

Federico Angelini\*, Zeyu Fu\*, Sergio A. Velastin†, Jonathon A. Chambers\*, Syed Mohsen Naqvi\*

\*Intelligent Sensing and Communications Research Group, Newcastle University, UK

†Department of Computer Science and Engineering, University Carlos III de Madrid, Spain

## 1 Overview

Given the high demand for automated systems for human action recognition, great efforts have been undertaken in recent decades to progress the field [1]. In this paper, we present frameworks for single and multi-viewpoints action recognition based on:

- Space-Time Volume (STV) of human silhouettes
- 3D-Histogram of Oriented Gradient (3D-HOG) Embedding [2]



Fig 1. Examples of RGB and Silhouettes data.

### Our contributions

- 3D-HOG Embedding [2] based frameworks exploiting local gestures analysis
- single and multi-viewpoints cases
- accuracy and robustness to appearance changes
- **outperforming** results on Weizmann and i3DPost datasets

## 2 Baseline Method

- Baseline method: 3D-HOG Embedding [2]
- It defines the basic data processing structure (Fig 2), also used in the Proposed Frameworks

### Key drawbacks

- **Attention problem** (Fig 3): it has not been addressed;
- **Performance stability**: affected by *randomly selected library* in the Embedding phase;
- **Action-labels-based local classifiers**, without considering cross-location local gestures relationships;

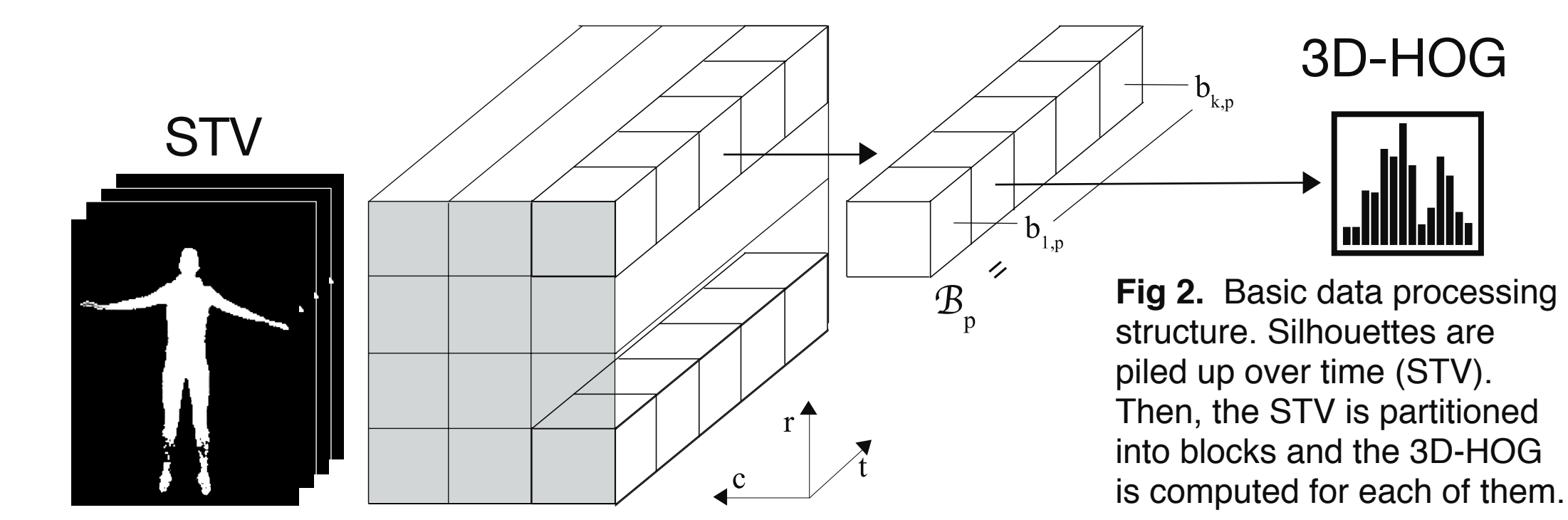


Fig 2. Basic data processing structure. Silhouettes are piled up over time (STV). Then, the STV is partitioned into blocks and the 3D-HOG is computed for each of them.

