

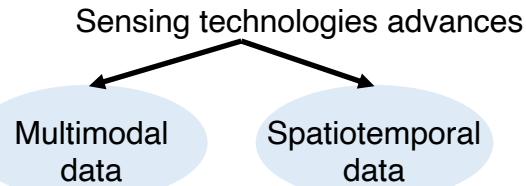
Space-Time Domain Tensor Neural Networks: An Application on Human Activity Classification



L-Università ta' Malta
Institute of Digital Games

Konstantinos Makantasis¹, Athanasios Voulodimos²,
Anastasios Doulamis³, Nikolaos Bakalos³, Nikolaos Doulamis³

INTRODUCTION



The current work has 2 **goals**:

1. Process and fuse multimodal spatiotemporal data
2. Create discriminative data representations for pattern recognition tasks

Application: Classification of human poses using 3D skeleton data

METHOD

Consider N sensors, J modalities and T time instances. Points from j -th modality are represented by a matrix $S_j \in \mathbb{R}^{N \times T}$.

Steps:

1. Project S_j in $\mathbb{R}^{M \times T}$ using $W \in \mathbb{R}^{N \times M}$

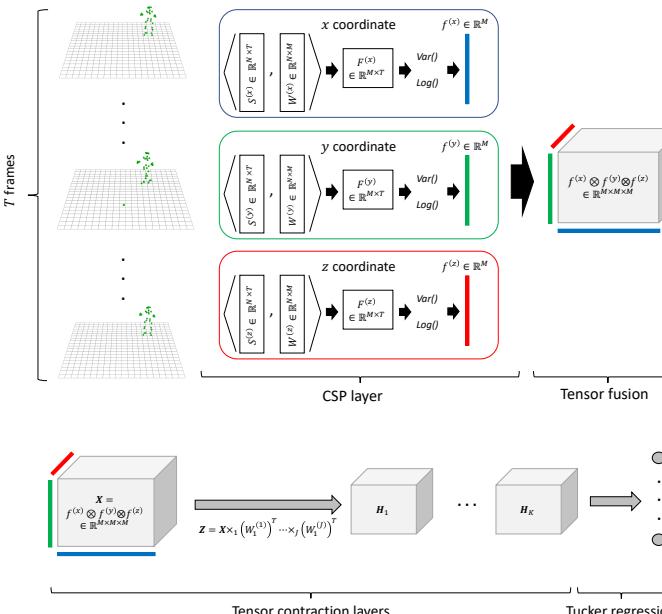
2. Compute feature

$$f_j = \log \left[\frac{\text{var}(W^T S_j^1)}{\sum_{i=1}^M \text{var}(W^T S_j^i)} \dots \frac{\text{var}(W^T S_j^M)}{\sum_{i=1}^M \text{var}(W^T S_j^i)} \right] \in \mathbb{R}^M$$

where $W^T S_j^i$ is the i -th row of $W^T S_j$

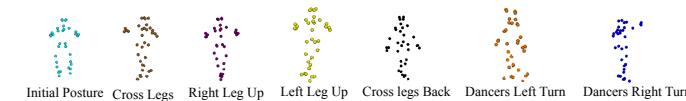
3. Fuse different modalities into a tensor $\mathcal{X} = f_1 \otimes \dots \otimes f_J \in \mathbb{R}^{M \times M \times \dots \times M}$
4. Classify \mathcal{X} using a tensor neural network

MODEL



RESULTS

Dataset:



- Kinect II 3D skeleton data
- 10 sessions, 3 dancers, 4 dances
- 10-fold cross validation

	Accuracy (%)	F1 Score (%)
LSTM	84.2%	82.0%
BOBi LSTM*	85.4%	80.7%
1D-CNN	91.1%	89.7%
Our approach	91.6%	90.9%

DISCUSSION

- The network is *end-to-end-trainable*
- State-of-the-art results with
 - 87 times less parameters than BOBi LSTM
 - 2 times less parameters than 1D-CNN

*I. Rallis et al. "Learning choreographic primitives through a bayesian optimized bi-directional lstm model," in IEEE ICIP, 2019, pp. 1940–1944.

¹ University of Malta, Malta

² University of West Attica, Greece

³ National Technical University of Athens, Greece