



# A Hierarchical Model of Shape and Appearance for Human Action Classification

Juan Carlos Niebles<sup>1,2</sup>

jnieble2@uiuc.edu

Li Fei-Fei<sup>1,3</sup>

feifeili@cs.princeton.edu



<sup>2</sup> Universidad del Norte, Colombia



<sup>3</sup> Princeton University

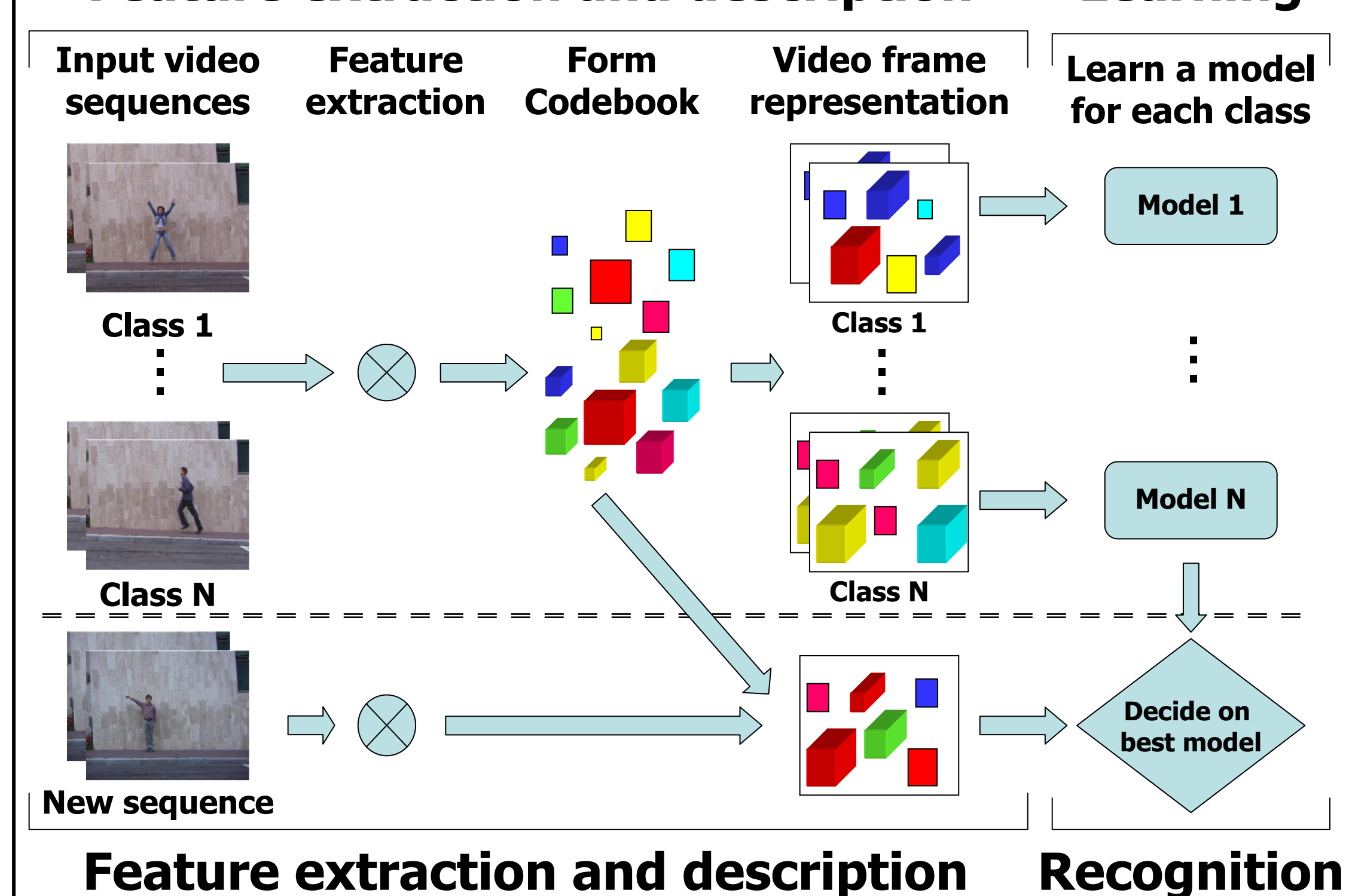
## Highlights and Summary

- A novel model for human action categorization from video sequences.
- Our model can be characterized as a constellation of bags-of-features.
- Use of hybrid features: combines both static shape and spatio-temporal features.

