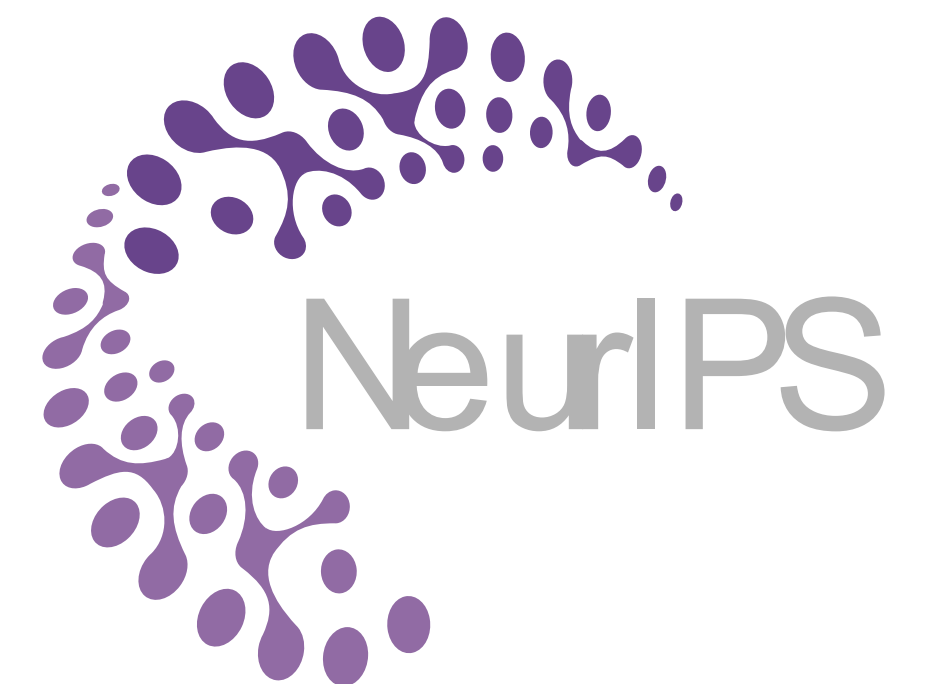


# Part-Aware Bottom-Up Group Reasoning for Fine-Grained Social Interaction Detection

Dongkeun Kim, Minsu Cho, Suha Kwak

Pohang University of Science and Technology (POSTECH)





