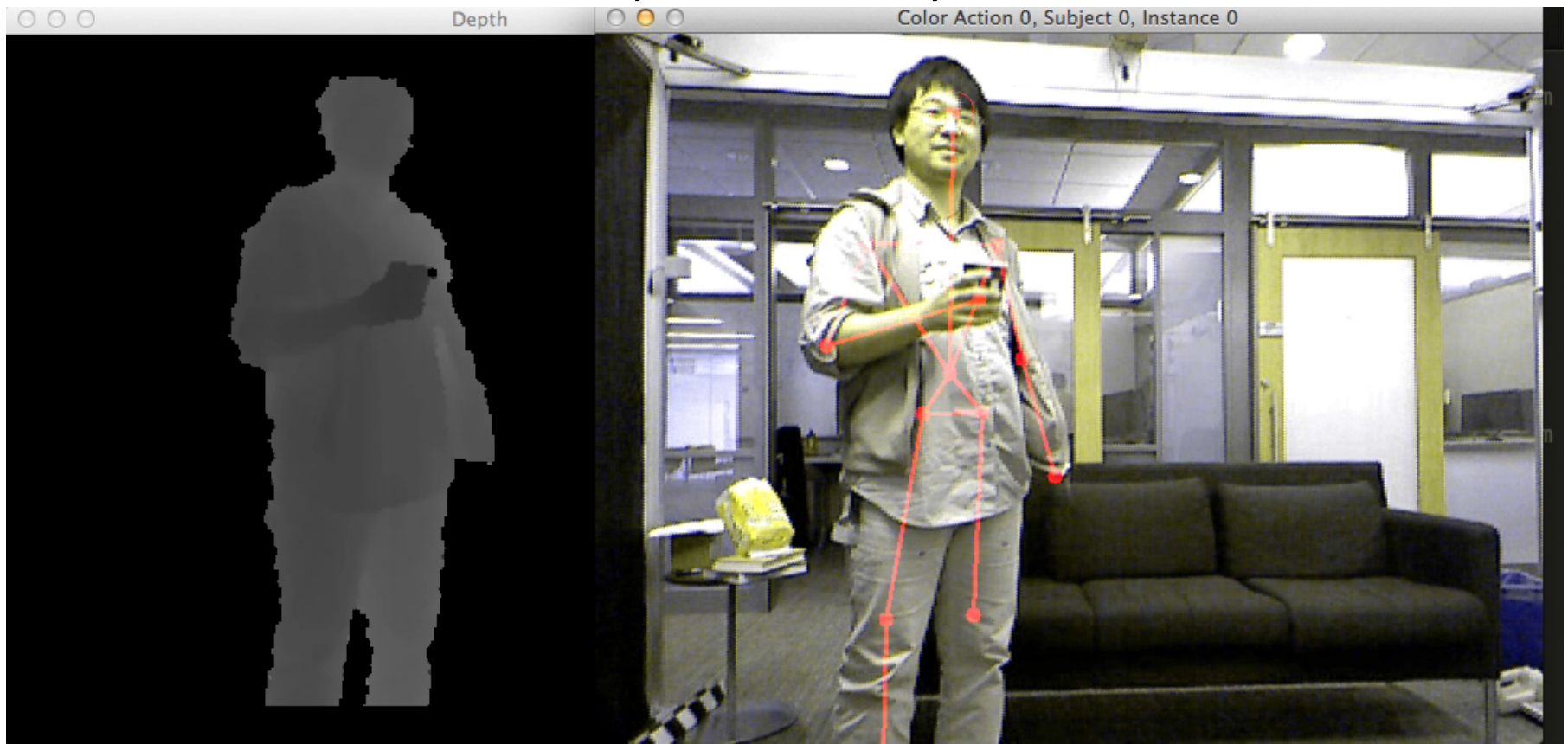# Mining Actionlet Ensemble for Action Recognition with Depth Cameras
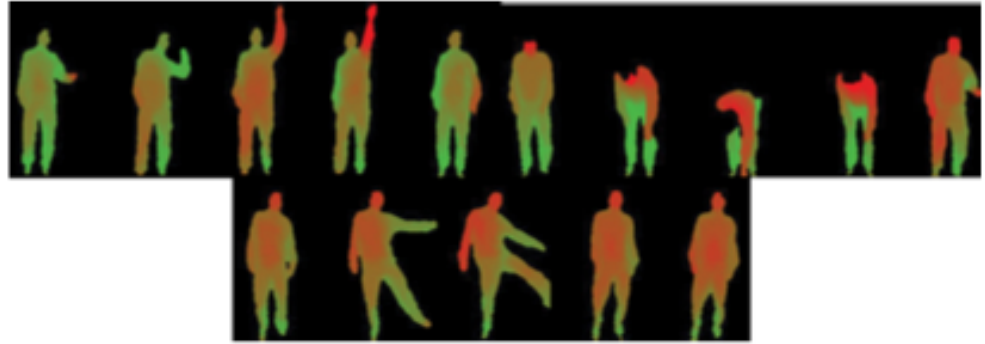
Jiang Wang[1]    Zicheng Liu[2]    Ying Wu[1]    Junsong Yuan[3]