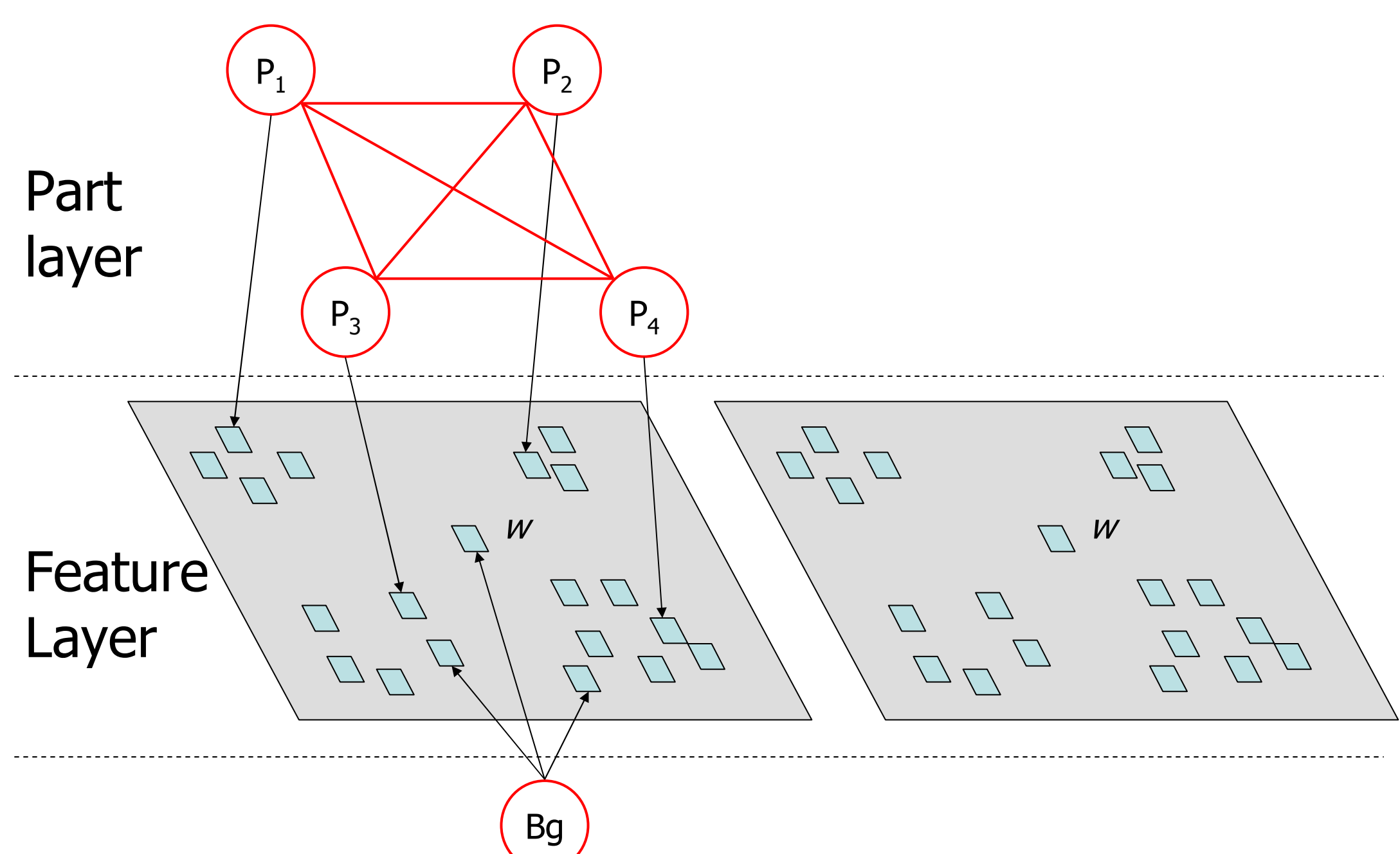
## Algorithm

### Feature extraction and description

### Learning



## Previous Works



### constellation

⊗ Small number of features

😊 Strong shape representation

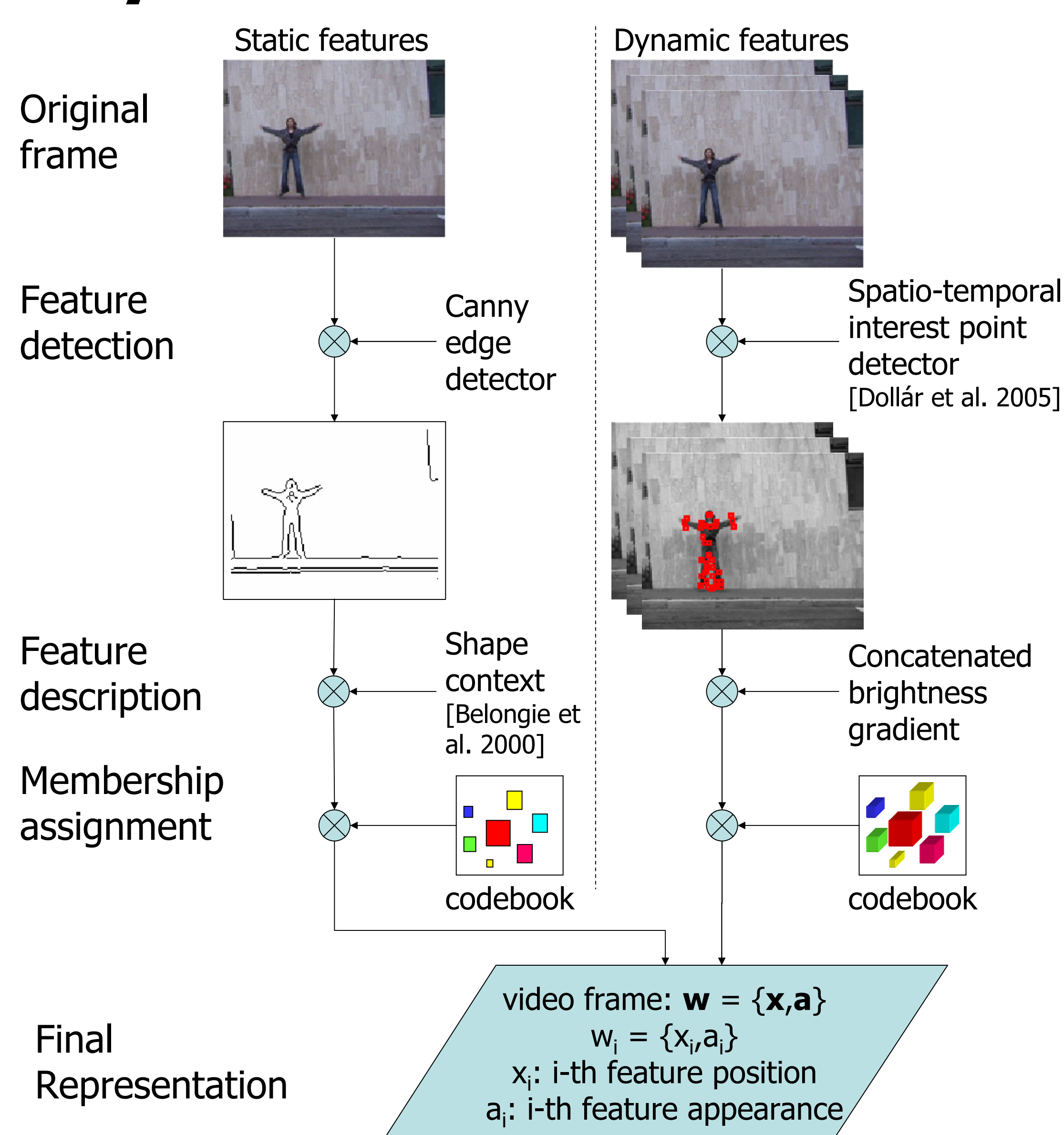
### bag-of-features

😊 Large number of features

⊗ No geometrical or shape information

[Weber et al. 2000, Csurka et al. 2004, Sudderth et al. 2006]

## Hybrid features



## Learning

Estimate model parameters using EM

$$\theta_\omega = \{\mu_{L,\omega}, \Sigma_{L,\omega}, \Sigma_{p,\omega}^X, \theta_{p,\omega}^A, \theta_0^X, \theta_0^A\} \quad p = 1 \dots P$$

