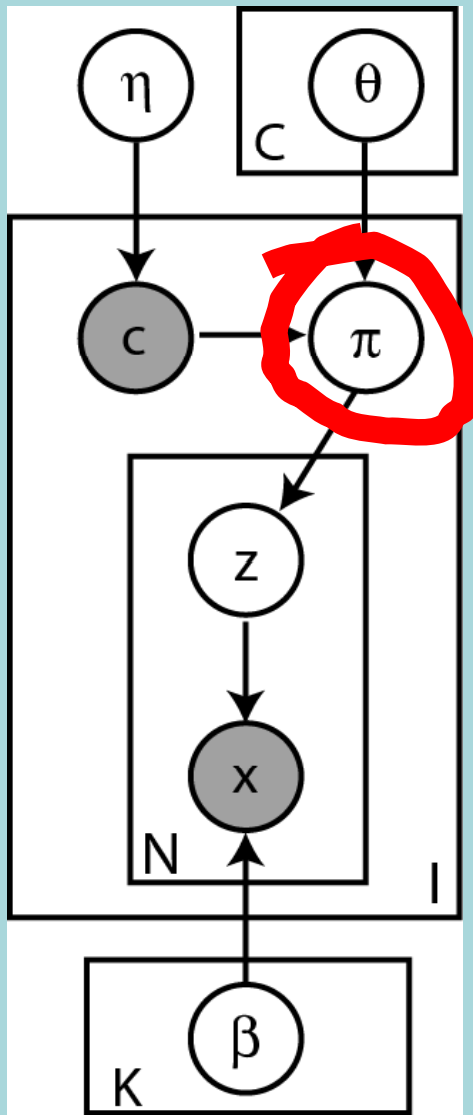discrete variable: $c \sim p(c|\eta)$

$p(c|\eta)$

forest    coast    kitchen    mountain    $c$

*Fei-Fei & Perona (CVPR 2005)*

**A Generative Model**

**mixing parameter for the latent topics**
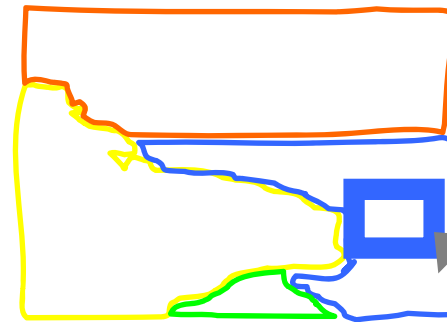
$$\pi \sim p(\pi | c, \theta)$$
$$\sim Dir(\pi | c, \theta)$$

where $\sum_{k=1}^{K} \pi_k = 1$

topic #13

topic #15

proportion of themes

topics

**A Generative Model**

**topic label**

discrete variable:

$$z \sim p(z|\pi)$$
$$\sim Mult\ (z|\pi)$$

topic #13

topic #15

proportion of themes

topics

**A Generative Model**

# patch label

discrete variable:

$$x \sim p(x \mid z, \beta)$$
$$\sim Mult(x \mid z, \beta)$$

**expected value of $\beta$ given 'z=13'**

codeword #265

**codewords**

## A Generative Model

# patch label

discrete variable:

$$x \sim p(x|z, \beta)$$
$$\sim Mult(x|z, \beta)$$

**expected value of $\beta$ given 'z=15'**

codewords

**A Generative Model**



# learning

Find the 'best' θ and β

**joint probability**

$$p(x, z, \pi | \theta, \beta, c) = p(\pi | c, \theta) \prod_n^N p(z_n | \pi) p(x_n | z_n, \beta)$$

$$p(x | \theta, \beta, c) = \int p(\pi | c, \theta) \left[ \prod_n^N \sum_{z_n} p(z_n | \pi) p(x_n | z_n, \beta) \right] d\pi$$

- exact inference is intractable
- use Variational Inference

*Fei-Fei & Perona (CVPR 2005)*

coast

forest

mountain

open
country

# model distance based on topic distribution



Fei-Fei & Perona (CVPR 2005)

Li, Socher, & Fei-Fei, *CVPR,* 2009

# Towards total scene understanding



Li, Socher, & Fei-Fei, *CVPR,* 2009

# A joint model for image classification, annotation & segmentation



class: Polo

$$p(C, \mathbf{O}, \mathbf{R}, \mathbf{X}, \mathbf{S}, \mathbf{T}, \mathbf{Z} | \eta, \alpha, \beta, \gamma, \theta, \varphi) = p(C) \cdot$$

| **Visual Component** |
|---|
| **Text Component** |

# A joint model for image classification, annotation & segmentation



class: Polo

**Visual**

C

O

N

$$p(C, \mathbf{O}, \mathbf{R}, \mathbf{X}, \mathbf{S}, \mathbf{T}, \mathbf{Z} | \eta, \alpha, \beta, \gamma, \theta, \varphi) = p(C) \cdot (\prod_{n=1}^{N_r} p(O_n | \eta, C))$$

**Text Component**

# A joint model for image classification, annotation & segmentation



class: Polo

**Visual**

C

O

R

**Color Location Texture Shape**

$$p(C, \mathbf{O}, \mathbf{R}, \mathbf{X}, \mathbf{S}, \mathbf{T}, \mathbf{Z}|\eta, \alpha, \beta, \gamma, \theta, \varphi) = p(C) \cdot (\prod_{n=1}^{N_r} p(O_n|\eta, C)) \prod_{n=1}^{N_r} ((\prod_{i=1}^{N_F} p(R_{ni}|O_n, \alpha_i))$$

**Text Component**

# A joint model for image classification, annotation & segmentation



class: Polo

**Visual**

C

O

R    X

N

$$p(C, \mathbf{O}, \mathbf{R}, \mathbf{X}, \mathbf{S}, \mathbf{T}, \mathbf{Z}|\eta, \alpha, \beta, \gamma, \theta, \varphi) = p(C) \cdot (\prod_{n=1}^{N_r} p(O_n|\eta, C)) \prod_{n=1}^{N_r} ((\prod_{i=1}^{N_F} p(R_{ni}|O_n, \alpha_i)) \cdot \prod_{r=1}^{A_r} p(X_{nr}|O_n, \beta)$$

