



Sparse Dictionary-based Representation and Recognition of Action Attributes

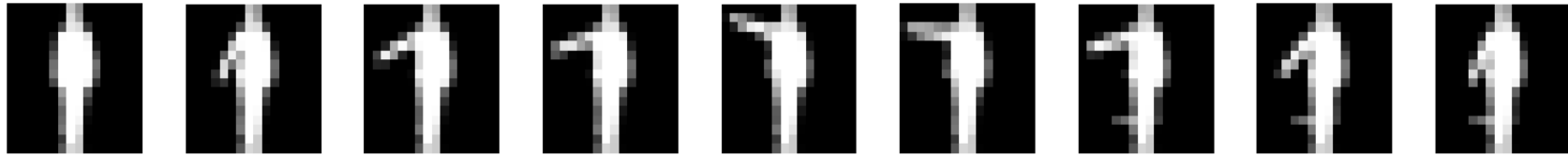
Qiang Qiu, Zhuolin Jiang, Rama Chellappa

Center for Automation Research,
Institute for Advanced Computer Studies
University of Maryland, College Park
qiu@cs.umd.edu, {zhuolin, rama}@umiacs.umd.edu

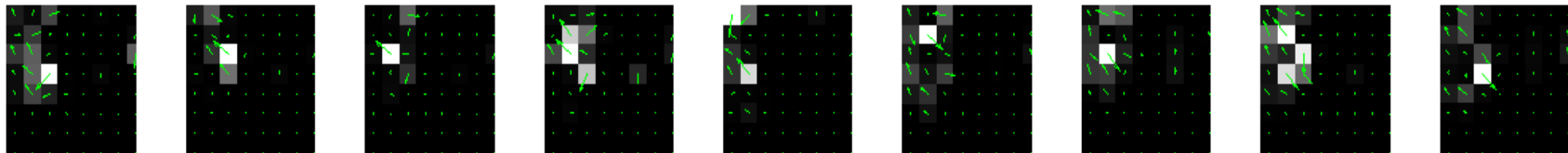
Action Feature Representation



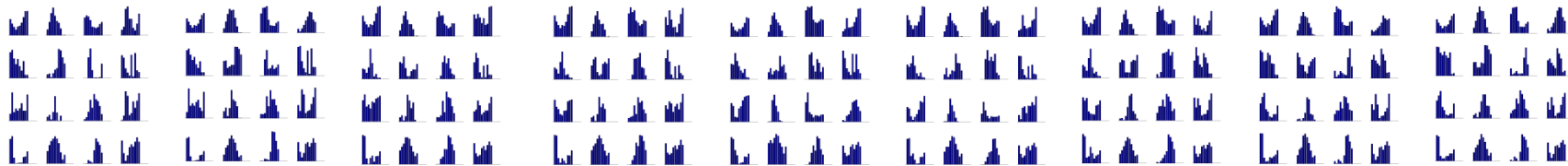
Shape



Motion



HOG



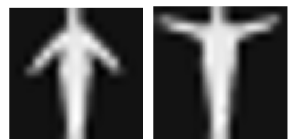
Action Sparse Representation



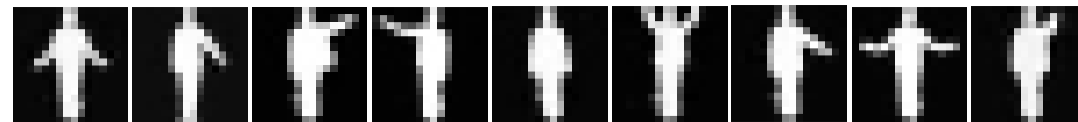
Sparse code

Action

Dictionary



=



<i>0.43</i>	<i>0</i>
<i>0.63</i>	<i>0</i>
<i>0</i>	<i>0.64</i>
<i>0</i>	<i>0.53</i>
<i>-0.33</i>	<i>-0.40</i>
<i>0</i>	<i>0.35</i>
<i>-0.36</i>	<i>0</i>
<i>0</i>	<i>0</i>
<i>0</i>	<i>0</i>

$$\text{Action 1} = 0.43 \times \text{Dict 1} + 0.63 \times \text{Dict 2} - 0.33 \times \text{Dict 3} - 0.36 \times \text{Dict 4}$$