[1]Northwestern University    [2]Microsoft Research    [3]Nanyang Technological University
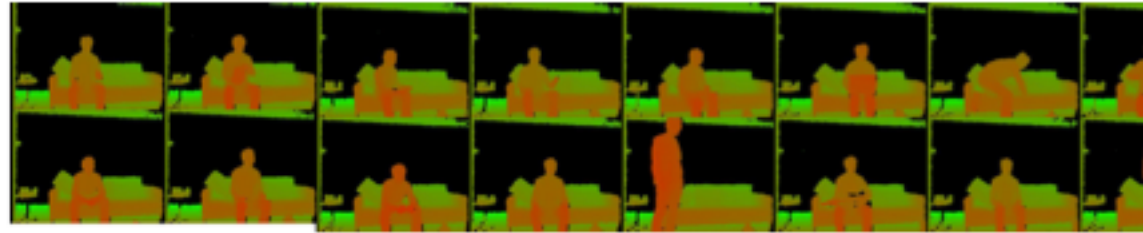
(CVPR 2012)

# Datasets



- **MSR Action 3D**

20 Actions (e.g. high arm wave, horizontal arm wave, forward punch, high throw)

10 people, 3x each - 26 minutes total footage, 402 total actions

15 FPS, 640x480

- **MSR Daily Activity 3D**



16 Actions (e.g. drink, eat, read book, call cellphone, cheer up, play guitar, lay down)

10 people, 2x each - 320 total actions

"Living room activities"

- **CMU MoCap dataset**

5 actions (walking, marching, dribbling, walking with stiff arms, walking with wild legs)
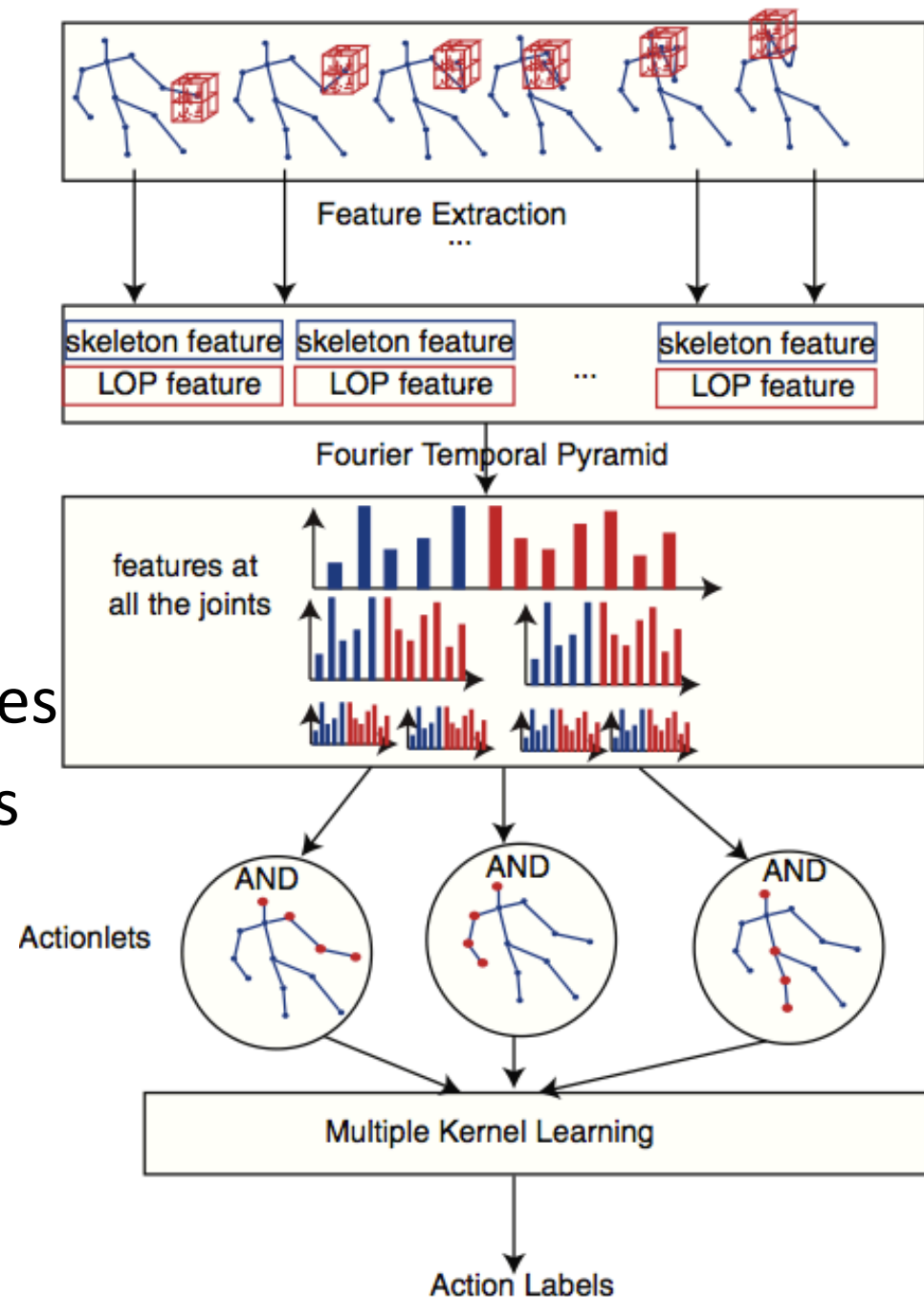
? People, ? iterations

Only skeletons

# Contributions

- **Actionlet ensemble**
  - **(their pipeline)**
- **Features**
  - **Relative positions**
  - **Local Occupancy Pattern**

