

Gradient-weighted Class Activation Mapping for spatio temporal graph convolutional network

Pratyusha Das, Antonio Ortega
University of Southern California (USA)

Introduction

- Motivation:**
 - Spatio-temporal Graph Convolutional Neural Network.
 - Captures simultaneously the spatial correlation and the temporal pattern in the data
 - Lacks interpretability
 - Unknown reason behind their prediction
- Contribution:**
 - Spatio-temporal Graph Grad CAM
 - Gradient based class activation maps for STGCN

Preliminaries

STGCN [1]

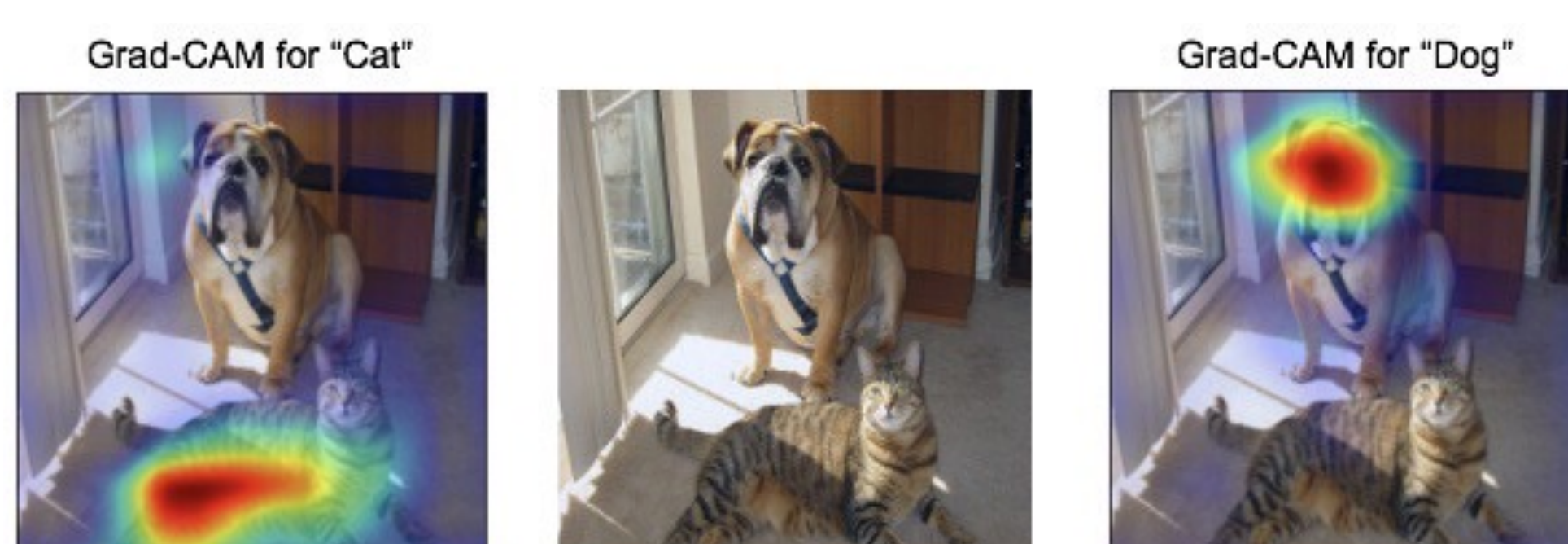
In each layer, ST-GCN is computed as

$$x_{out} = (\tilde{A} \odot Q)x_{in}W$$

- x_{in} : input feature map of size (C, V, T)
- C : Channel, V : Vertices, T : Temporal length.
- x_{out} : output feature map.
- \tilde{A} : adjacency matrix of spatial graph.
- Q : edge weight matrix of the spatial graph.
- W : stacked weight vector of the multiple output channel.