**Text Component**

# A joint model for image classification, annotation & segmentation



class: Polo

**Visual**

**Text**

C

O

T

R

X

N

Athlete
Horse
Grass
Trees
Sky
Saddle

$$p(C, \mathbf{O}, \mathbf{R}, \mathbf{X}, \mathbf{S}, \mathbf{T}, \mathbf{Z}|\eta, \alpha, \beta, \gamma, \theta, \varphi) = p(C) \cdot (\prod_{n=1}^{N_r} p(O_n|\eta, C)) \prod_{n=1}^{N_r} ((\prod_{i=1}^{N_F} p(R_{ni}|O_n, \alpha_i)) \cdot \prod_{r=1}^{A_r} p(X_{nr}|O_n, \beta)$$

$$\cdot \qquad p(T_m|O_{Z_m}, S_m, \theta, C, \varphi)$$

# A joint model for image classification, annotation & segmentation



class: Polo

**Visual**

**Text**

C

O

T

R

X

Z

N

Athlete
Horse
Grass
Trees
Sky
Saddle
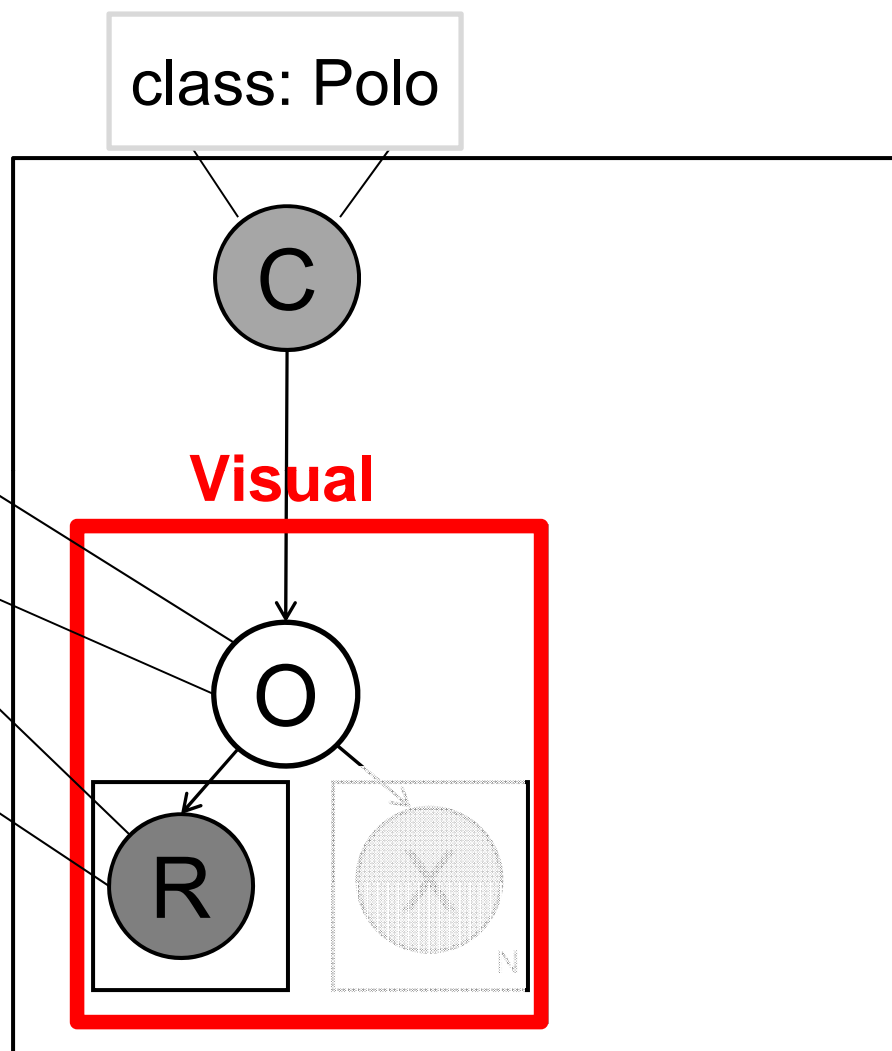
$$p(C, \mathbf{O}, \mathbf{R}, \mathbf{X}, \mathbf{S}, \mathbf{T}, \mathbf{Z}|\eta, \alpha, \beta, \gamma, \theta, \varphi) = p(C) \cdot (\prod_{n=1}^{N_r} p(O_n|\eta, C)) \prod_{n=1}^{N_r} ((\prod_{i=1}^{N_F} p(R_{ni}|O_n, \alpha_i)) \cdot \prod_{r=1}^{A_r} p(X_{nr}|O_n, \beta))$$

$$\cdot \prod_{m=1}^{N_t} p(Z_m|N_r) \qquad p(T_m|O_{Z_m}, S_m, \theta, C, \varphi)$$