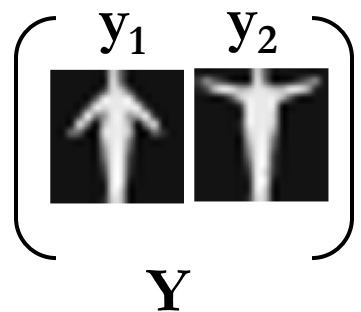
$$\text{Action 2} = 0.64 \times \text{Dict 5} + 0.53 \times \text{Dict 6} - 0.40 \times \text{Dict 7} + 0.35 \times \text{Dict 8}$$

K-SVD

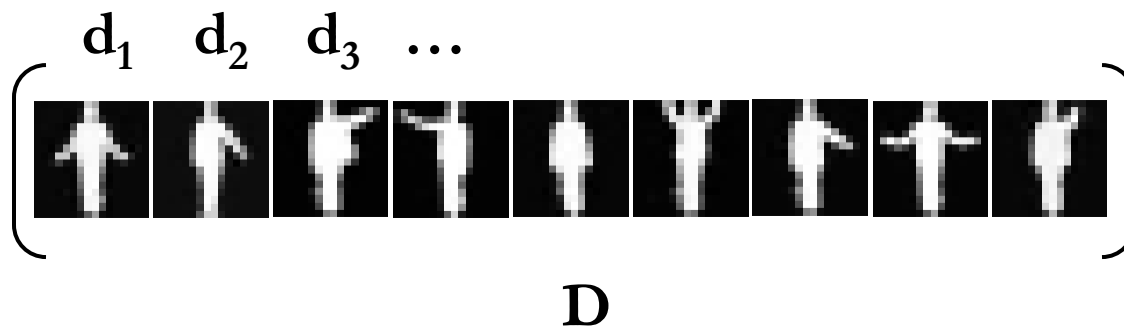


Sparse codes

Input signals



Dictionary



x_1	x_2
0.43	0
0.63	0
0	0.64
0	0.53
-0.33	-0.40
0	0.35
-0.36	0
0	0
0	0

X

$$\arg \min_{D, X} \|Y - DX\|^2 \quad \text{s.t.} \quad \forall i, \|x_i\|_0 \leq T$$

■ K-SVD [1]

- Input: signals Y , dictionary size, sparsity T
- Output: dictionary D , sparse codes X

Objective



Learn a

Compact and **Discriminative**

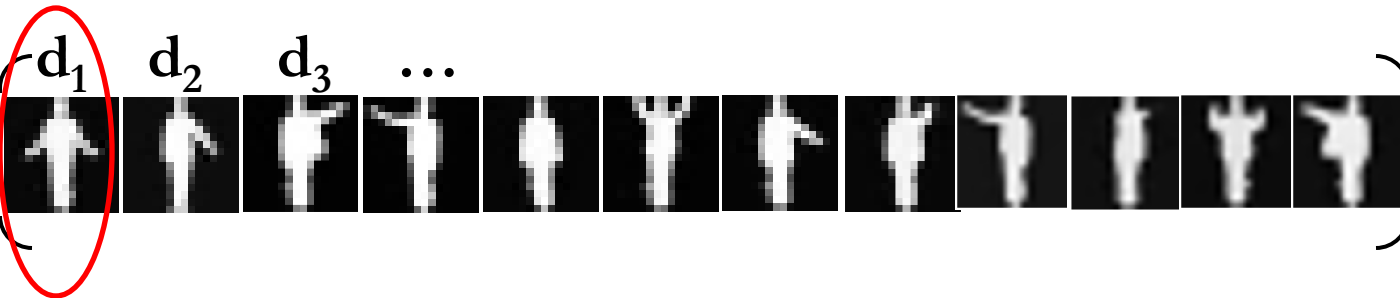
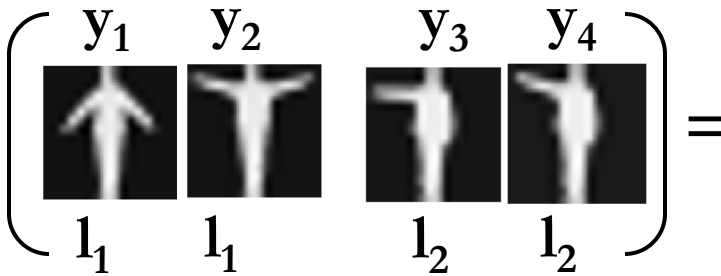
Dictionary.