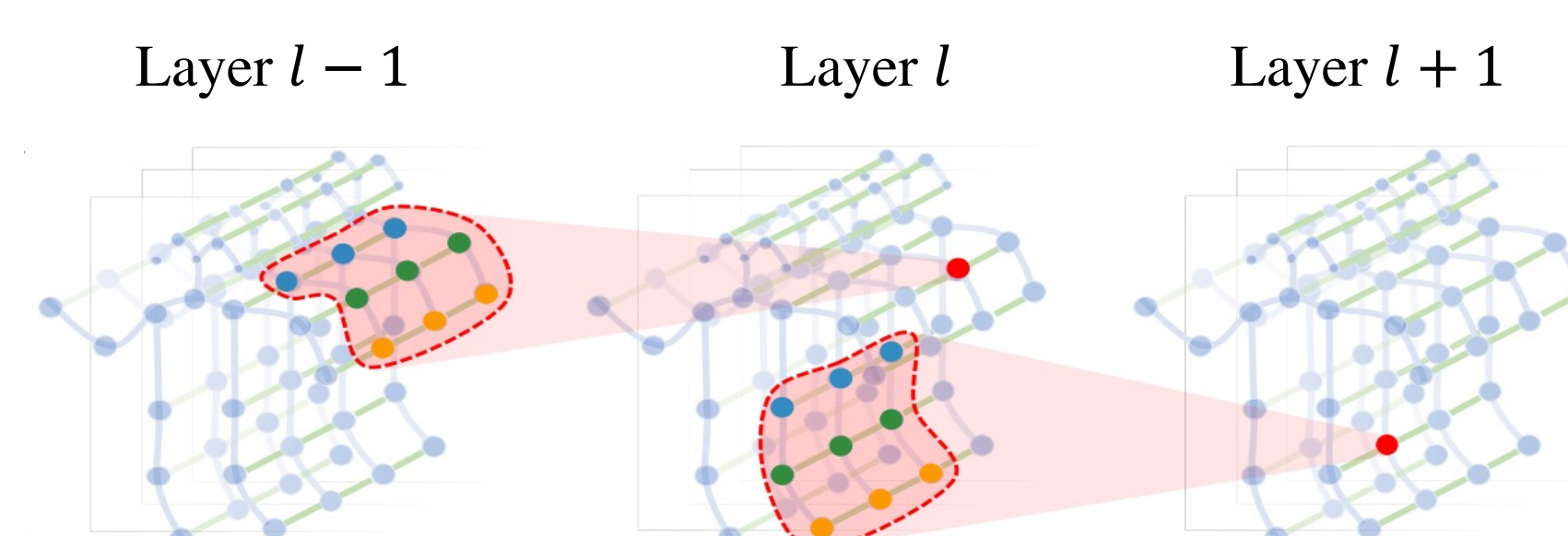
Grad-CAM [2]

- CNN uses 2-D spatial convolutional filter.
- Gradient based class activation maps (Grad-CAM) detects the spatial image important region for the network



Spatio-emporal graph Grad CAM

- Adjacency matrix as filter
 $((A + I)x)_i = \sum_{j \in N(i)} w_{i,j}x_j + w_{i,i}x_i$
- The spatio-temporal GCNN is computed as
 $F_k^l(X, A) = \sigma(\tilde{A}F^{l-1}(X, A)W_k^l)$
- Here, normalized adjacency matrix \tilde{A} , F_k^l represents the k -th feature at the l -th layer.
- Features are **localized spatially** and **temporally**.



- The heatmap $H_{ST}^{c,l}$ for joint-time importance for c^{th} class and l^{th} layer is computed using

$$H_{ST}^{c,l} = \text{ReLU}\left(\sum_k \alpha_k^{c,l} F_k^l\right)$$
- $\alpha_k^{c,l}$: Weights for k -th feature, F_k^l : k -th feature map.
- ReLU considers the positive value which contributed to the final decision
- Weights are calculated based on the gradients

$$\alpha_k^{c,l} = \frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \frac{\partial y^c}{\partial F_{k,n,t}^l}$$
- Here, y^c represents c -th class score.

Evaluation

Faithfulness : Measures the impact of occlusions to the graph nodes

$$\beta_{faithful} = \sum_i m_i \left(\frac{|\alpha_0 - \alpha_{m_i}|}{\alpha_0} \right)$$