# A joint model for image classification, annotation & segmentation



class: Polo

**Visual**

**Text**

C

O

T

R    X

Z

N

Athlete
Horse
Grass
Trees
Sky
Saddle
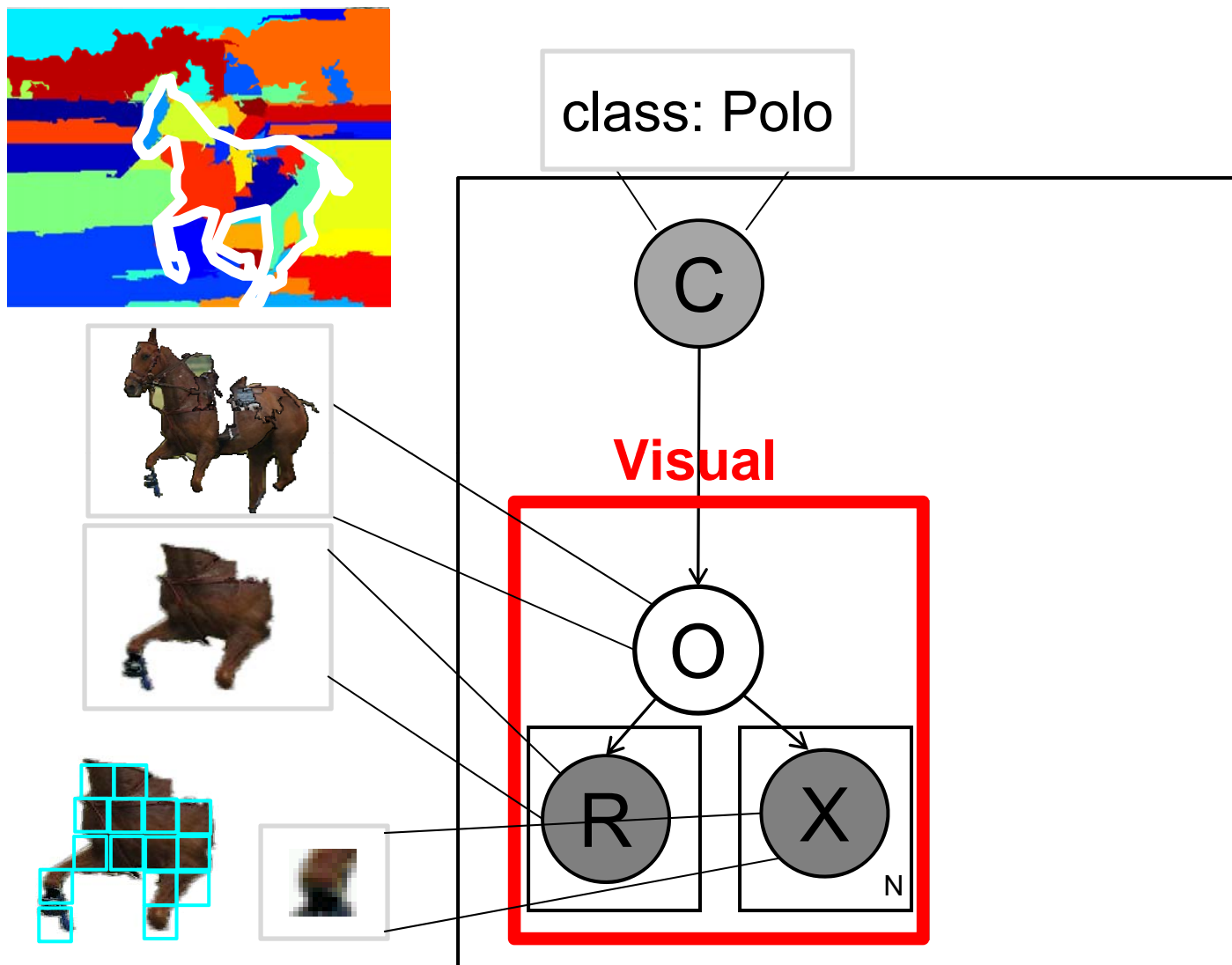
$$p(C, \mathbf{O}, \mathbf{R}, \mathbf{X}, \mathbf{S}, \mathbf{T}, \mathbf{Z} | \eta, \alpha, \beta, \gamma, \theta, \varphi) = p(C) \cdot \left( \prod_{n=1}^{N_r} p(O_n | \eta, C) \right) \prod_{n=1}^{N_r} \left( \left( \prod_{i=1}^{N_F} p(R_{ni} | O_n, \alpha_i) \right) \cdot \prod_{r=1}^{A_r} p(X_{nr} | O_n, \beta) \right)$$
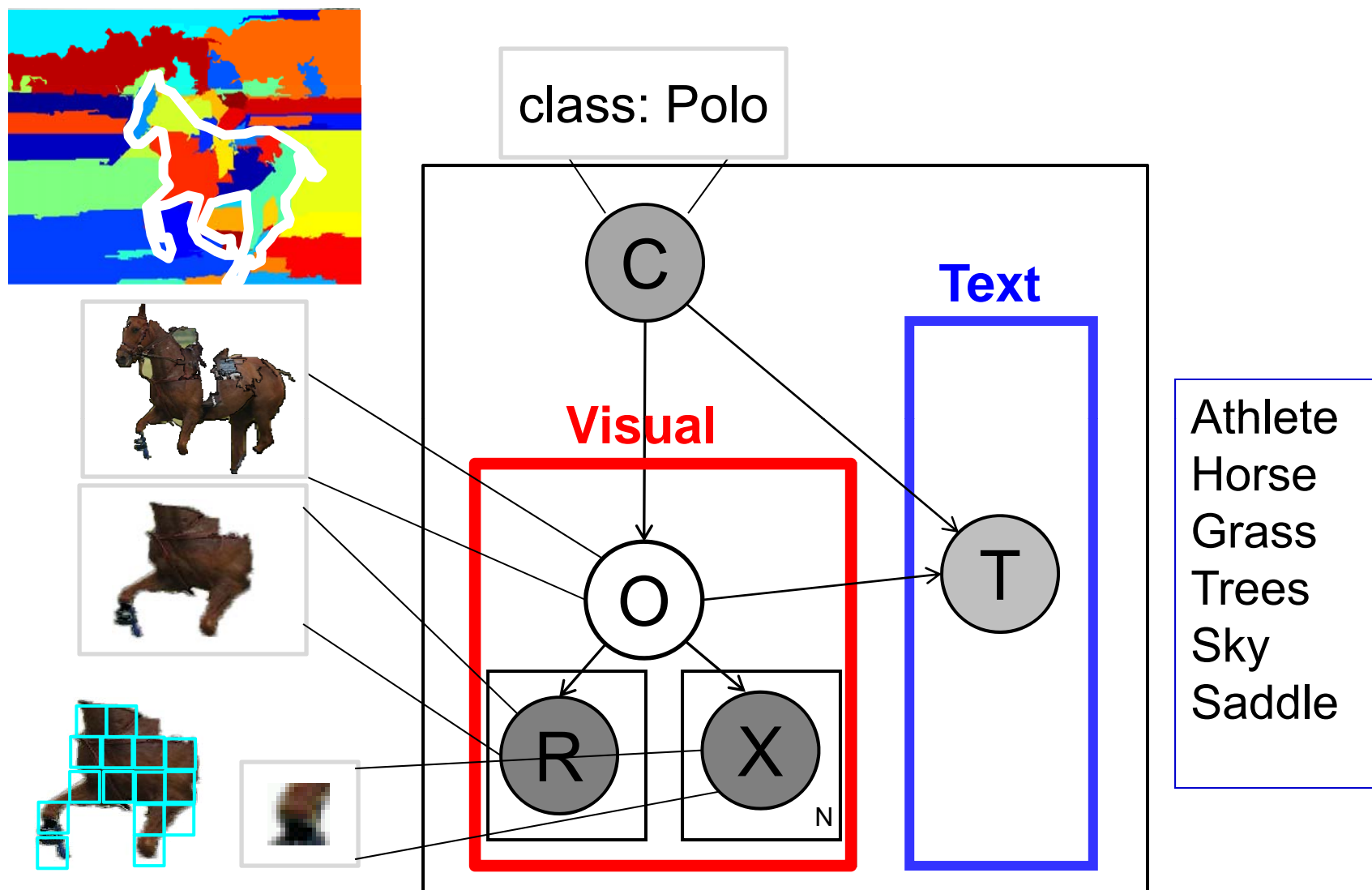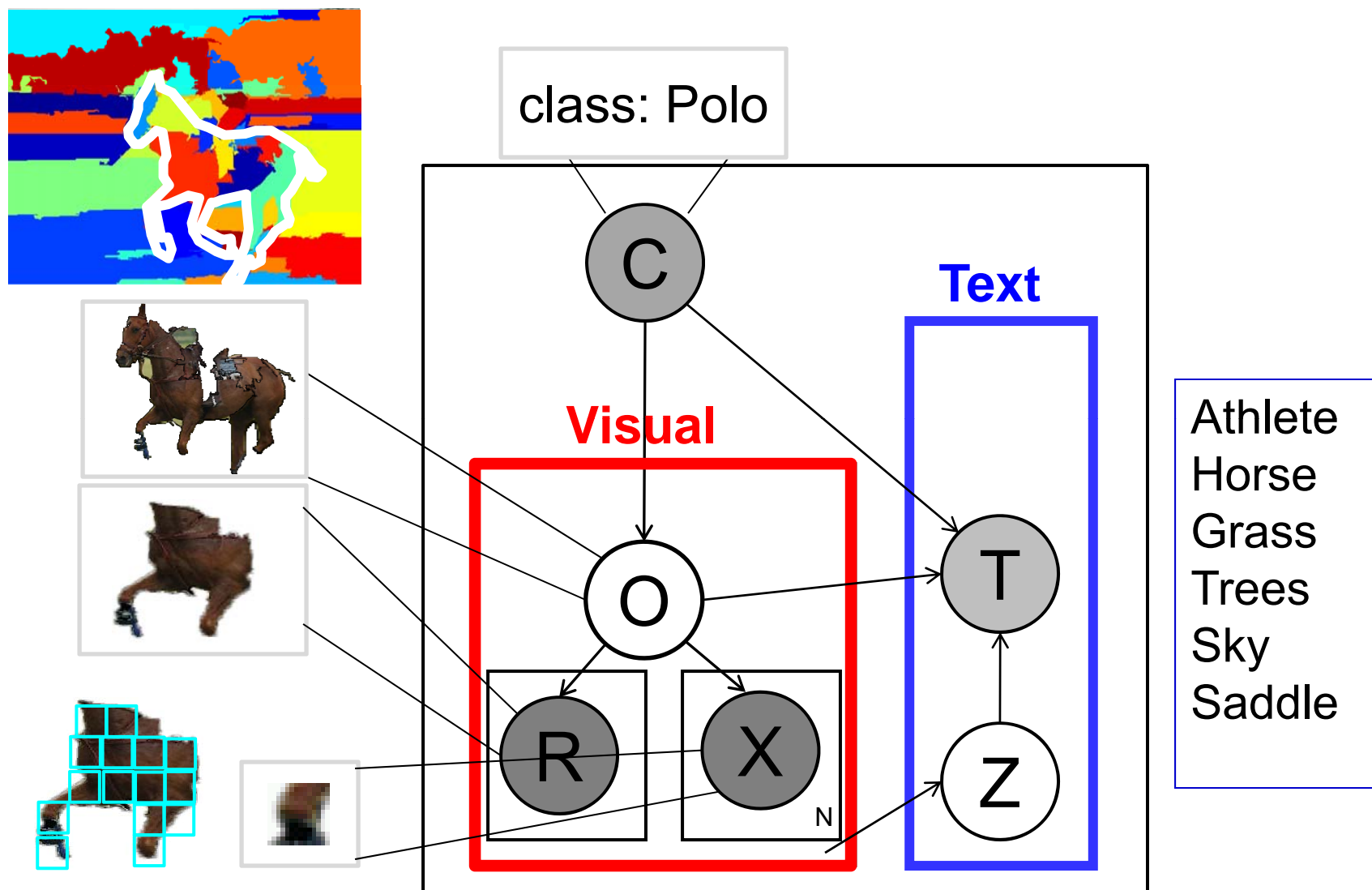
$$\cdot \prod_{m=1}^{N_t} p(Z_m | N_r) \qquad p(T_m | O_{Z_m}, S_m, \theta, C, \varphi)$$