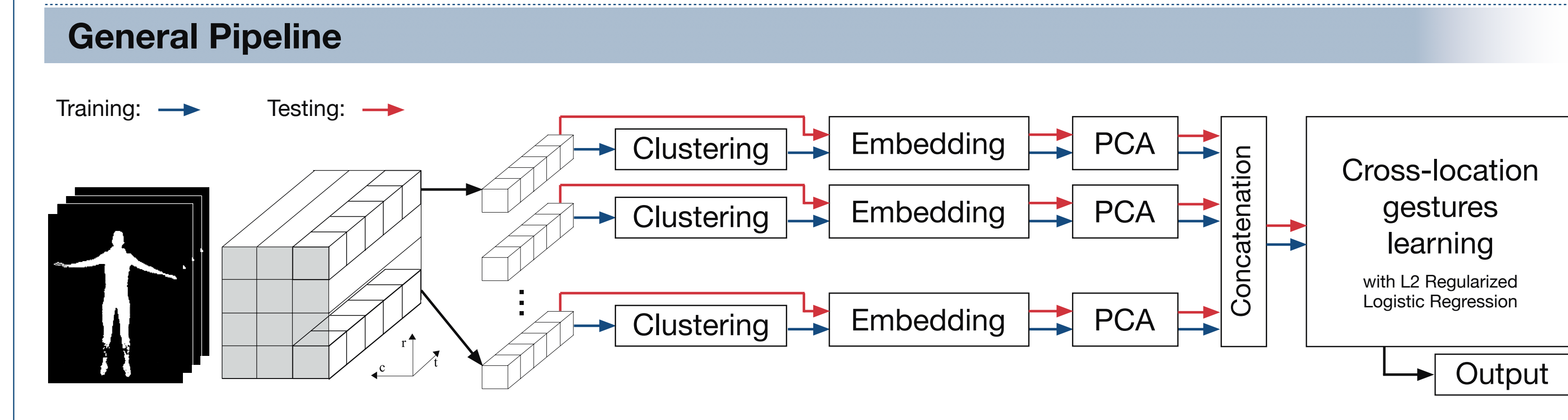
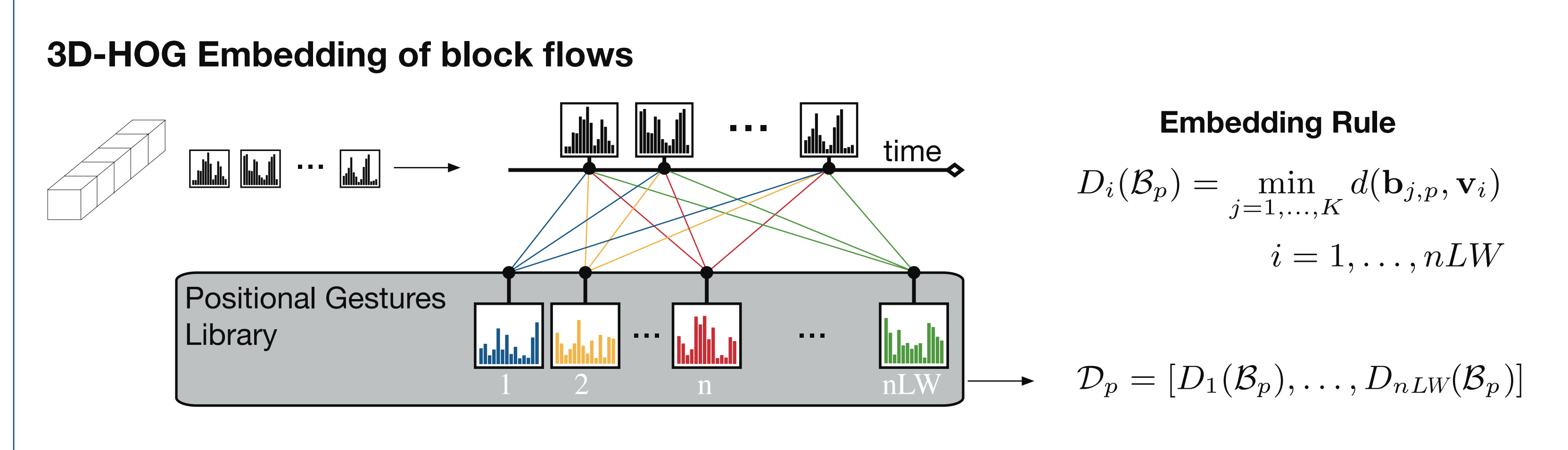