- **Temporal representation**
  - **Fourier temporal pyramid**

# Overview

1) Extract features:

   3D Local Binary Patterns ("LOP")

   Skeletal Features

2) Combine features

3) Look at different time scales

4) Combine *top* joint features

5) MKL classification

# Features (1/2)

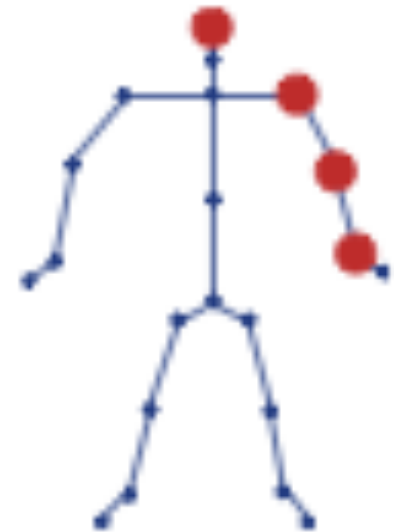**Invariant 3D joint positions**

- Pairwise-relative positions

pairwise

body part feature

$$\boldsymbol{p}_{ij} = \boldsymbol{p}_i - \boldsymbol{p}_{j:} \qquad \boldsymbol{p}_i = \{\boldsymbol{p}_{ij} | i \neq j\}$$

- Normalize

- Invariances: Translation, body size

- Robust to noise & temporal misalignment

# Features (2/2)

**Local Occupancy Patterns**

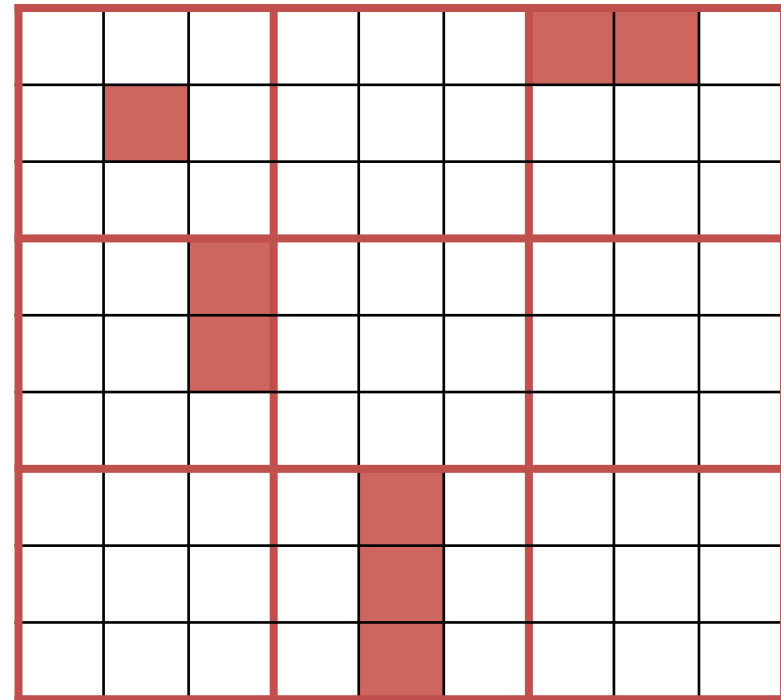- Related Work: SIFT, STIPS, HOG, Cuboids, LBP …

- Model interactions

For each joint:

1) Partition into $N_x$ x $N_y$ x $N_z$ grid.

2) Sum all entries in each bin

3) Apply sigmoid

4) Concatenate value at all bins

$$o_{xyz} = \delta(\sum_{q \in \text{bin}_{xyz}} I_q)$$

$$\delta(x) = \frac{1}{1 + e^{-\beta x}}$$

X = [ 1 0 2 2 0 0 0 3 0]

