

COMPOSER: Compositional Reasoning of Group Activity in Videos with Keypoint-Only Modality

Honglu Zhou¹, Asim Kadav², Aviv Shamsian³, Shijie Geng¹, Farley Lai², Long Zhao⁴, Ting Liu⁴, Mubbasis Kapadia¹, and Hans Peter Graf²¹ Rutgers University & ² NEC Laboratories America, Inc. & ⁴ Google Research & ³ Bar-Ilan UniversityGitHub: <https://github.com/hongluzhou/composer>

Google Research

Introduction

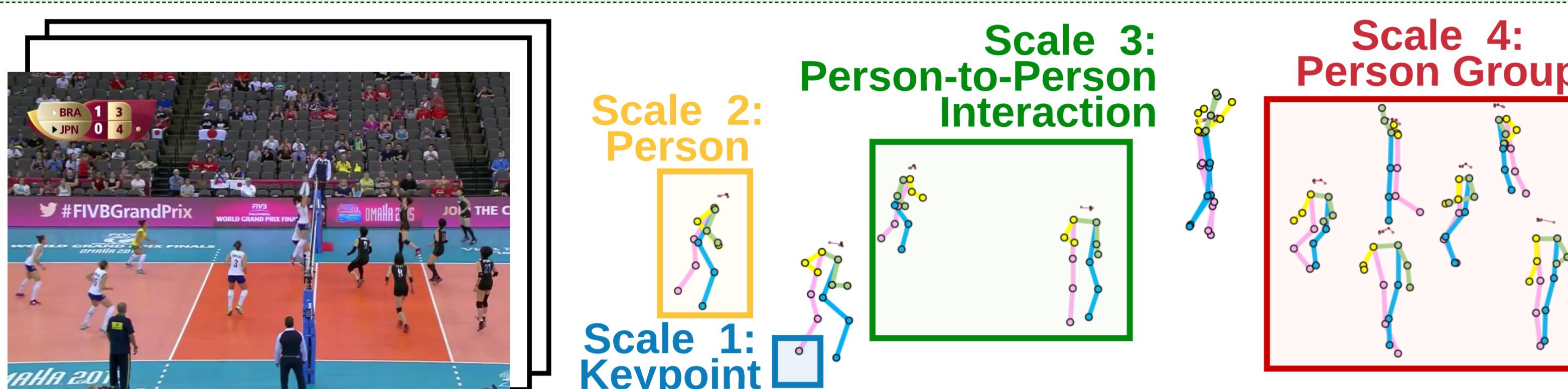
Group Activity Recognition (GAR) detects the activity collectively performed by a group of actors in a short video.

Motivations:

1. GAR requires compositional reasoning of actors and objects.
2. Prior works suffer from scene biases with privacy and ethical concerns.

Our Solutions:

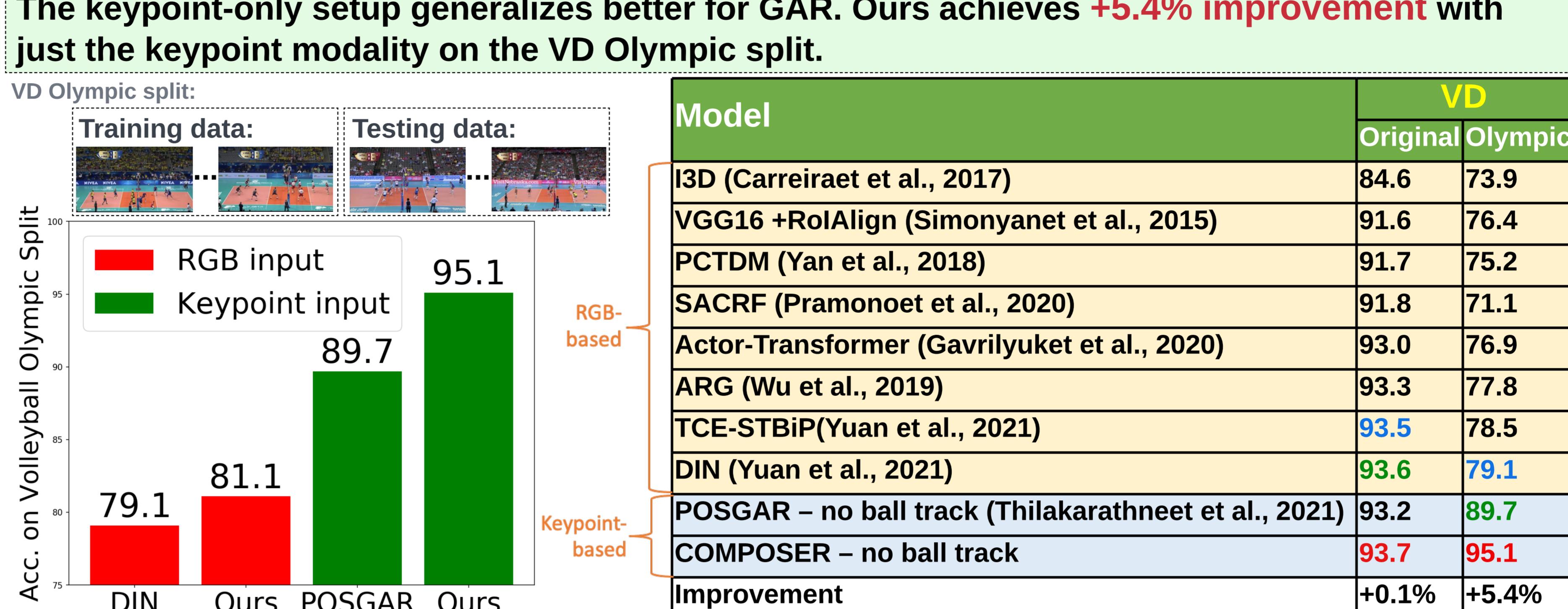
1. Model the video as tokens that represent the multi-scale semantic concepts.
2. Use only the keypoint modality.