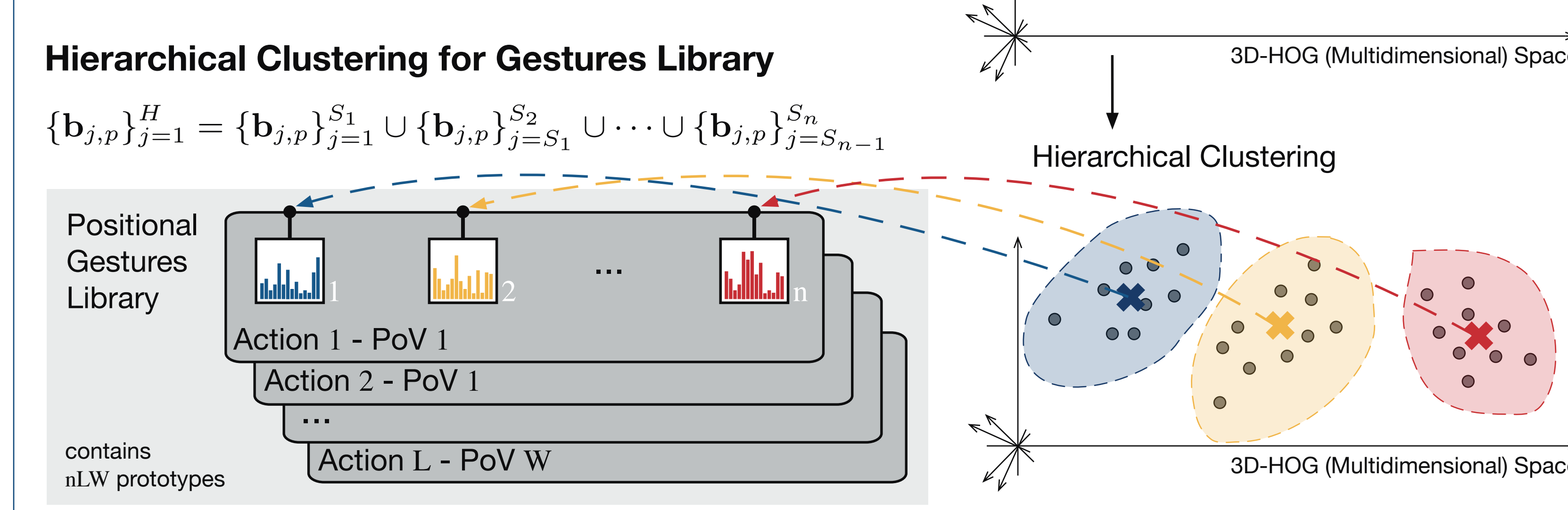
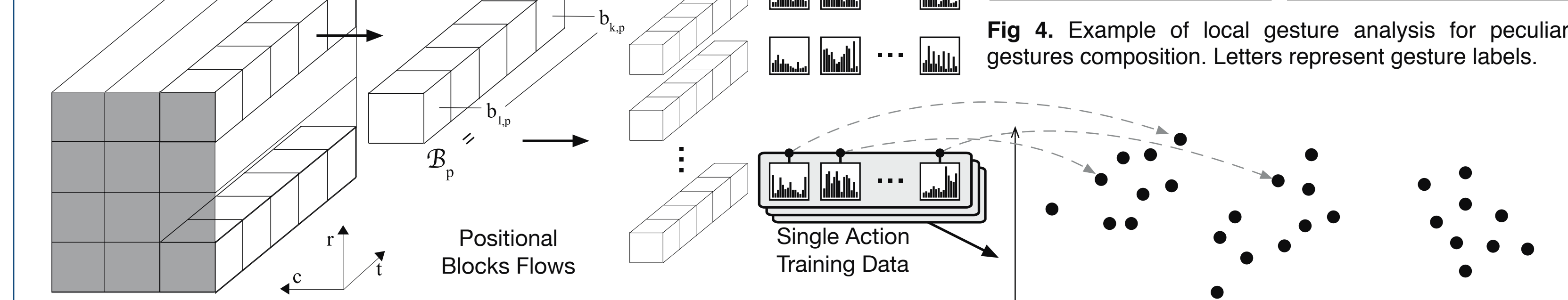
### Example of Attention Problem

