





1. Summary

- We proposed a method to deal with the short effective distance problem of RGB-D cameras to take advantage of RGB-D cameras which are not online accessible. We illustrated it with the application to action recognition.
- Our approach
- > We use Kinect to offline collect a multi-modal database.
- > Our goal is to augment actions to be recognized with proper depth and skeleton features.
- We do domain adaptation (DA), by which the actions to be recognized can be concisely reconstructed by entries in the auxiliary database.
- Each action can be augmented with additional depth and skeleton images retrieved from the auxiliary database, and we use multiple kernel learning (MKL) to fuse these three kinds of features for action recognition.

2. Related Work

2.1

Feature Descriptor

- Designing powerful feature descriptors against intra-class variations has gained significant progress.
 - Global: sensitive to occlusion and deformation
 - Local: less discriminant power
- Descriptors based on RGB data are vulnerable to camera perspectives and partial occlusions.
- RGB-D cameras can provide depth data, which affords informative clues for action recognition.
- *OpenNI library* provides the positions of key joints on the human body, *i.e.*, the skeleton, which is also helpful for action recognition.
- The short ranges of effective distances still make RGB-D cameras in real-world applications, e.g., surveillance.

2.2

Transfer Learning

- Training data acquisition is also a challenge.
- > Complete training data may not be available.
- Labeling huge data is expensive.
- Transfer learning can alleviate the above problems.
- Our approach also utilizes knowledge transferred from an auxiliary database via DA and data reconstruction for improving action recognition.
- Difference: we borrow visual features across different video modalities, and resolves the problems caused by the absence of RGB-D cameras.

2.3

Multiple Kernel Learning

- MKL derives an optimal kernel over a given convex set of kernels, which means it can fuse heterogeneous features and lead to boosted performance.
- In our case, we fuse the original RGB features as well as the augmented depth and skeleton features.

Depth and Skeleton Associated Action Recognition without Online Accessible RGB-D Cameras

Yen-Yu Lin¹, Ju-Hsuan Hua², Nick C. Tang¹, Min-Hung Chen¹, and Hong-Yuan Mark Liao¹ ¹Academia Sinica, Taiwan; ²Carnegie Mellon University, USA



- Goal: borrow depth and skeleton information from an auxiliary, multi-modal database.
- Symbols:
- ▶ Target actions: $D = \{x_i, y_i\}_{i=1}^N, x_i$: RGB data, y_i : label $\in \{1, 2, ..., C\}$
- > Auxiliary database: $\widetilde{D} = \{\widetilde{x}_i, \widetilde{d}_i, \widetilde{s}_i\}_{i=1}^M, \widetilde{d}_i$: depth data, \widetilde{s}_i : skeleton data
- We focus on associating each action $x_i \in D$ with proper depth map d_i and skeleton structure s_i , so that the absence of RGB-D cameras is compensated.

Domain Adaptation

• We adopt reconstruction-based DA to tackle the inter-database variations, which means we model D and \widetilde{D} by a linear transformation and reconstruction.

$$WX = \tilde{X}A + E \tag{1}$$

• With discriminant learning and outlier handling, we obtain

$$\min_{W,A,E} \sum_{c=1}^{S} rank(A^{c}) + \lambda \|E\|_{2,1}$$
(2)
s. t. $WX = \tilde{X}A + E$ and $WW^{T} = I$

- $J_1 C_1 \vee M = M H + L H H + V = V$ • For optimization, we adopt some modifications
- > Serve nuclear norm as a convex approximation of rank minimization
- > Introduce an auxiliary variable $F = [F^1 \cdots F^c]$
- > Orthogonalize *W* afterwards with orthogonality preserving method

$$\min_{W,A,E} \sum_{c=1}^{c} \|F^{c}\|_{*} + \lambda \|E\|_{2,1}$$
(3)
s.t. $WX = \tilde{X}A + E$ and $A = F$

3. The Proposed Method

3.2

Optimization

• We optimize (3) by the inexact Augmented Lagrange Multiplier (ALM) method, which minimizes the augmented Lagrange function of (3):

$$\min_{W,A,F,E,U,V} \sum_{c=1}^{c} \|F^{c}\|_{*} + \lambda \|E\|_{2,1} + \langle U, A - F \rangle + \frac{\mu}{2} \|A - F\|_{F}^{2} + \langle V, WX - \tilde{X}A - E \rangle + \frac{\mu}{2} \|WX - \tilde{X}A - E\|_{F}^{2}$$
(4)

- Alternate optimization for {*W*, *A*, *F*, *E*}
- > 1. F: use singular value shrinkage operator
- > 2. W: closed-form solution
- > 3. Apply QR-decomposition to orthogonalize W
- > 4. *E*: use the analytical solution
- > 5. A: closed-form solution
- > 6. update Lagrange multipliers and penalty parameters
- > 7. check convergence conditions

3.3

Feature Augmentation

• After obtaining transformation W, we optimize reconstruction coefficients for both training and unseen testing data by solving $\boldsymbol{\alpha} = \arg \min \| W \boldsymbol{x} - \tilde{X} \boldsymbol{\alpha} \|^2 + \gamma \| \boldsymbol{\alpha} \|^2$ (5)

• Use the obtained
$$\alpha$$
 to reconstruct depth and skeleton features:
 $d \leftarrow [\widetilde{d}_1 \cdots \widetilde{d}_M] \alpha$ and $s \leftarrow [\widetilde{s}_1 \cdots \widetilde{s}_M] \alpha$ (6)

• We compile an kernel matrix for actions in each modality, and adopt simpleMKL [Rakotomamonjy et al. JMLR 2008] to learn classifiers that optimally combine the three types of features.



🛨 Ours: d

- 1NN-Bo

KSDA

4 5 6 Number of Actors in Training

- 1NN-Bor

KSDA

4 5 6 Number of Actors in Training

🛨 Ours: d

1NN-Bor

KSDA

3 4 5 6 7 8 Number of Actors in Training