# Fourier Temporal Pyramid

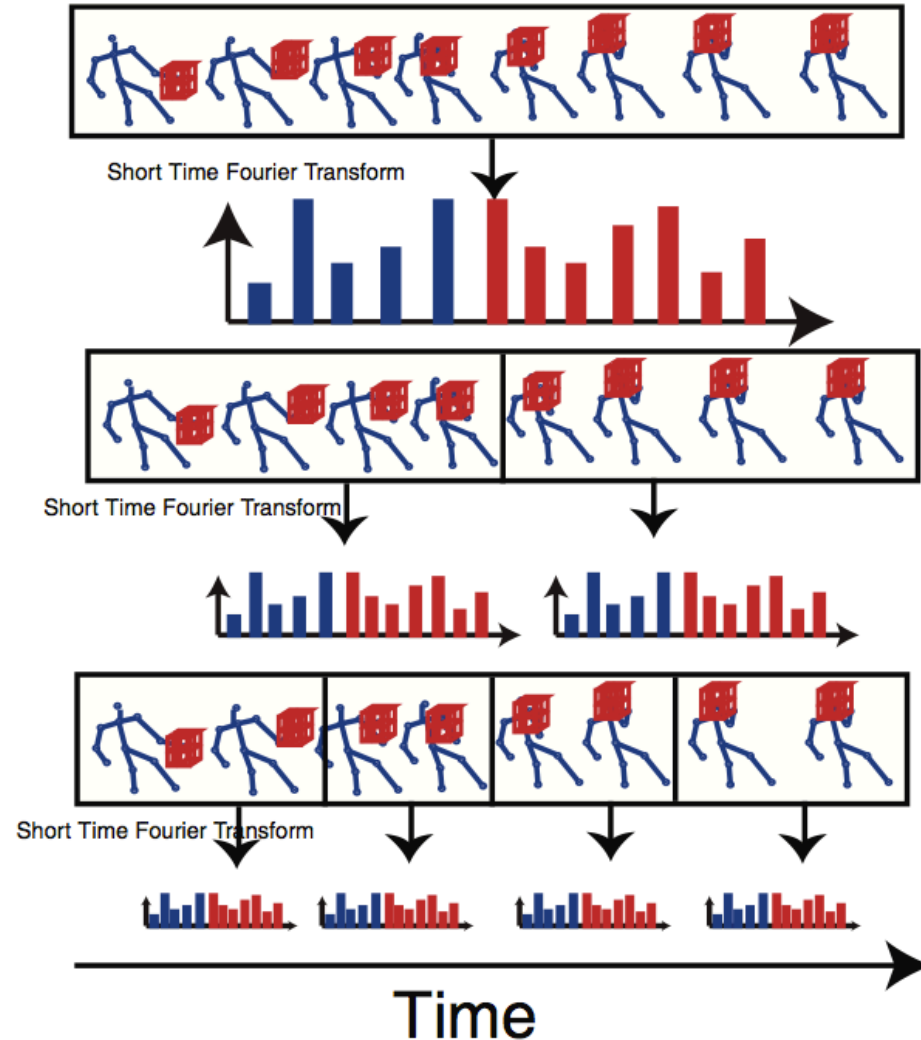For periodic motion, techniques like DTW produce large misalignment

Idea: recursively partition pyramid

Apply Short Time Fourier Transform
to feature set   g[t] = [ Pos[t], LOP[t]]

Use low frequency coeffs as features

Feature $G_i$: concatenate all segments

e.g. all 7 sets ->

In results: 3 levels, cutoff=1/4 length

# Actionlet ensemble (1/4)

Correlation of body parts is important!

actionlet := conjunctive (AND) structure on base features

denoted by $S = \{1,2,\dots,N_j\}$

base feature := a fourier pyramid of one joint

Learn discriminatively which joints should be modeled together

Use AND/OR structure:

Prediction (y) is labeled (c) iff all joint features ($x_j$) are labeled c

$$P_S(y^{(j)} = c | \boldsymbol{x}^{(j)}) = \prod_{i \in S} P_i(y^{(j)} = c | \boldsymbol{x}^{(j)})$$

Define $\mathcal{X}_c$ as $\{j : t^{(j)} = c\}$

# Actionlet ensemble (2/4)

Maximize confidence

Minimize ambiguity

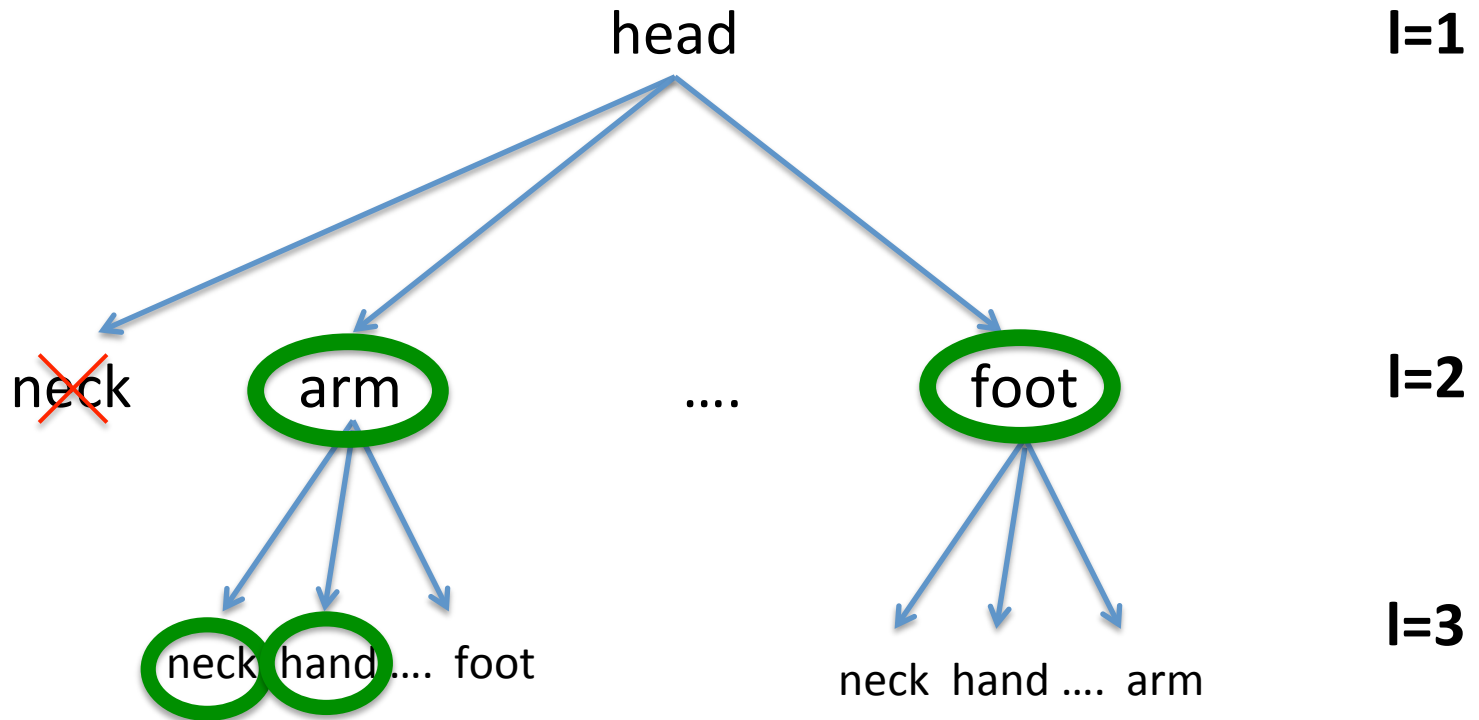$$\text{Conf}_S = \max_{j \in \mathcal{X}_c} \log P_S(y^{(j)} = c | \boldsymbol{x}^{(j)})$$