Probabilistic Model for Sparse Representation



- **A Gaussian Process**
- **Dictionary Class Distribution**

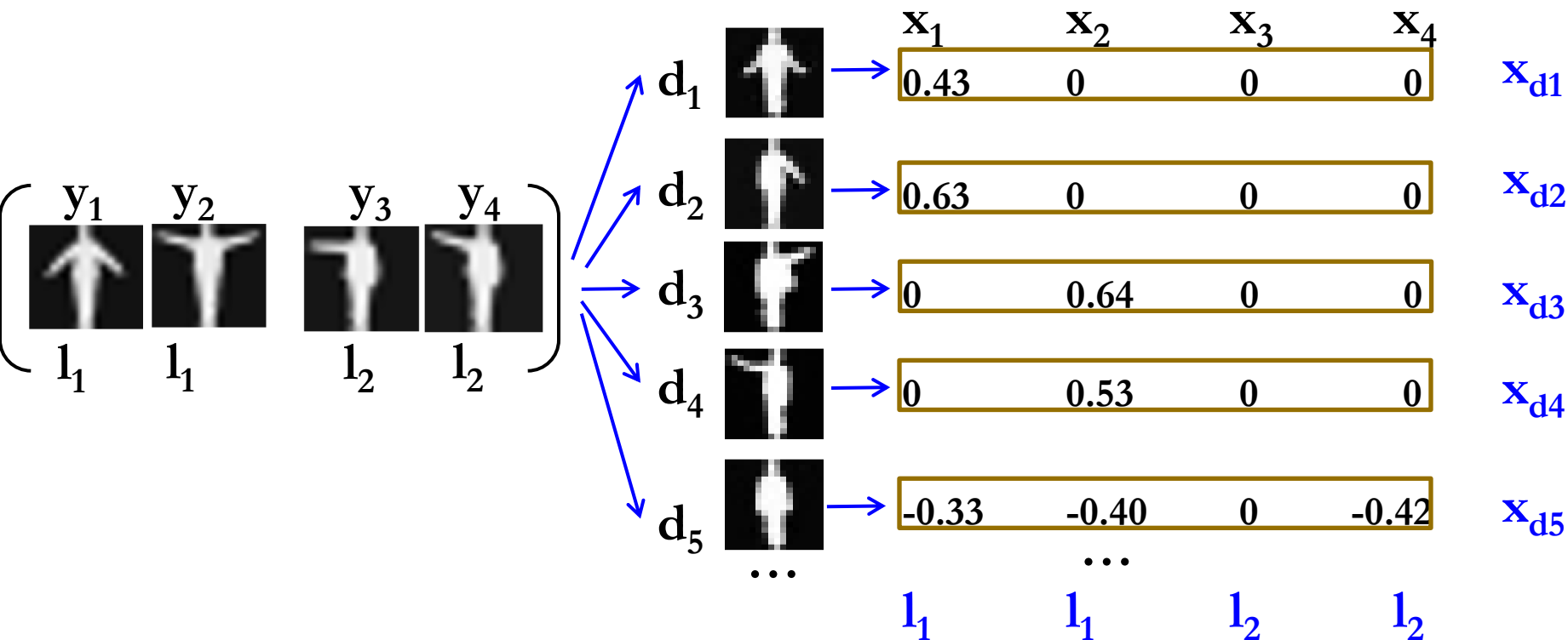
More Views of Sparse Representation



x_{d1}

x_1	x_2	x_3	x_4
0.43	0	0	0
0.63	0	0	0
0	0.64	0	0
0	0.53	0	0
-0.33	-0.40	0	-0.42
0	0.35	0	0
-0.36	0	0	0
0	0	0	0
0	0	-0.28	0
0	0	0.698	0.42
0	0	0.37	0.47
0	0	0.25	0
0	0	0	0.32
l_1	l_1	l_2	l_2