$$\omega = 1 \dots \Omega$$

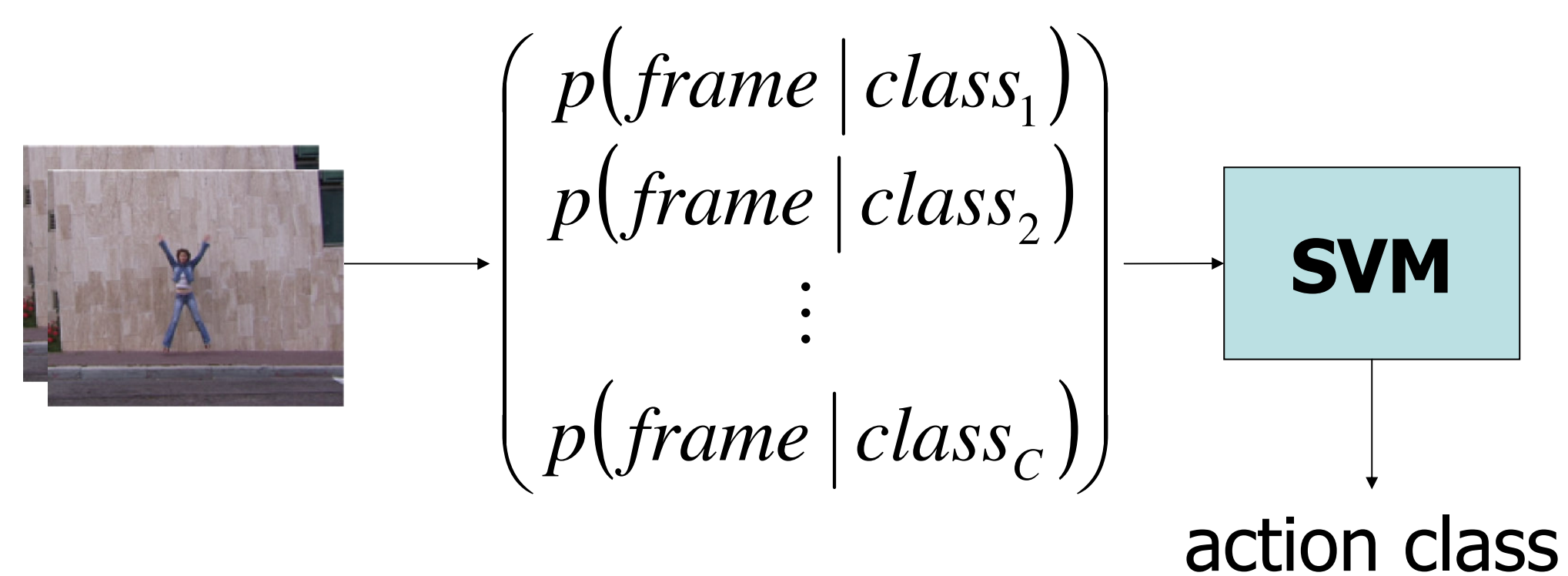
### • E-step:

$$p(\mathbf{h}, \omega | \mathbf{w}, \mathbf{Y}, \theta^{old}) \approx \frac{\pi_\omega p(\mathbf{Y} | \mathbf{h}, \theta_\omega^{old}) p(\mathbf{h} | \theta_\omega^{old}) p(\mathbf{w} | \mathbf{Y}, \mathbf{h}, \mathbf{m}^*, \theta_\omega^{old})}{p(\mathbf{w}, \mathbf{Y} | \theta_\omega^{old})}$$

