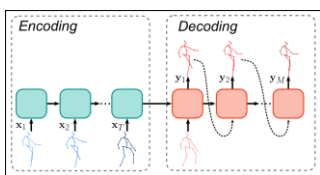


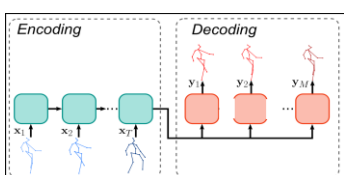
## Motivations

- Sequence-to-sequence motion prediction
- AR decoding: potential error accumulation
- AR decoding: computationally costly

### AR decoding



### Parallel decoding



## Contributions

- Parallel decoding: efficient inference
- Non-autoregressive Transformer architecture
- Single model: activity and motion prediction

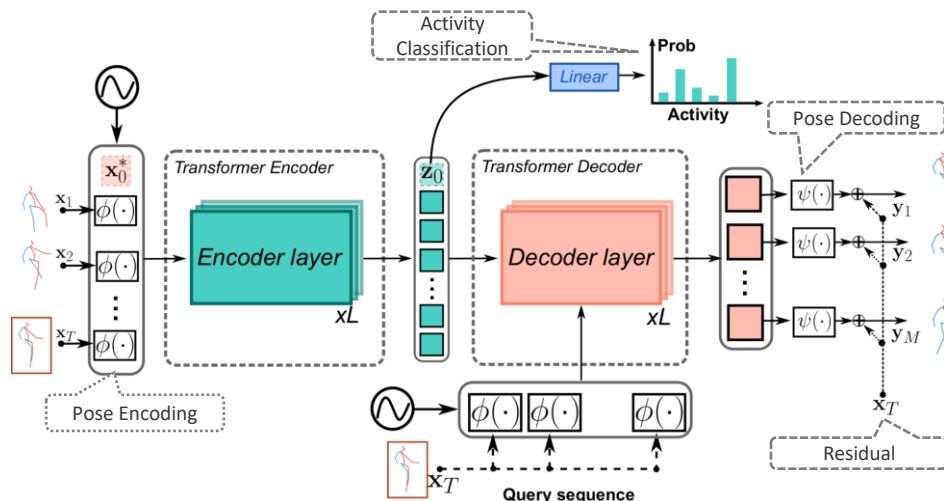
## Related Work

- [1] Martínez et al, On Human motion prediction with recurrent neural networks, CVPR 17
- [2] Carion et al, Detection Transformers, ECCV 2020
- [3] Dosovitskiy et al, Vision Transformers, ICLR 2021
- [4] Mao et al, History repeats itself, ECCV 2020
- [5] Gu et al, Non-Autoregressive Machine Translation, ICRL 2018

## Acknowledgements



## Method

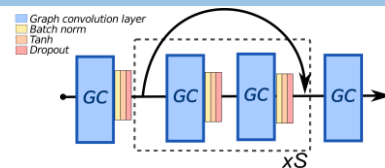


- Query sequence: repeat  $x_T$
- Residual and parallel motion decoding
- Activity token – encode activity from motion
- Speed – **Non-AR: 149.2 SPS**; AR: 8.9 SPS

## Pose Encoding & Decoding

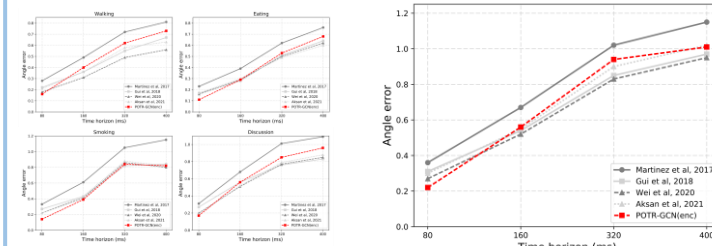
Investigated architectures

- $\phi$  and  $\psi$  are linear layers
- Graph Convolutional Network



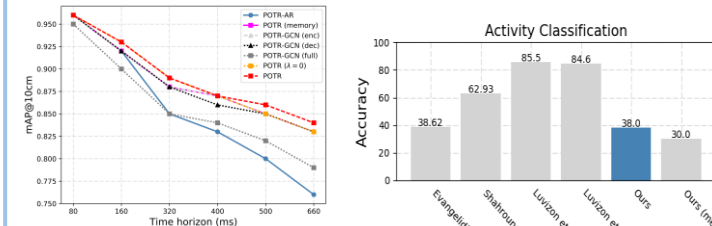
Code & Models  
<https://github.com/idiap/potr>

## H36M – Motion



- Angle error – lower is better
- Best in shorter horizons

## NTURGB – Motion & Activity



- mAP – higher is better
- Skeleton-based classification – low accuracy

## Attention Visualization