A Gaussian Process

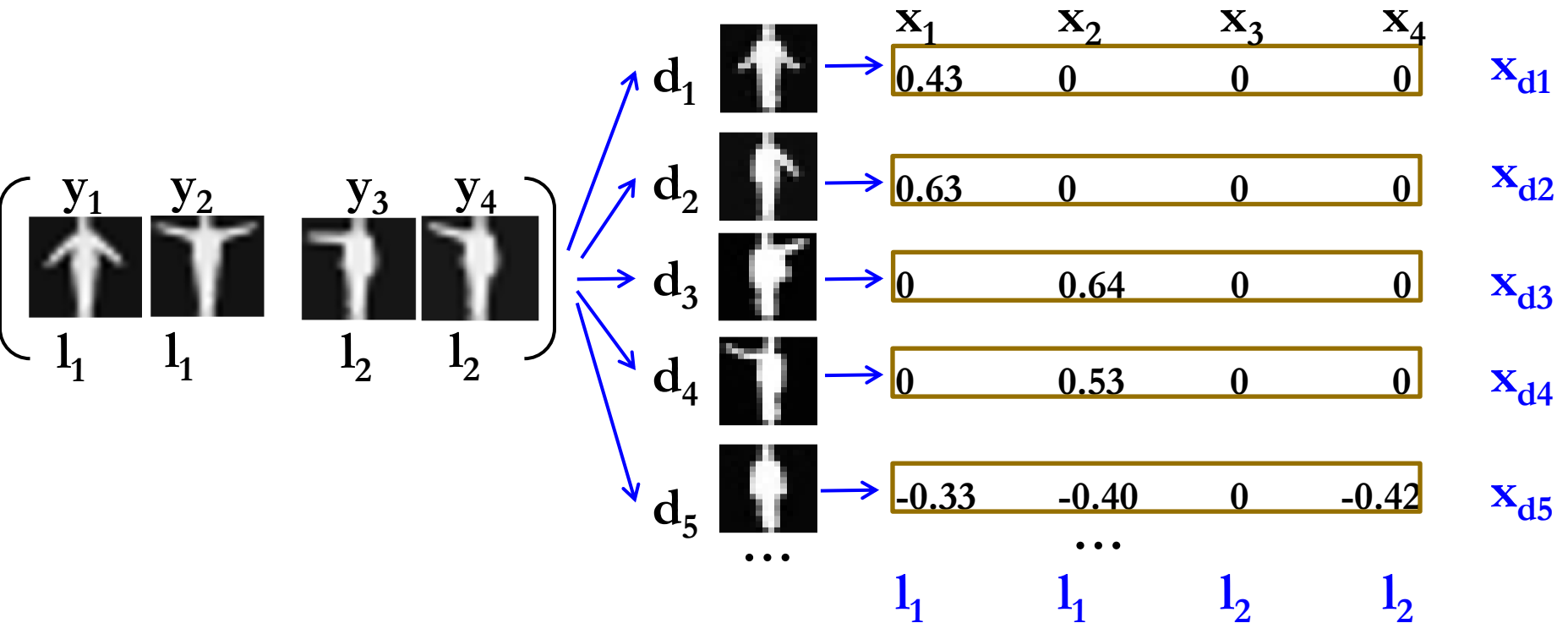


A Gaussian Process

- Covariance function entry: $K(i,j) = \text{cov}(x_{di}, x_{dj})$
- $P(X_{d^*} | X_{D^*})$ is a Gaussian with a closed-form conditional variance

$$\underline{\mathbb{V}(d^* | D^*) = \mathcal{K}_{(d^*, d^*)} - \mathcal{K}_{(d^*, D^*)}^T \mathcal{K}_{(D^*, D^*)}^{-1} \mathcal{K}_{(d^*, D^*)}} \underline{\hspace{10em}}$$

Dictionary Class Distribution



Dictionary Class Distribution

- $P(L | d_j), L \in [1, M]$
 - aggregate $|x_{d_i}|$ based on class labels to obtain a M sized vector
 - $P(L=l_1 | d_5) = (0.33+0.40)/(0.33+0.40+0.42) = 0.6348$
 - $P(L=l_2 | d_5) = (0+0.42)/(0.33+0.40+0.42) = 0.37$

Dictionary Learning Approaches



- Maximization of Joint Entropy (ME)
- Maximization of Mutual Information (MMI)
 - Unsupervised Learning (MMI-1)
 - Supervised Learning (MMI-2)

Maximization of Joint Entropy (ME)

- Initialize dictionary using k-SVD

$$D^0 = \left(\begin{array}{c} \text{[Grid of 15 grayscale images of handwritten digits]} \end{array} \right)$$

- Start with $D^* = \phi$
- Untill $|D^*| = k$, iteratively choose d^* from $D^0 \setminus D^*$,

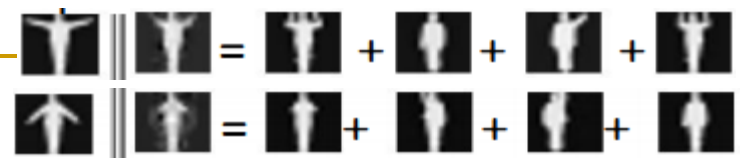
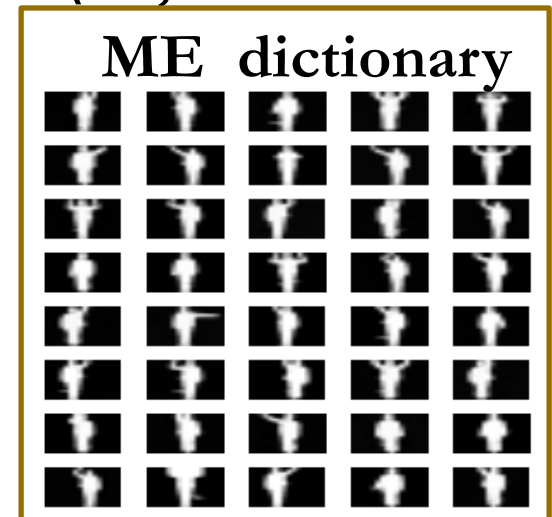
$$d^* = \arg \max_{d} H(d | D^*)$$