■ locations where Action 1 and Action 2 look different  
 ■ locations where Action 1 and Action 2 look similar

## 3 Proposed Frameworks

### Main ideas

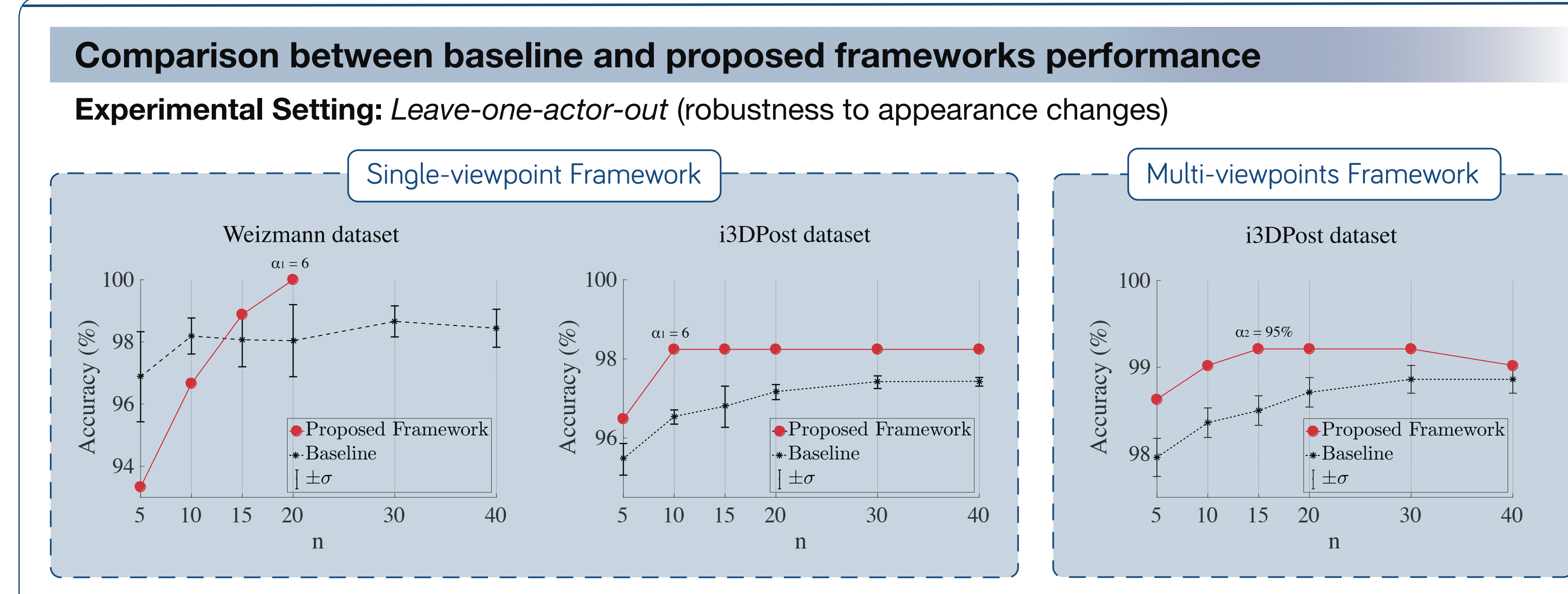
- Locally, actions look like **simpler gestures**
- Globally, an action can be seen as a **particular combination** of local gestures (Fig 4)