# A joint model for image classification, annotation & segmentation



class: Polo

**Visual**

**Text**

"Switch variable"

○—< Visible

Not visible

Athlete
Horse
Grass
Trees
Sky
Saddle
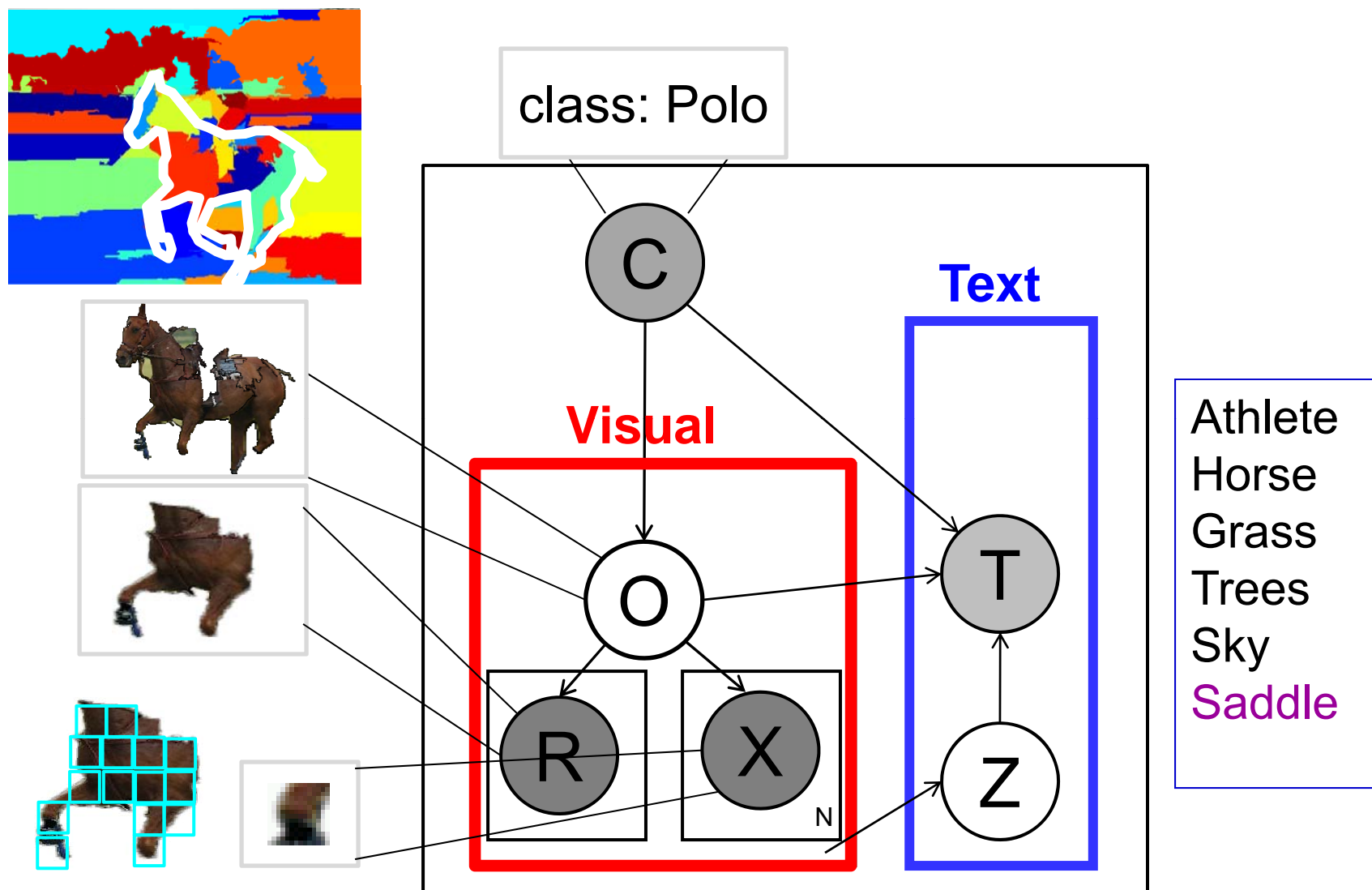
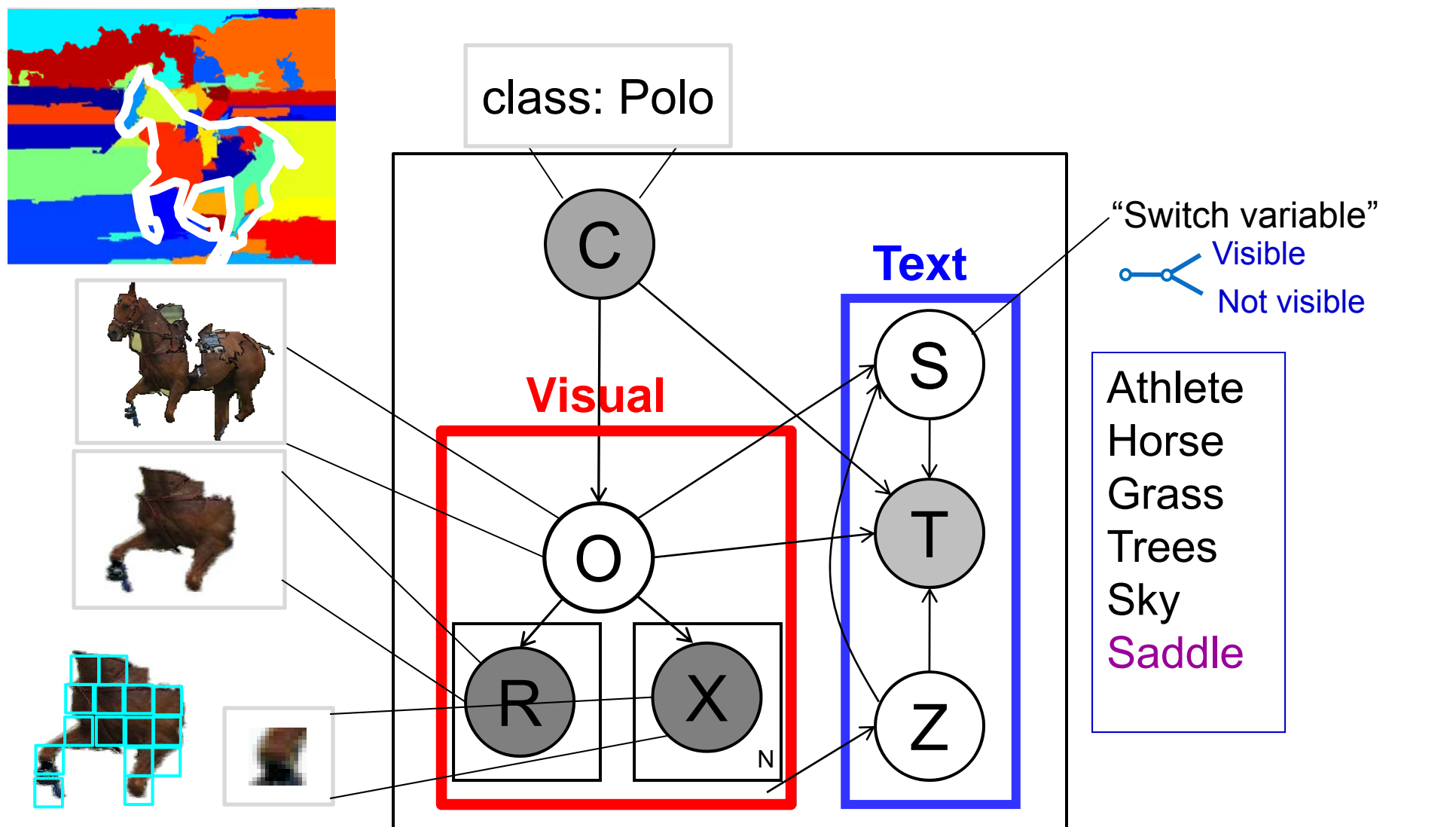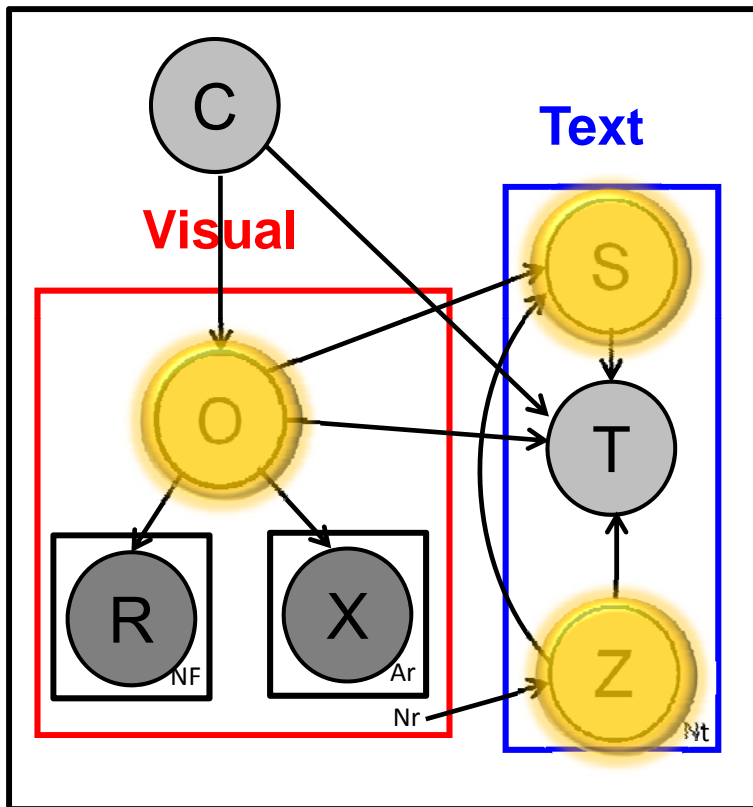$$p(C, \mathbf{O}, \mathbf{R}, \mathbf{X}, \mathbf{S}, \mathbf{T}, \mathbf{Z} | \eta, \alpha, \beta, \gamma, \theta, \varphi) = p(C) \cdot \left( \prod_{n=1}^{N_r} p(O_n | \eta, C) \right) \prod_{n=1}^{N_r} \left( \left( \prod_{i=1}^{N_F} p(R_{ni} | O_n, \alpha_i) \right) \cdot \prod_{r=1}^{A_r} p(X_{nr} | O_n, \beta) \right)$$