Where

$$H(d^* | D^*) = \frac{1}{2} \log(2\pi e \mathbb{V}(d^* | D^*))$$

- A good approximation to ME criteria

$$\arg \max_{D} H(D)$$



Maximization of Mutual Information for Unsupervised Learning (MMI-1)



- Initialize dictionary using k-SVD

$$D^{\circ} = \begin{pmatrix} \text{[Grid of 100 small images]} \end{pmatrix}$$

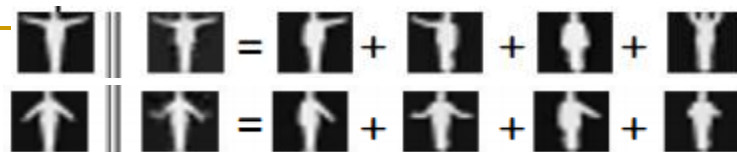
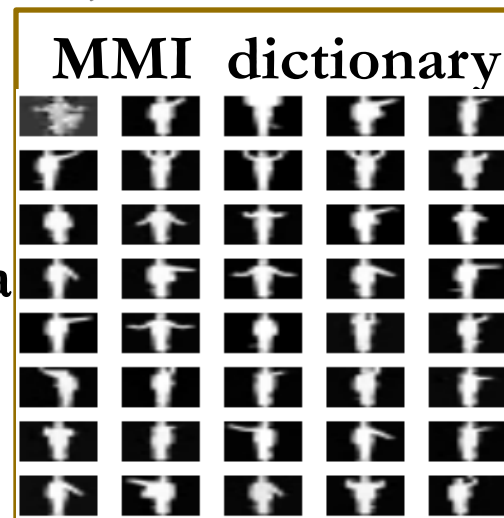
- Start with $D^* = \phi$
- Until $|D^*| = k$, iteratively choose d^* from $D^{\circ} \setminus D^*$,

$$d^* = \arg \max_d H(d | D^*) - H(d | D^{\circ} \setminus (D^* \cup d))$$

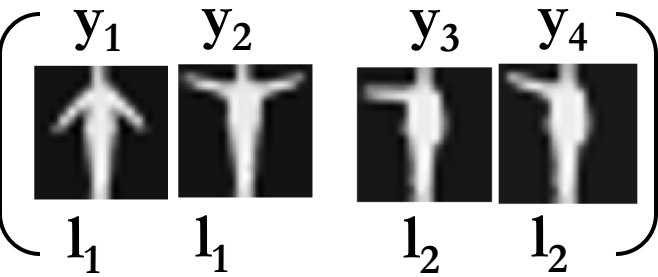
- A near-optimal approximation to MMI criteria

$$\arg \max_D I(D; D^{\circ} \setminus D)$$

Within $(1-1/e)$ of the optimum



Revisit



	x_1	x_2	x_3	x_4	
d_1	0.43	0	0	0	x_{d1}
d_2	0.63	0	0	0	x_{d2}
d_3	0	0.64	0	0	x_{d3}
d_4	0	0.53	0	0	x_{d4}
d_5	-0.33	-0.40	0	-0.42	x_{d5}
	...				
	l_1	l_1	l_2	l_2	

Dictionary Class Distribution

- $P(L | d_i), L \in [1, M]$
 - aggregate $|x_{di}|$ based on class labels to obtain a M sized vector
 - $P(l_1 | d_5) = (0.33+0.40)/(0.33+0.40+0.42) = 0.6348$
 - $P(l_2 | d_5) = (0+0.42)/(0.33+0.40+0.42) = 0.37$
- $P(L_D) = P(L/d)$
- $P(L_D) = P(L/D)$, where $P(L|D^*) = \frac{1}{|D^*|} \sum_{d_i \in D^*} P(L|d_i)$

Maximization of Mutual Information for Supervised Learning (MMI-2)



- Initialize dictionary using k-SVD

$$D^0 = \left(\begin{array}{cccccccccccccccc} \text{[Grid of 16x16 small images]} \end{array} \right)$$

- Start with $D^* = \phi$
- Untill $|D^*| = k$, iteratively choose d^* from $D^0 \setminus D^*$,

$$d^* = \arg \max_d [H(d | D^*) - H(d | D^0 \setminus (D^* \cup d))] + \lambda [H(L_d | L_{D^*}) - H(L_d | L_{D^0 \setminus (D^* \cup d)})]$$

- MMI-1 is a special case of MMI-2 with $\lambda=0$.