$$\text{Amb}_S = \sum_{j \notin \mathcal{X}_c} \log P_S(y^{(j)} = c | \boldsymbol{x}^{(j)})$$

$P_c$ := Discriminative actionlet pool for class c

**Algorithm:**

```
1) For c in 1…C:   // Each class
2)      Pc <- {}       // Initialize pool
3)      l <- 1
4)      Do:
5)          Generate l-actionlets by adding one joint into
            each (l-1) actionlet in Pc
6)          Add l-actionlets if conf > Tconf
7)          l++
8)      Until no actionlet is added
9)      Remove actionlets if Ambig > Tamb in Pc
Output: all actionlets that meet the criteria
```

# Actionlet ensemble (3/4)
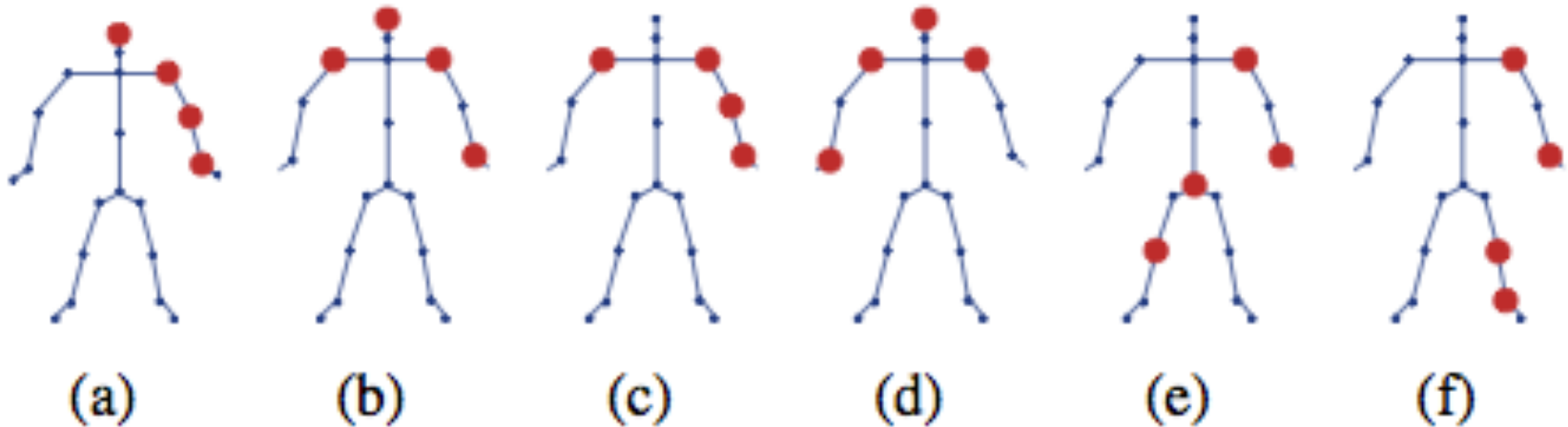
# Actionlet ensemble (4/4)



Figure 9. Examples of the mined actionlets. The joints contained in each actionlet are marked as red. (a), (b) are actionlets for "drink" (c), (d) are actionlets for "call". (e), (f) are actionlets for "walk".

# Multiple Kernel Learning

- Multiclass-MKL
  - One versus all

kernel weight      kernel

$$f_{\text{final}}(\boldsymbol{x}, y) = \sum_{k=1}^{p} [\beta_k \langle \boldsymbol{w}_k, \Phi_k(\boldsymbol{x}, y) \rangle + b_k]$$

L1 regularizer on beta so small number of actionlets are learned

$$: \Omega(\boldsymbol{\beta}) = \|\boldsymbol{\beta}\|_1^2$$

$$\min_{\boldsymbol{\beta}, \boldsymbol{w}, \boldsymbol{b}, \xi} \frac{1}{2} \Omega(\boldsymbol{\beta}) + C \sum_{i=1}^{n} \xi_i$$

$$\text{s.t. } \forall i : \xi_i = \max_{u \neq y_i} l(f_{\text{final}}(\boldsymbol{x}^{(i)}, y^{(i)}) - f_{\text{final}}(\boldsymbol{x}^{(i)}, u))$$

Solve by iteratively:

   1) optimizing beta with fixed w, b with linear programming

   2) optimization w,b with fixed beta using generic SVM solver
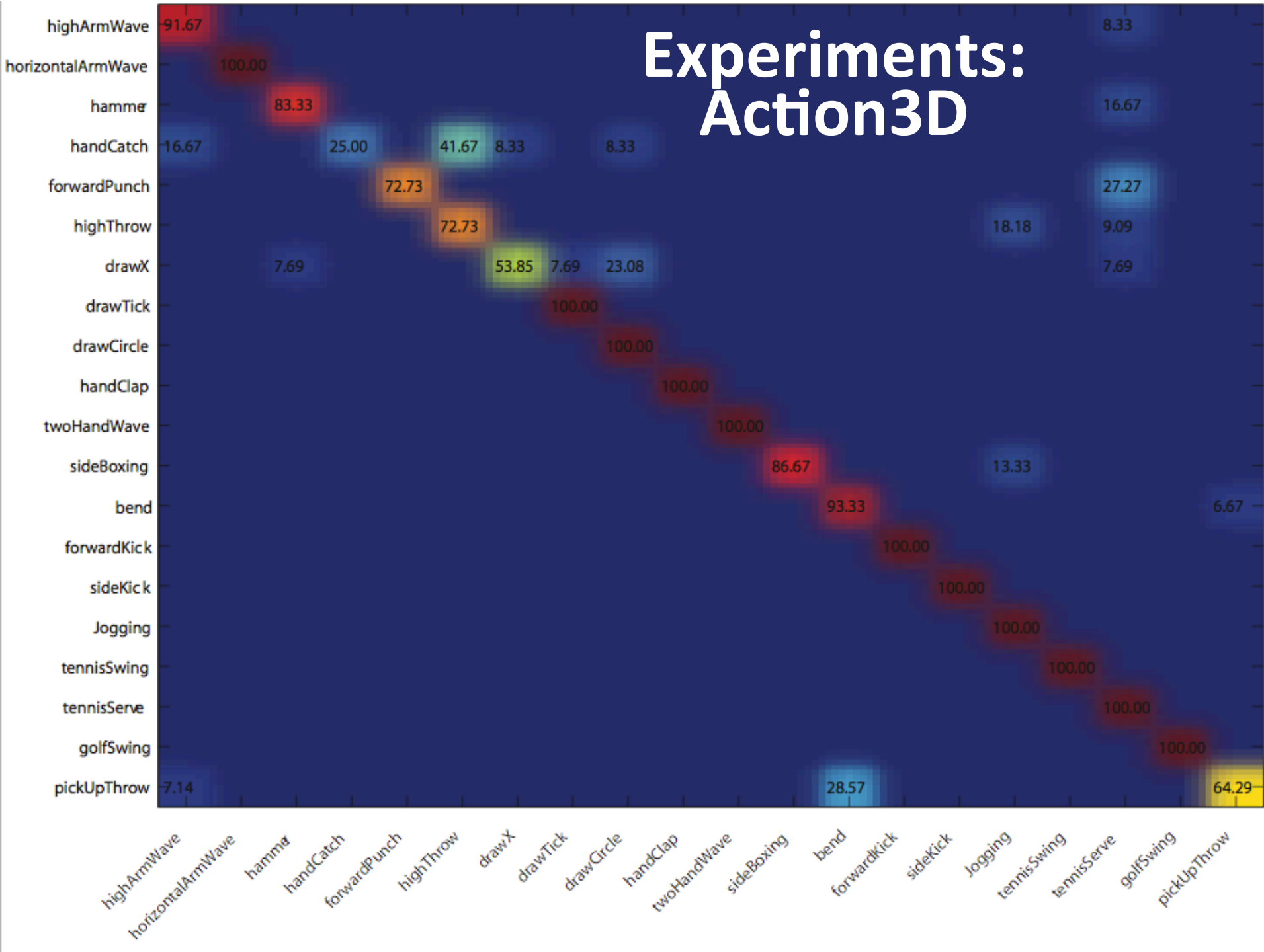
# MSR Action3D

Half of subjects used as training data

Errors when actions are too similar ("hand catch" and "high throw")

Other models thrown by temporal misalignment

| Method | Accuracy |
|---|---|
| Recurrent Neural Network [16] | 0.425 |
| Dynamic Temporal Warping [17] | 0.54 |
| Hidden Markov Model [15] | 0.63 |
| Action Graph on Bag of 3D Points [14] | 0.747 |
| **Proposed Method** | **0.882** |