$$\cdot \prod_{m=1}^{N_t} p(Z_m | N_r) \quad p(S_m | O_{Z_m}, \gamma) \quad p(T_m | O_{Z_m}, S_m, \theta, C, \varphi)$$
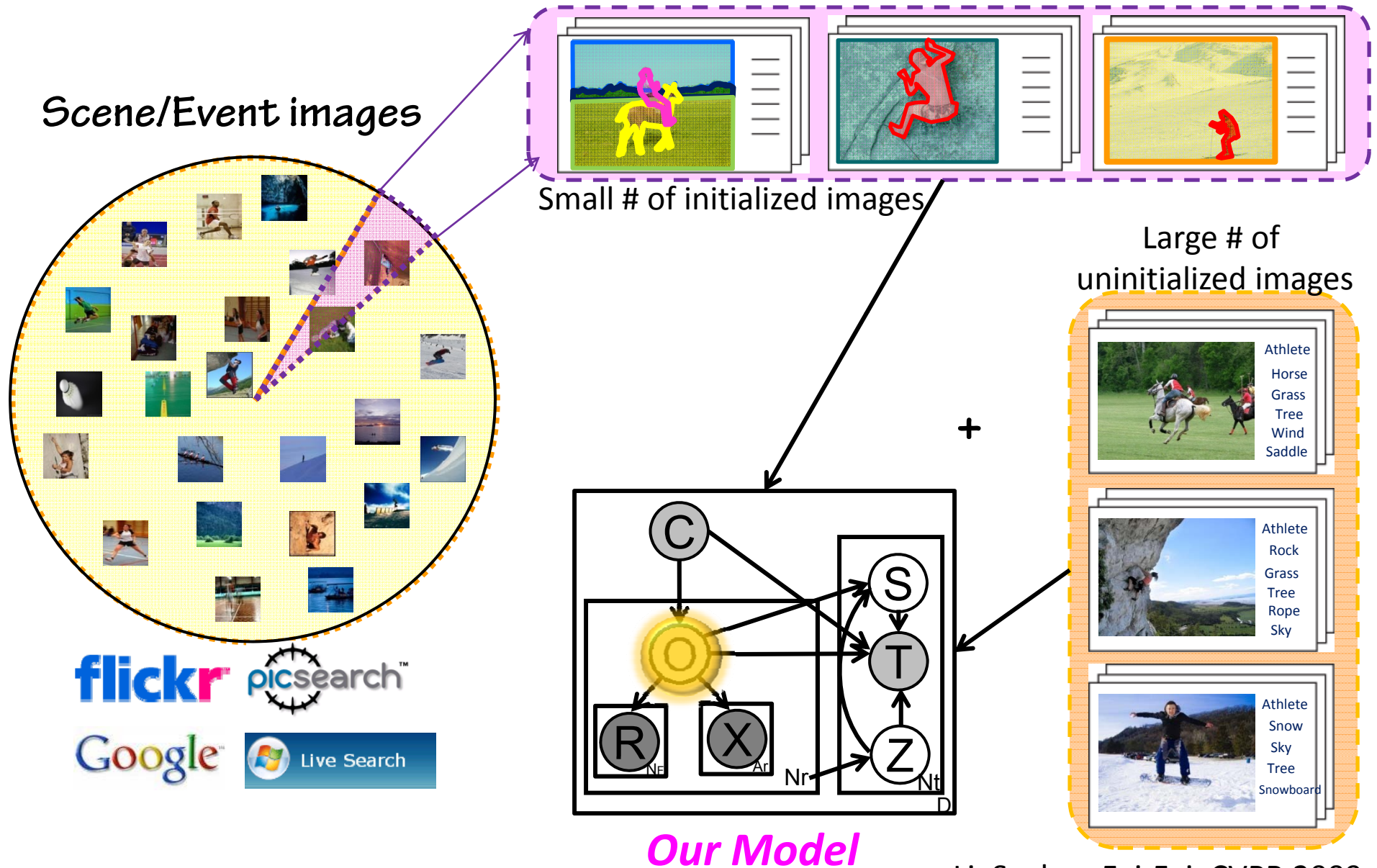
# Learning