Other learning methods



- **K-means**
- **Liu-shah [1]**

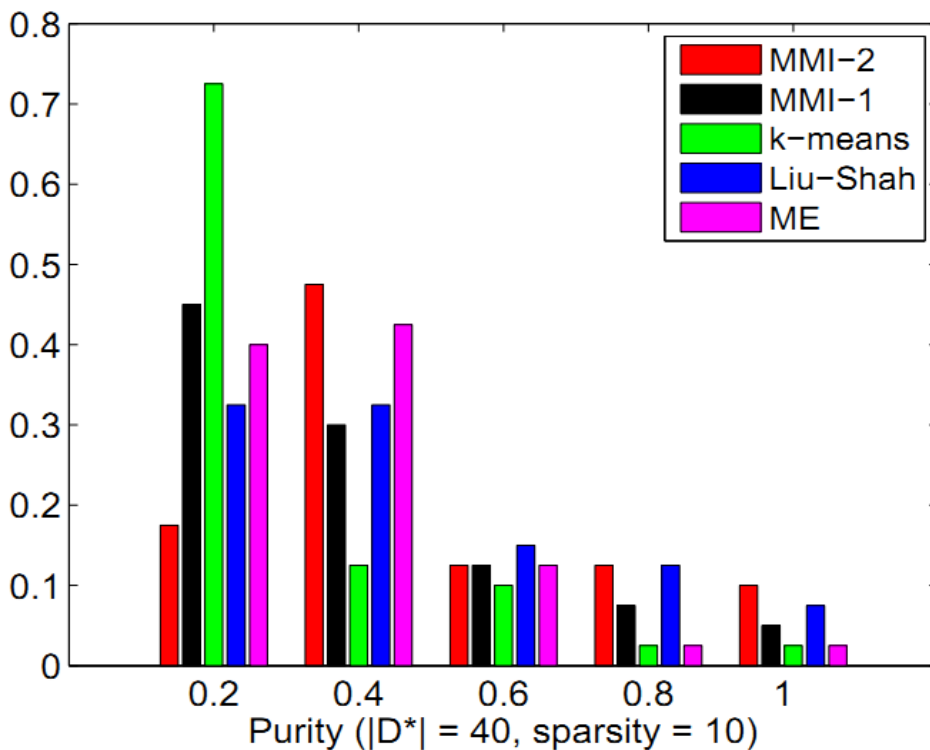
$$\Delta I(d_1, d_2) = \sum_{L \in [1, M], i=1,2} p(d_i) p(L|d_i) \log p(L|d_i) - p(d_i) p(L|d_i) \log p(L|d^*)$$

where

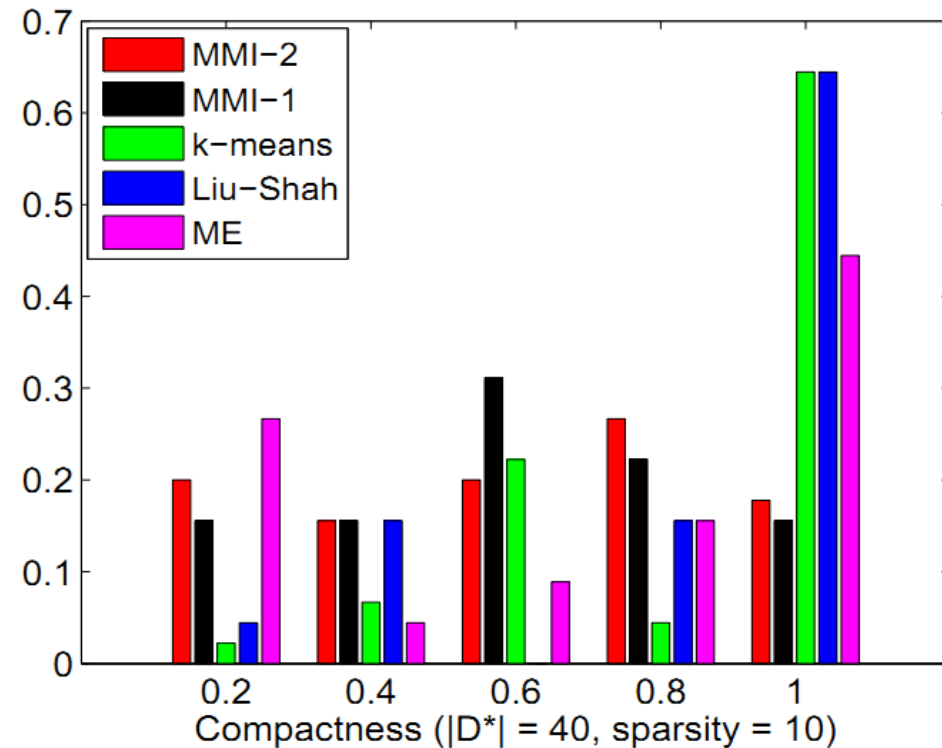
$$p(L|d^*) = \frac{p(d_1)}{p(d^*)} p(L|d_1) + \frac{p(d_2)}{p(d^*)} p(L|d_2)$$