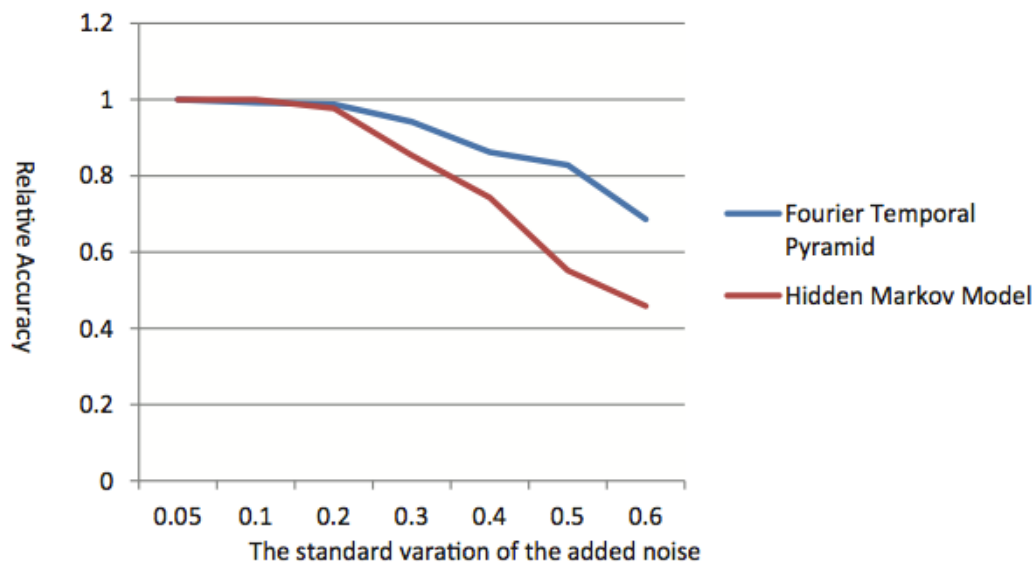Table 1. Recognition Accuracy Comparison for MSR-Action3D dataset.
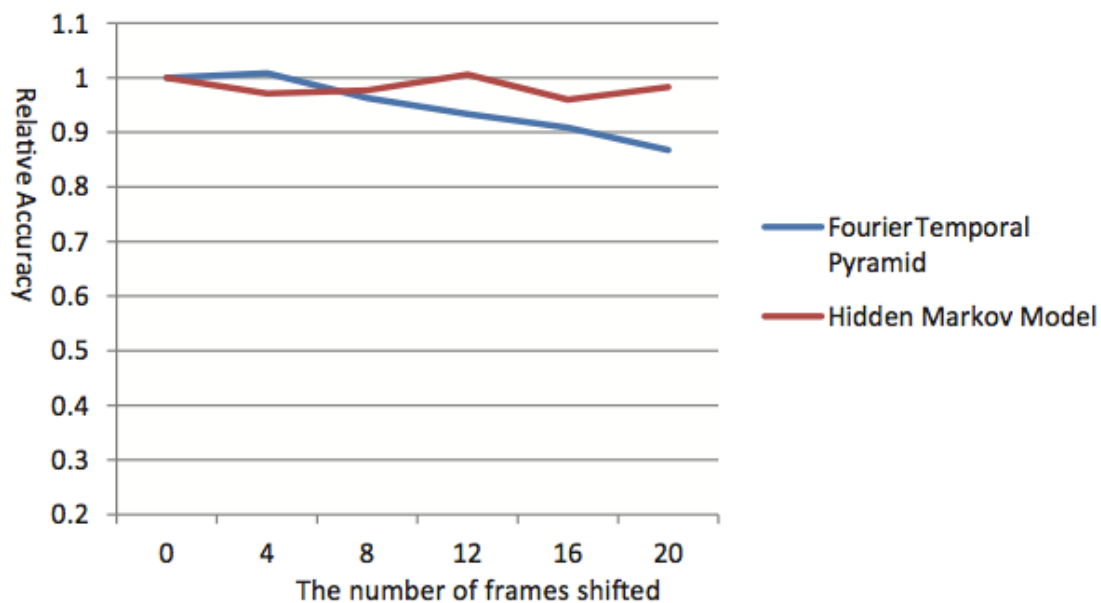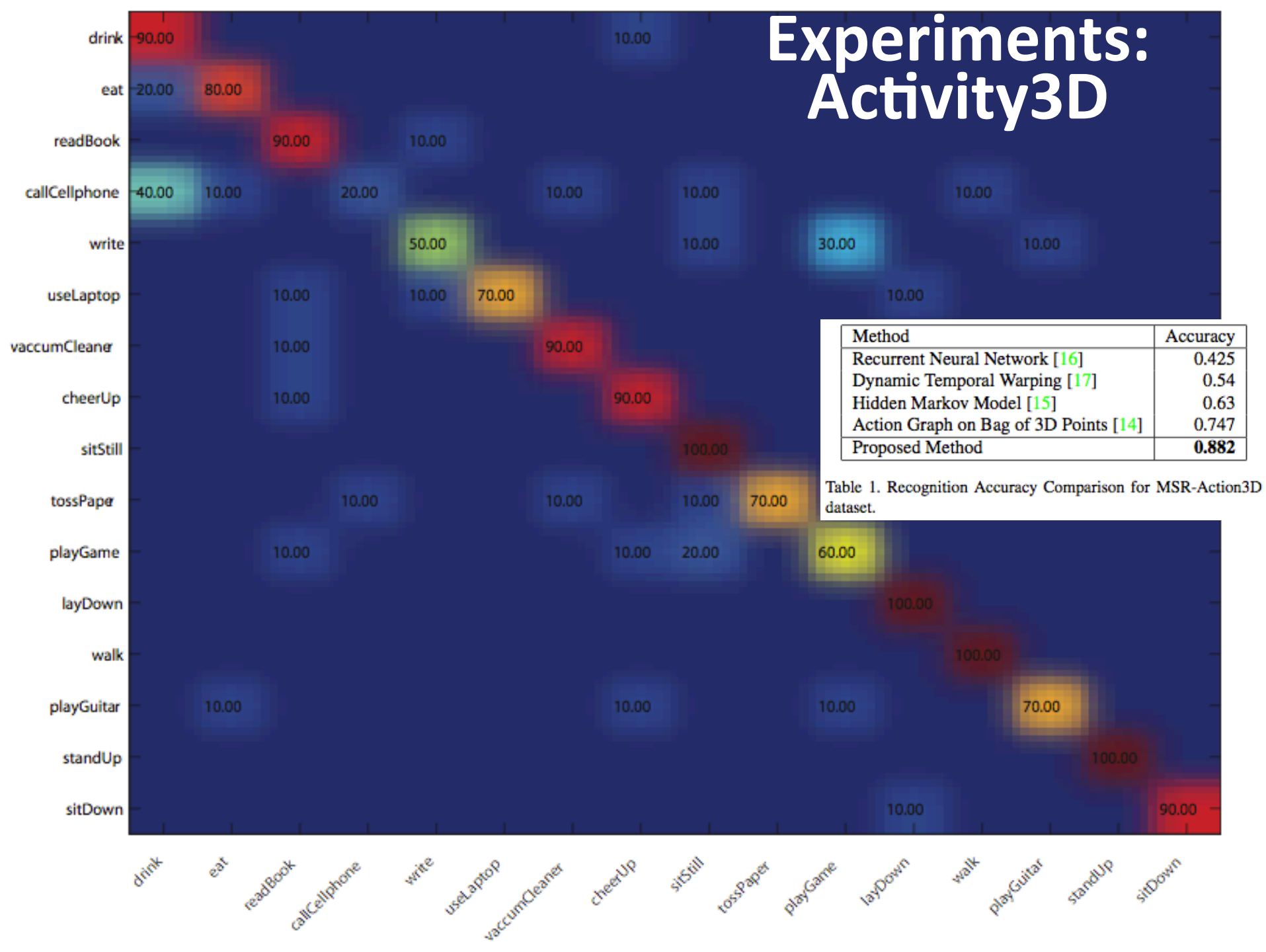
Experiments: Action3D