Exact Inference is Intractable!

Relationship of the random variables

$$p(C, \mathbf{O}, \mathbf{R}, \mathbf{X}, \mathbf{S}, \mathbf{T}, \mathbf{Z} | \eta, \alpha, \beta, \gamma, \theta, \varphi) = p(C) \cdot (\prod_{n=1}^{N_r} p(O_n | \eta, C)) \prod_{n=1}^{N_r} ((\prod_{i=1}^{N_F} p(R_{ni} | O_n, \alpha_i)) \cdot \prod_{r=1}^{A_r} p(X_{nr} | O_n, \beta)$$
$$\prod_{m=1}^{N_t} p(Z_m | N_r) \, p(S_m | O_{Z_m}, \gamma) \, p(T_m | O_{Z_m}, S_m, \theta, C, \varphi)$$

Li, Socher, Fei-Fei, CVPR 2009

# Auto-semi-supervised learning:

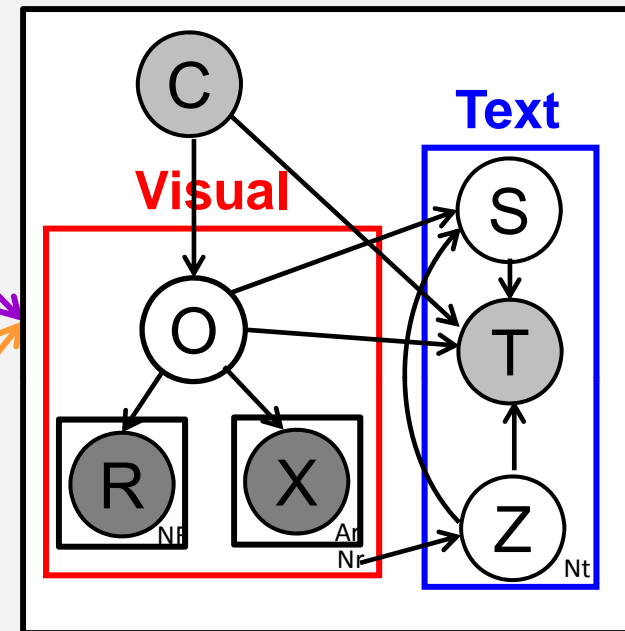Small # of initialized images + Large # of uninitialized images



Scene/Event images

Small # of initialized images

Large # of uninitialized images

Athlete
Horse
Grass
Tree
Wind
Saddle

Athlete
Rock
Grass
Tree
Rope
Sky

Athlete
Snow
Sky
Tree
Snowboard

*Our Model*

Li, Socher, Fei-Fei, CVPR 2009

# Outline

## Learning

flickr Small # of automatically initialized images

Athlete
Horse
Grass
Tree
Wind
Saddle

Athlete
Rock
Grass
Tree
Rope
Sky

Athlete
Snow
Sky
Tree
Snowboard

Large # of uninitialized images

## Model

Visual

Text

C

S

O

T

R    X

Z

## Recognition & Experiment
- Dataset
- Learned Model
- Results

Li, Socher, Fei-Fei, CVPR 2009
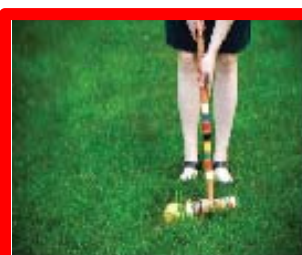
**flickr** **8 Event/Scene Classes**