Results

Dataset: VD (Volleyball) & CAD (Collective Activity)
Metric: Accuracy

Best, Second, Third



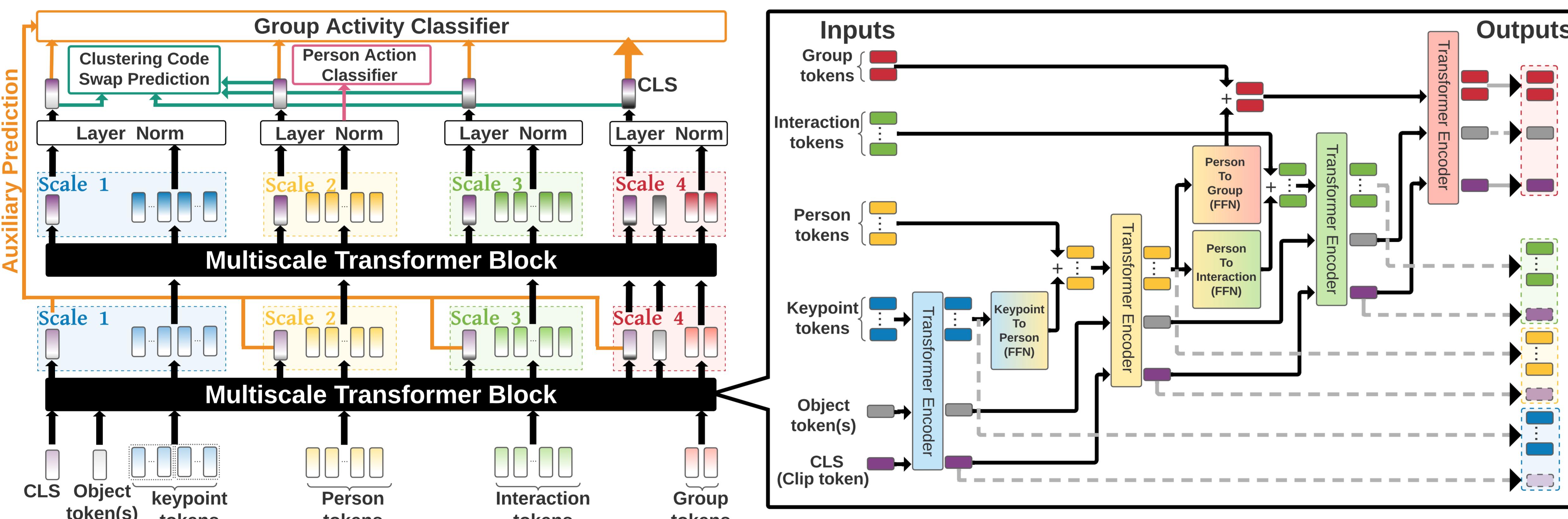
COMPOSER outperforms the GAR methods that use a single modality, and performs favorably compared with methods that exploit multiple expensive modalities (COMPOSER use less computations!).

Model	VD	CAD
	Original	Olympic
I3D (Carreira et al., 2017)	84.6	73.9
VGG16 +RoIAlign (Simonyan et al., 2015)	91.6	76.4
PCTDM (Yan et al., 2018)	91.7	75.2
SACRF (Pramonoet et al., 2020)	91.8	71.1
Actor-Transformer (Gavrilyuk et al., 2020)	93.0	76.9
ARG (Wu et al., 2019)	93.3	77.8
HiGCIN	91.5	93.4
Ehsanpour et al.	93.1	89.4
DIN	93.6	N/A
Actor-Transformer	93.0	92.8
SACRF	95.0	95.2
GroupFormer	95.7	96.3
COMPOSER (ours)	94.6	96.2