masking percentage m , corresponding accuracy α_{m_i} ,

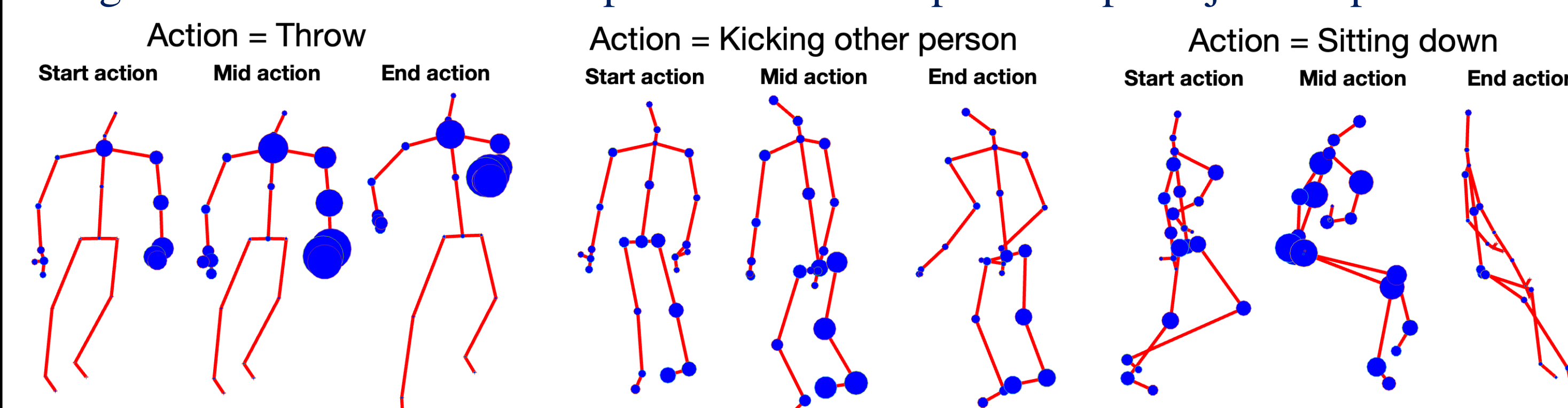
α_0 is accuracy without masking

- Achieved $\beta_{faithful} = 89\%$

Experiments and Results

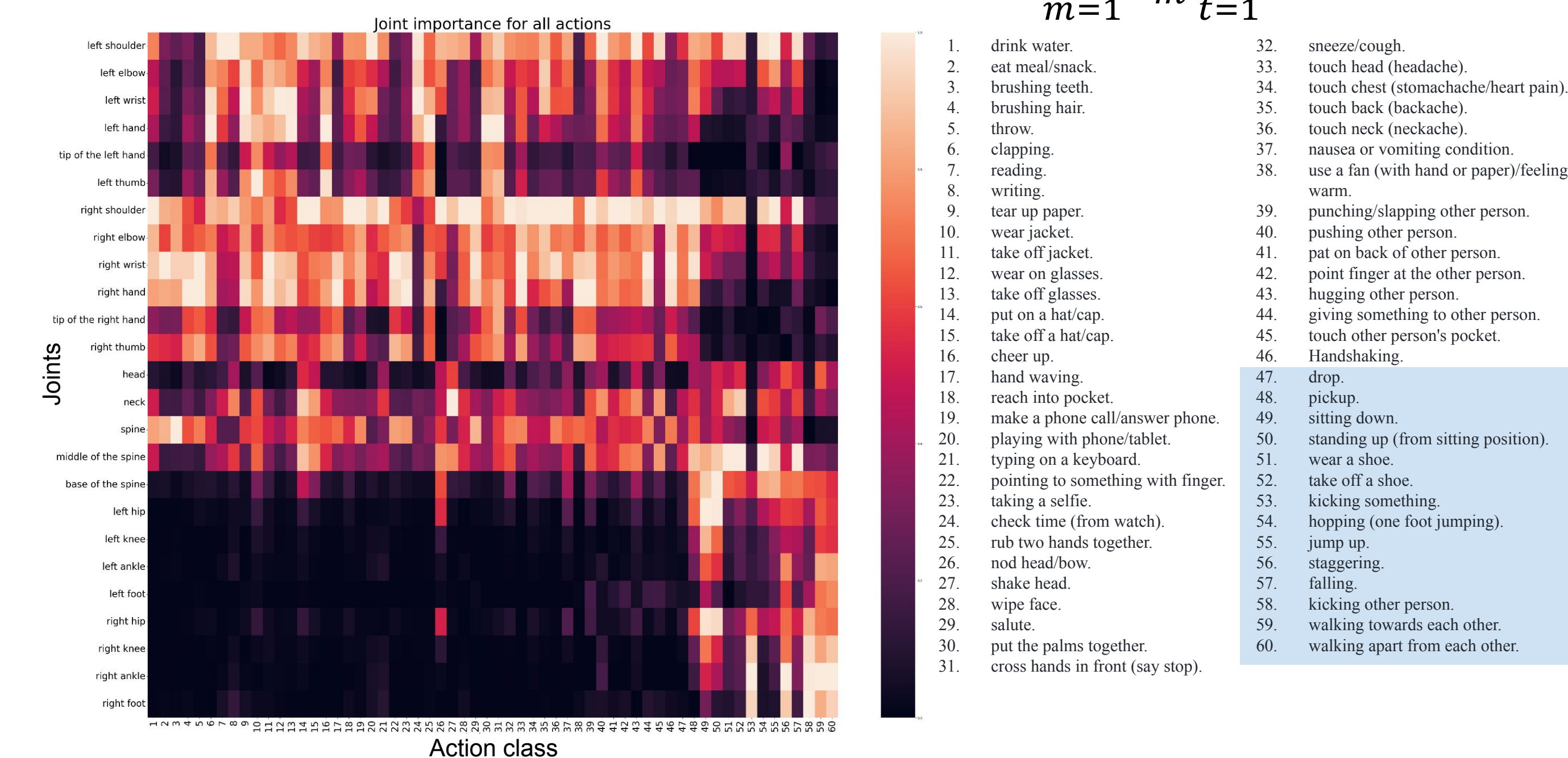
- Task: Skeleton based activity recognition**
- Model: STGCN** 10 layers of STGCN, Trained on NTU-RGB dataset, No of classes 60, Accuracy 81.5%

Fig. Actions at different time points with their spatio temporal joint importance



*Bigger size of the node represents higher importance of that joint.

$$\text{Joint importance map} \rightarrow \psi_n^c = \frac{1}{M^c} \sum_{m=1}^{M^c} \frac{1}{T_m^c} \sum_{t=1}^{T_m^c} H_{ST}^{c,L,n,t}$$



Accuracy achieved by STGCN for different amount of masking of the data of the joints based on their importance

Node type		Important nodes			Non-important nodes		
Masking (in %)	No masking	10%	50%	90%	10%	50%	90%
Cross subject -Accuracy (%)	81.5	80.3	64.52	20.59	80.3	72.5	66.5
Cross Camera view -Accuracy (%)	88.3	88.1	64.31	14.82	88.1	81.09	66.81

Future work

- Interpretability of STGCN
- Generic method – change the graph to implement if for other applications.