Badminton

Bocce

Croquet

Polo

Remark: Tags are not used during testing
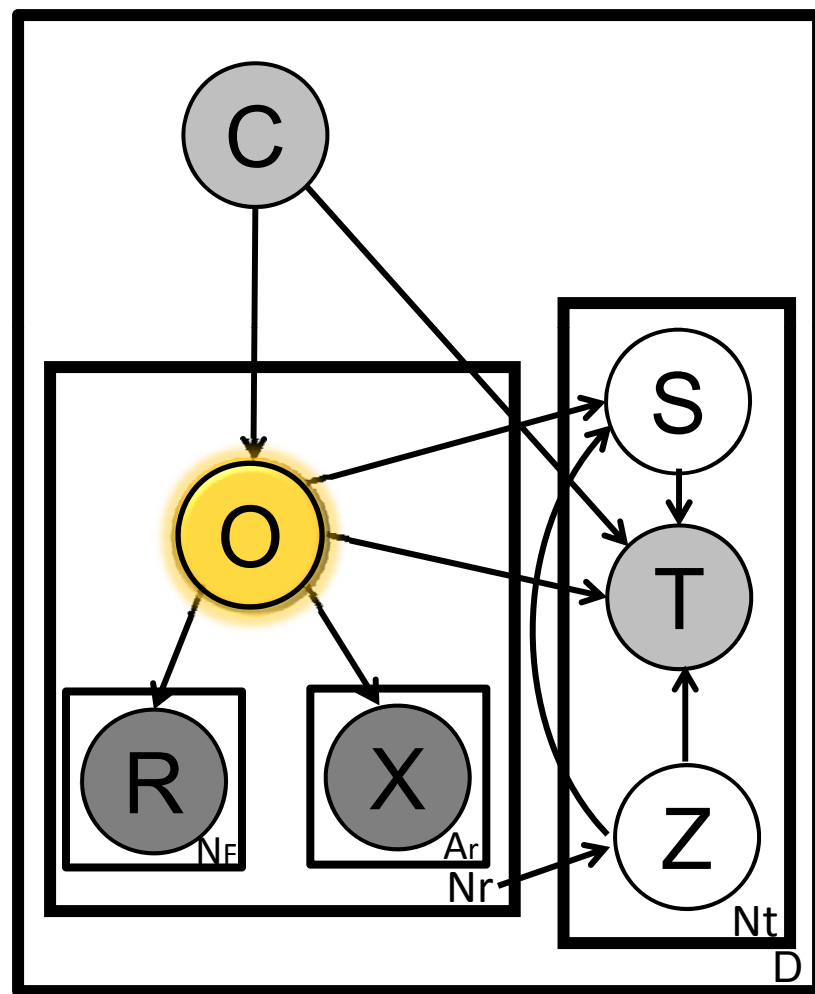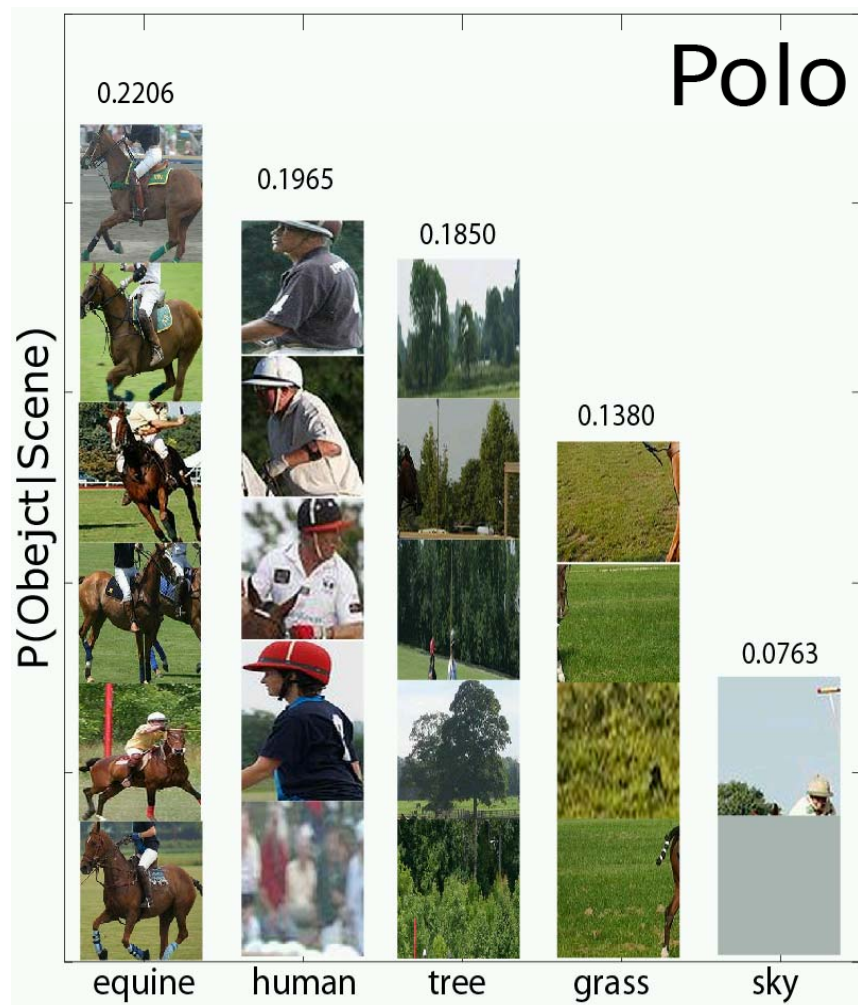
**flickr** 8 Event/Scene Classes

Rock climbing

Rowing

Sailing

Snow boarding

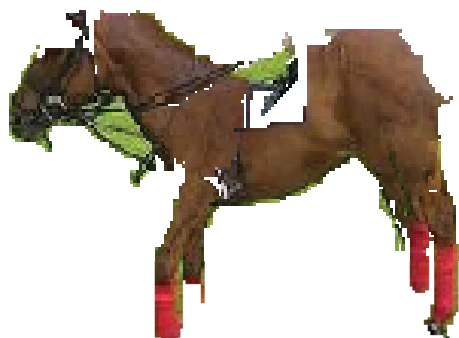Li, Socher, Fei-Fei, CVPR 2009

# Learned model: O

Li, Socher, Fei-Fei, CVPR 2009

Athlete

Grass

Horse

# Learned model: R

Li, Socher, Fei-Fei, CVPR 2009

Class: Badminton

Sail

Net

Human

Court

Class: Snowboarding

Cloud

Sky

Human

Snowboard

Tree

Rock

Snow

Class: Rowing

Human

Boat

Oar

Water

Class: Rock Climbing

Sky

Human

Mountain

Rock

Tree

Water