### • M-Step:

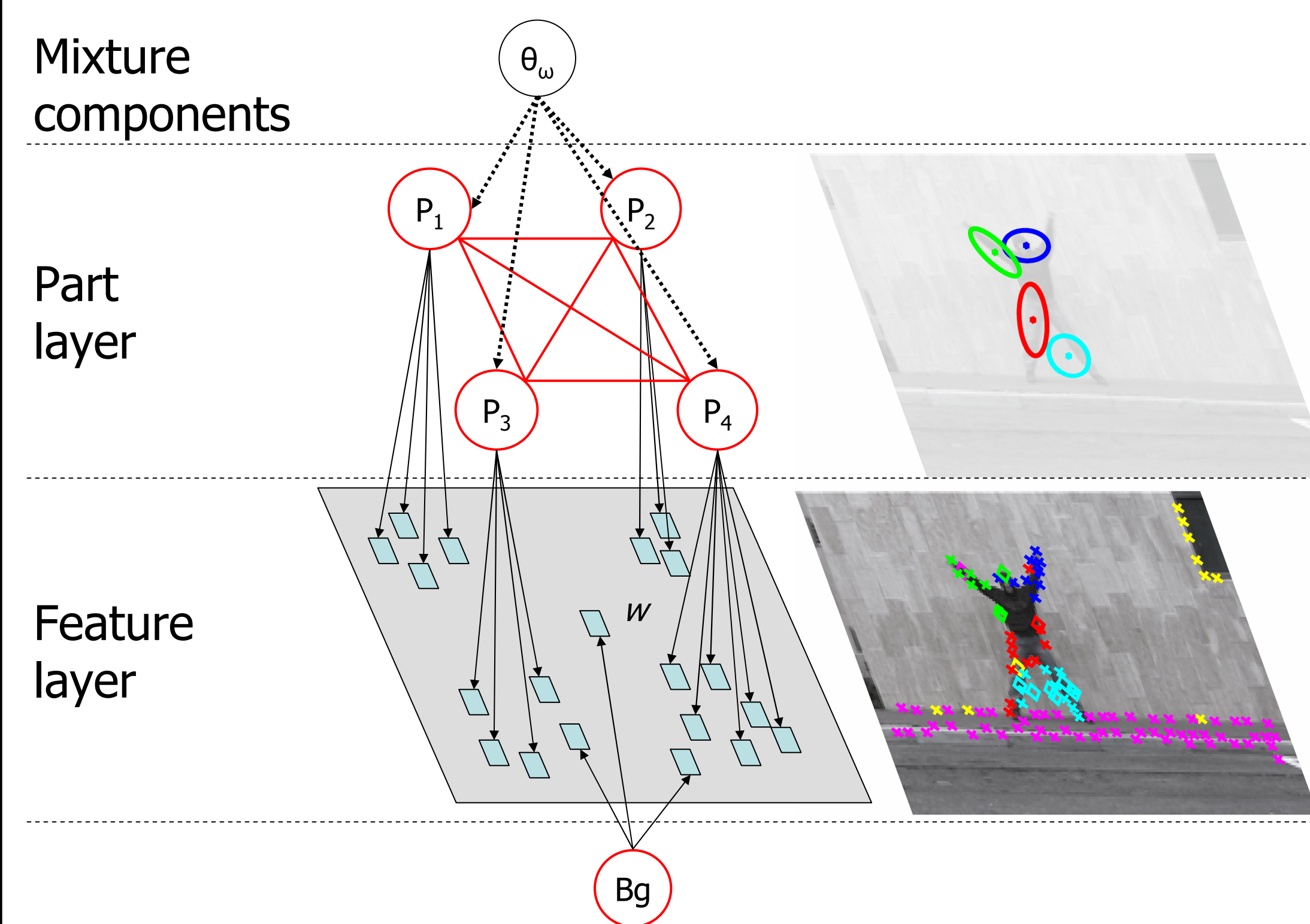
$$\theta^{new} = \arg \max_{\theta} \sum_{\mathbf{h}} p(\mathbf{h}, \omega | \mathbf{w}, \mathbf{Y}, \theta^{old}) \ln p(\mathbf{w}, \mathbf{Y}, \mathbf{h}, \omega | \theta)$$

## Recognition



- Classify actions in both frame based and video based manner
- Video classification based on majority votes of frames

## Hierarchical model



- Large number of features from the bag-of-features model
- Strong shape representation from the constellation model.

Approximated data likelihood:

$$p(\mathbf{w}, \mathbf{Y} | \theta) \approx \sum_{\omega=1}^{\Omega} \pi_\omega \sum_{\mathbf{h} \in H} p(\mathbf{h} | \theta_\omega) p(\mathbf{Y} | \mathbf{h}, \theta_\omega) p(\mathbf{w} | \mathbf{Y}, \mathbf{m}^*, \mathbf{h}, \theta_\omega)$$

Part layer term:

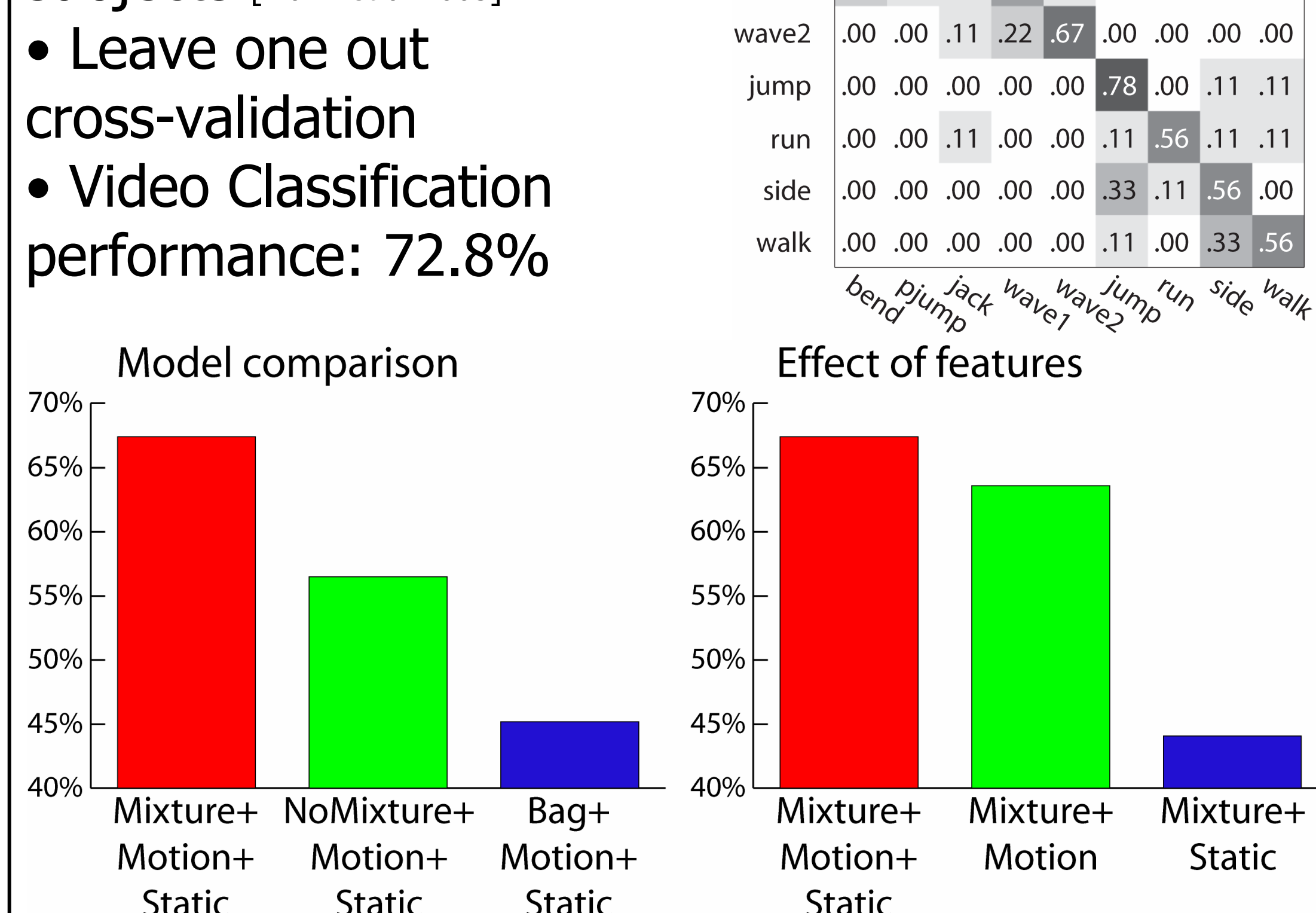
$$p(\mathbf{Y} | \mathbf{h}, \theta_\omega) = N(\mathbf{Y}_T(\mathbf{h}) | \mu_{L,\omega}, \Sigma_{L,\omega})$$

Local feature layer term:

$$p(\mathbf{w} | \mathbf{Y}, \mathbf{m}^*, \mathbf{h}, \theta_\omega) = \prod_{\mathbf{w}_j \in B_g} \underbrace{p(x_j^r | \theta_0^X)}_{Bg \text{ Shape}} \underbrace{p(a_j | \theta_0^A)}_{Bg \text{ Appearance}} \prod_{p=1}^P \prod_{\mathbf{w}_i \in P_p} \underbrace{p(x_i^r | \mathbf{Y}, h_p, \theta_p^X)}_{Part \text{ Shape}} \underbrace{p(a_i | \theta_p^A)}_{Part \text{ Appearance}}$$

## Experimental Results

- 9 action classes, performed by 9 subjects [Blank et al 2005]
- Leave one out cross-validation
- Video Classification performance: 72.8%



## Action Models



## Conclusions

- The constellation of bags-of-features is able to capture semantic information of human action classes.
- Combines hybrid features: static shape features and dynamic motion features.
- Capable of classifying in both frame based and video based manner.

**Ref:** J.C. Niebles & L. Fei-Fei. A hierarchical model of shape and appearance for human action classification. *CVPR* 2007. Minneapolis, USA.