$$p(d^*) = p(d_1) + p(d_2)$$

Purity and Compactness

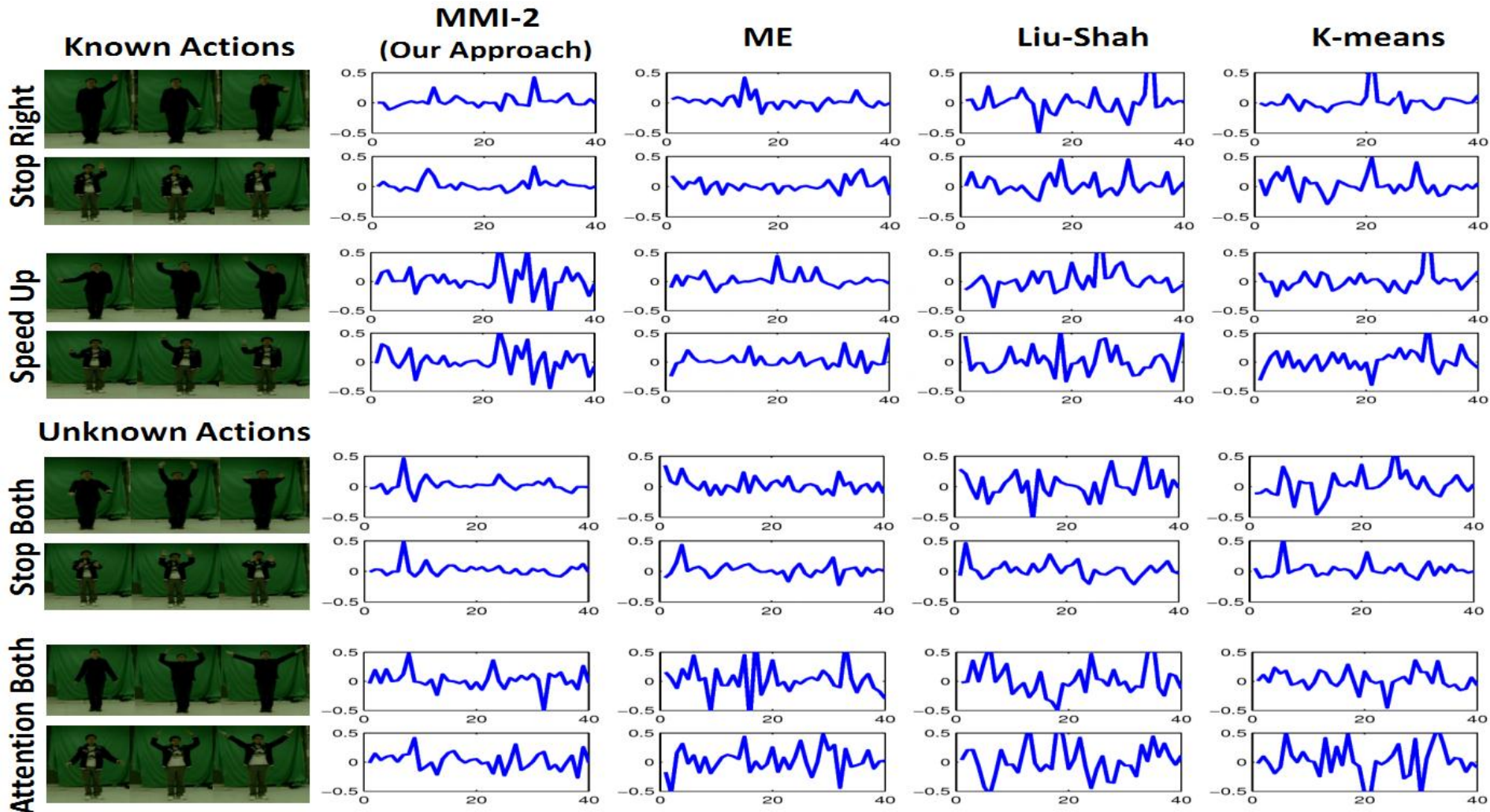


(a) Purity



(b) Compactness

Representation Consistency



Keck gesture dataset



Turn left



Turn right



Attention left



Attention right



Attention both



Stop left



Stop right



Stop both



Flap



Start



Go back



Close distance