### Key facts

- **Gesture analysis**: to leverage local information for global action information retrieval
- **Searching local gestures**: 3D-HOG features clustered with *Hierarchical Clustering*
- No longer randomly chosen **Local (Gestures) Library**
- **Cross-location Gesture Composition Learning**: with *L2 Regularized Logistic Regression*

## 4 Results



### Legend

- n: clustering parameter  $\propto$  library size
- $\sigma$ : standard deviation
- $\alpha_1$ : number of principal components (PCA)
- $\alpha_2$ : explained variance percentage (PCA)

### Conclusions

- **Outperforming results** in all studied cases
- **Stable performance** over different trainings
- Higher accuracy for smaller n (best values)

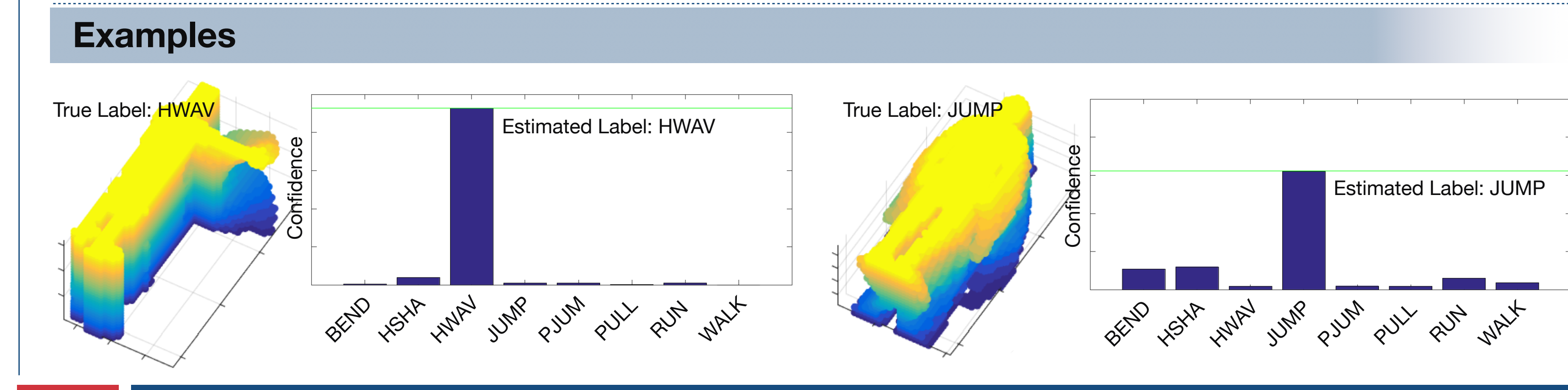
### Comparison between dataset state-of-arts and proposed frameworks performance

**Experimental Setting:** Leave-one-actor-out (robustness to appearance changes)

Weizmann Dataset					i3DPost Dataset					
Method	Actions	Accuracy	n	$\alpha_1$	Method	Actions	POV	Accuracy	n	$\alpha_2$
Proposed Framework	10	<b>100%</b>	20	6	Proposed Framework	8	8	<b>99.60%</b>	30	95%
Gorelick et al.	10	<b>100%</b>	-	-	Castro et al.	6	2	99.00%	-	-
Jiang et al.	10	<b>100%</b>	-	-	Iosifidis et al.	6	8	98.16%	-	-
C. Li et al.	9	97.53%	-	-	Iosifidis et al.	8	8	96.34%	-	-
Ahsan et al.	9	97.5%	-	-	Hilsenbeck et al.	6	8	92.42%	-	-
Ahsan et al.	10	94.26%	-	-						

### Conclusions

**State-of-art results (Weizmann) and outperforming results (i3DPost)**



## 5 References

[1] S. Herath, M. Harandi, F. Porikli, "Going deeper into action recognition: A survey", Image and Vision Computing, vol 60, pp 4-21, 2017.

[2] D. Weinland, M. Özuysal, P. Fua, "Making action recognition robust to occlusions and viewpoint changes", LNCS - Lecture Notes in Computer Science, vol 6313, no. part 3, pp. 635-648, 2010.