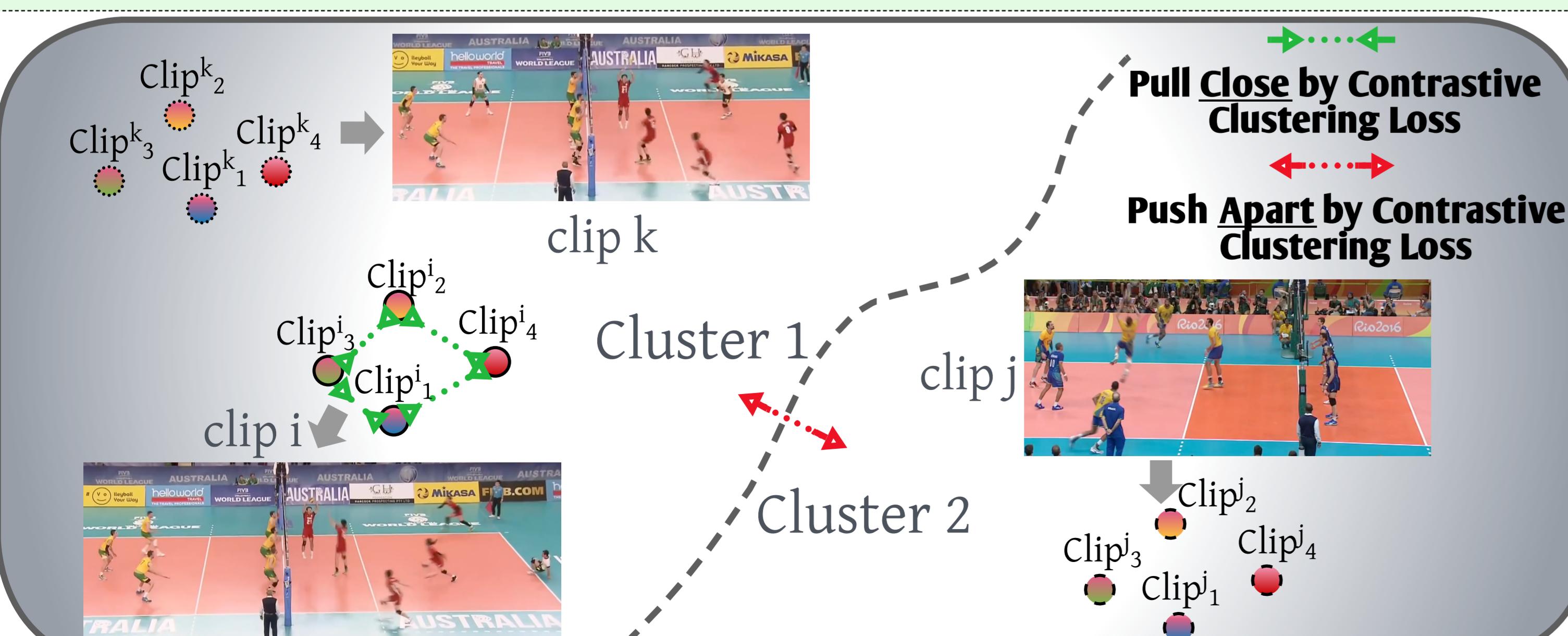
Method

Research Question: How to learn multi-scale representations to aid GAR using the keypoint-only modality?

Novelty 1: COMPOSER is proposed, a **Multiscale Transformer** based architecture that performs attention-based reasoning over tokens at each scale and learns group activity compositionally.



Novelty 2: Contrastive Multiscale Representation Learning, achieved by **clustering** the intermediate scale representations, while **maintaining consistent cluster assignments between scales of the same video clip**.



Novelty 3: Auxiliary Prediction & Data Augmentations, tailored to the keypoint signals.

- Actor Dropout
- Horizontal Flip
- Horizontal Move
- Vertical Move

Data Augmentations

$$\mathcal{L}_{\text{total}} = \sum_{m=1}^{M-1} \mathcal{L}_{\text{groupAux}} + \lambda (\mathcal{L}_{\text{groupLast}} + \mathcal{L}_{\text{person}} + \mathcal{L}_{\text{cluster}})$$

m is the index of the Multiscale Transformer block

Auxiliary Multi-task Learning

Qualitative analysis – showcasing attention matrices. Tokens that the model has mostly attended to at each scale are highlighted.

- COMPOSER is able to attend to relevant information across different scales.
- COMPOSER can produce interpretable results.

Conclusion

Contributions:

1. Compositional and relational reasoning at different scales for GAR.
2. Contrastive clustering, auxiliary prediction, data augmentation techniques to improve the intermediate representations.
3. Extraordinary scene generalization capability using only the keypoint modality.

Future work:

- More complex scenarios, e.g., crowd understanding.
- Use additional modalities like RGB while maintaining scene generalization.

Interested to know more?

Check out the supplementary video!



SCAN ME