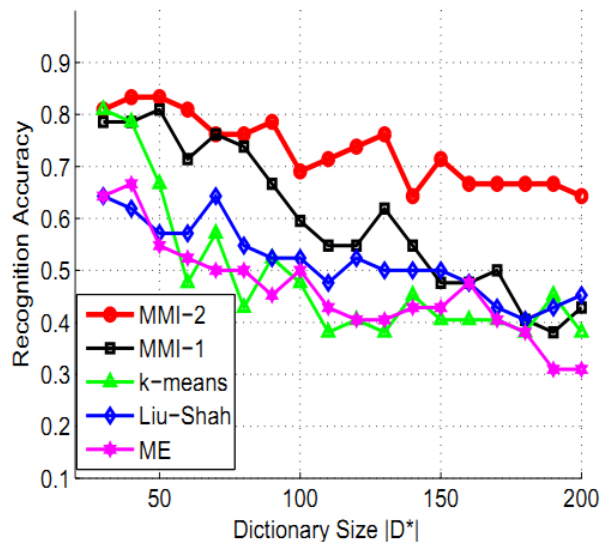


Speed up

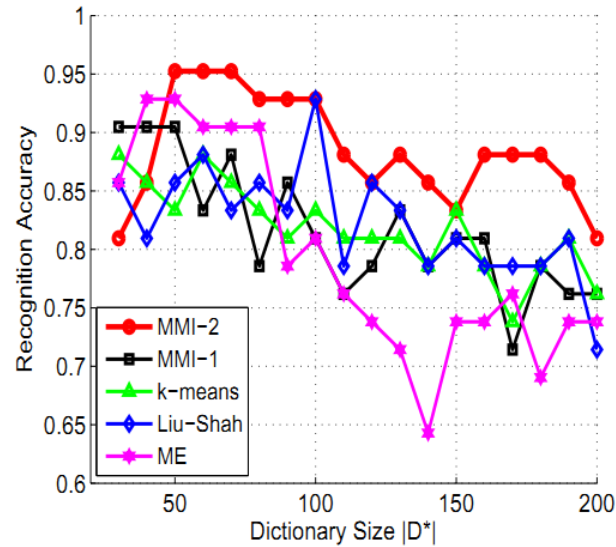


Come near

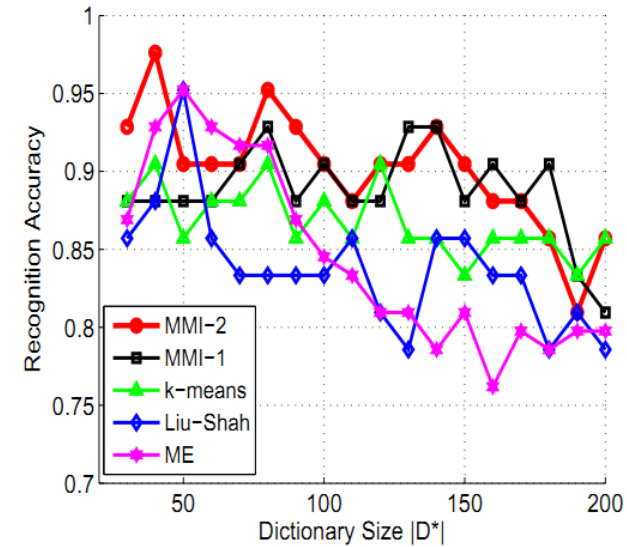
Recognition Accuracy



(a) Shape ($|D^o| = 600$)



(b) Motion ($|D^o| = 600$)



(c) Shape and Motion ($|D^o| = 1200$)

The recognition accuracy using initial dictionary D^o : (a) 0.23 (b) 0.42 (c) 0.71