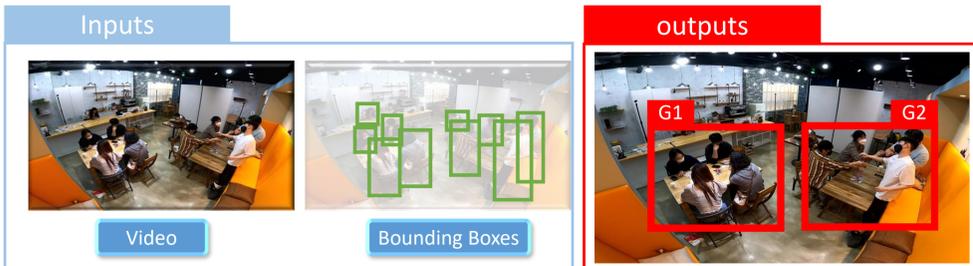


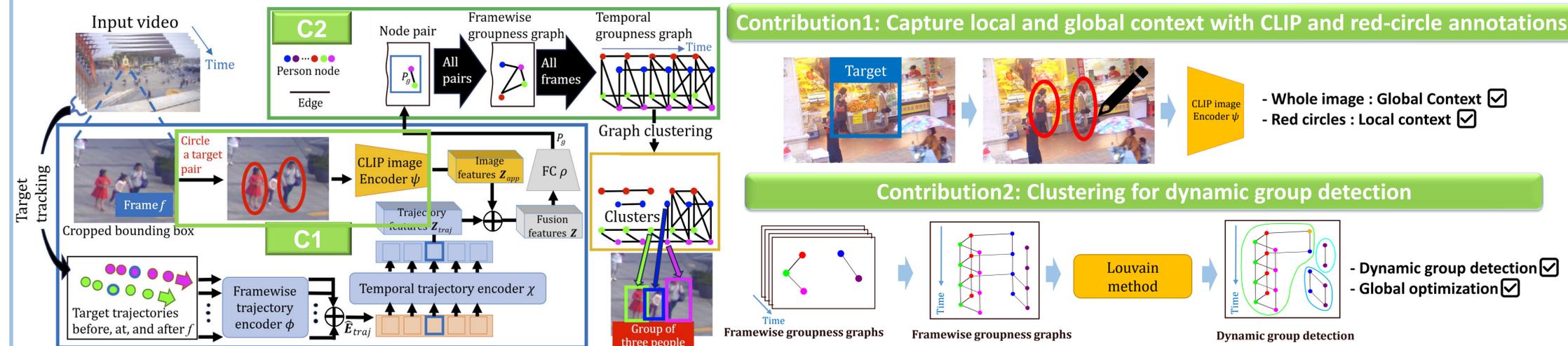


## Task

Detecting the **groups** of people in videos



## Proposed Method



## Limitations of Previous methods

### Limitation 1 : Ignore global context

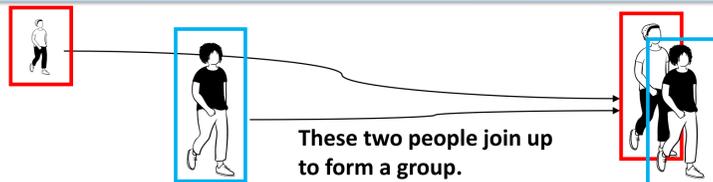
Without Global Context



With Global Context



### Limitation 2 : Cannot detect the temporal changes of groups



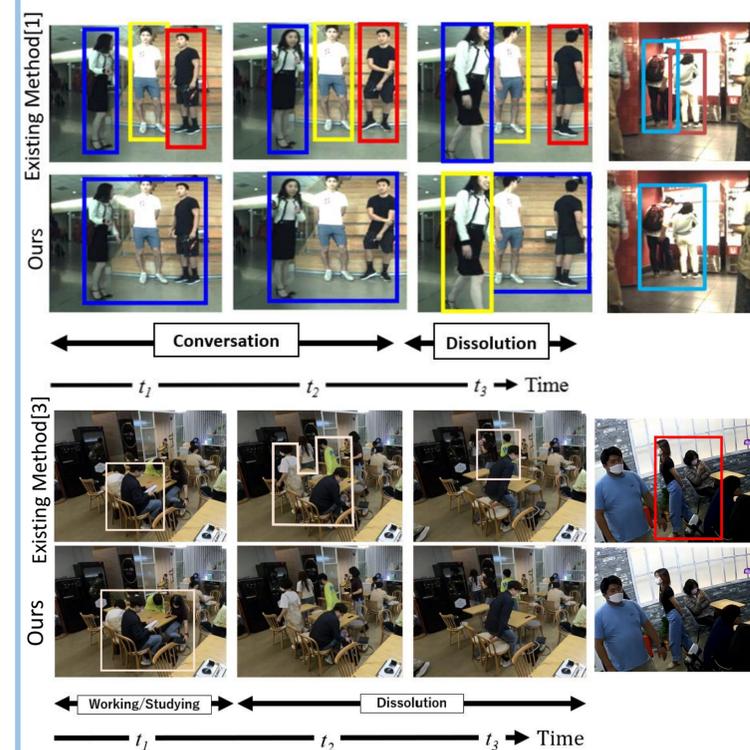
Correct

NOT Correct

Note: People enclosed in bounding boxes of the same color are identified as being in the same group.

## Experiments

### Qualitative results



### References

- J. Zhang et al. Toward Grouping in Large Scenes With Occlusion-Aware Spatio-Temporal Transformers. TCSVT2024.
- R. Han et al. Panoramic Human Activity Recognition. ECCV2022.
- J. Li et al. Self-supervised Social Relation Representation for Human Group Detection. ECCV2022.

### Quantitative results

Dynamic group detection on JRDB

Method	Precision	Recall	F1
S3R2 [3]	0.522	0.662	0.584
P-HAR [2]	0.568	0.697	0.626
GroupTrans [1]	0.514	0.534	0.524
Ours	<b>0.724</b>	<b>0.820</b>	<b>0.769</b>

Dynamic group detection on Cafe

Method	Precision	Recall	F1
S3R2 [3]	0.576	0.700	0.631
GroupTrans [1]	0.263	0.305	0.283
Ours	<b>0.681</b>	<b>0.904</b>	<b>0.776</b>

### Component Analysis: Encoder

Using other image encoders (Trained only on images)

Dataset (Task)	Model	Group detection		
		Precision	Recall	F1
JRDB (dynamic)	ResNet50	0.686	0.609	0.645
	ViT-L/16	0.652	0.666	0.659
	Ours	<b>0.724</b>	<b>0.820</b>	<b>0.769</b>

### Component Analysis: Visual prompts

Using other visual prompts



Dataset (Task)	Model	Group detection		
		Precision	Recall	F1
JRDB (dynamic)	No prompt	0.700	0.689	0.695
	Mask	0.584	0.750	0.657
	A circle	0.640	0.666	0.653
	Two circles	<b>0.724</b>	<b>0.820</b>	<b>0.769</b>

### Limitations