# Sensitivity Action3D

+Gaussian noise to joints



Temporal shift

# Experiments: Activity3D

| Method | Accuracy |
|---|---|
| Recurrent Neural Network [16] | 0.425 |
| Dynamic Temporal Warping [17] | 0.54 |
| Hidden Markov Model [15] | 0.63 |
| Action Graph on Bag of 3D Points [14] | 0.747 |
| Proposed Method | **0.882** |

Table 1. Recognition Accuracy Comparison for MSR-Action3D dataset.
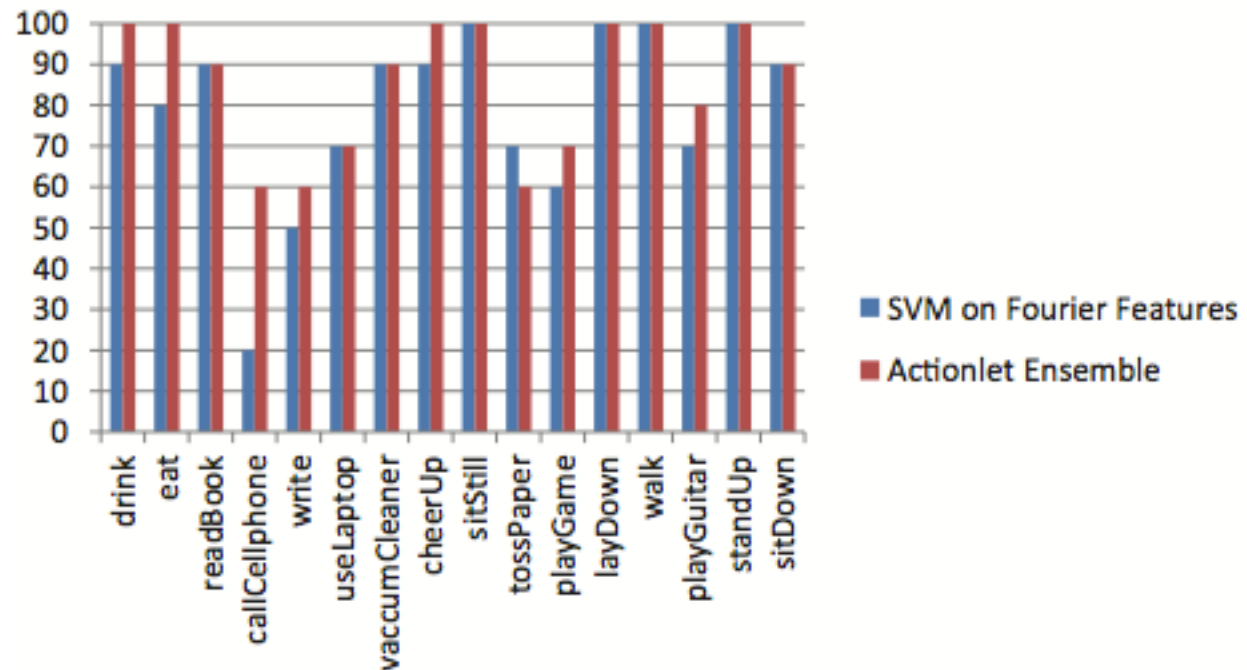
# Experiments: DailyActivities



Figure 8. The comparison between the accuracy of the proposed actionlet ensemble method and that of the support vector machine on the Fourier Temporal Pyramid features.

# Experiments: CMU Mocap

Data is much cleaner than from the Kinect

| Method | Accuracy |
|---|---|
| CRF with learned manifold space [9] | 0.9827 |
| Proposed Method | **0.9813** |

Table 3. Recognition Accuracy Comparison for CMU MoCap dataset.