# Understanding Social Interactions

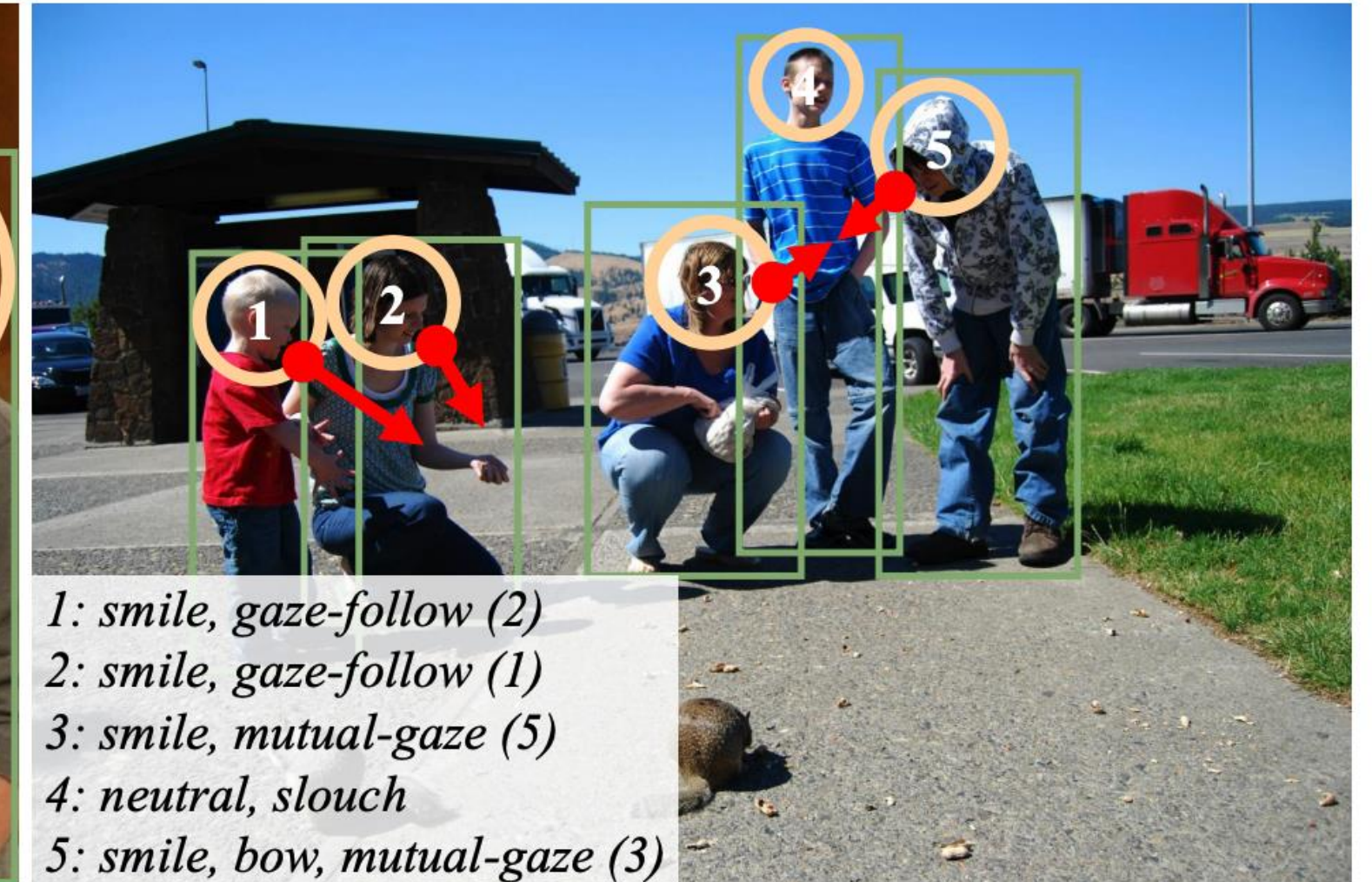
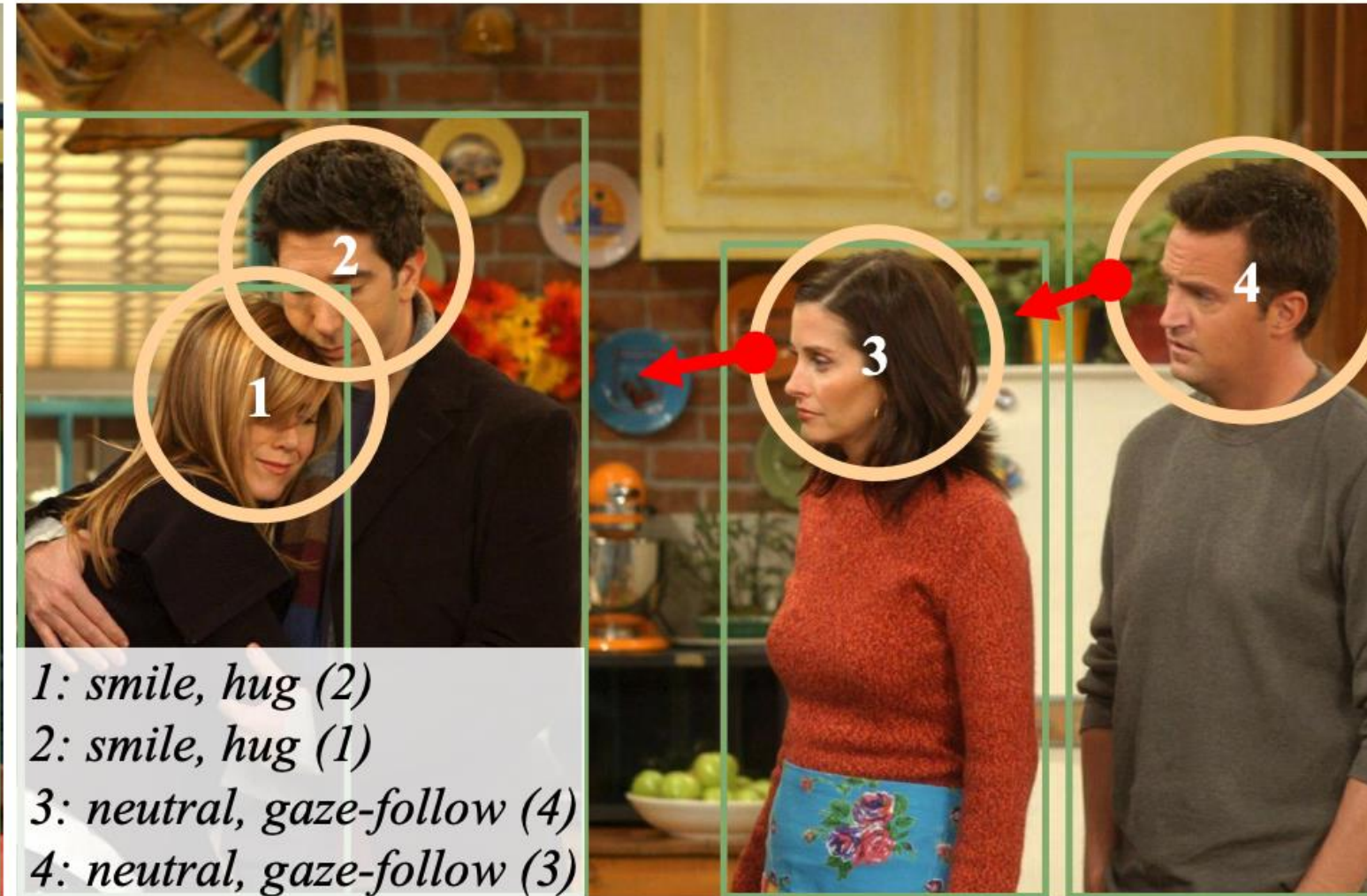
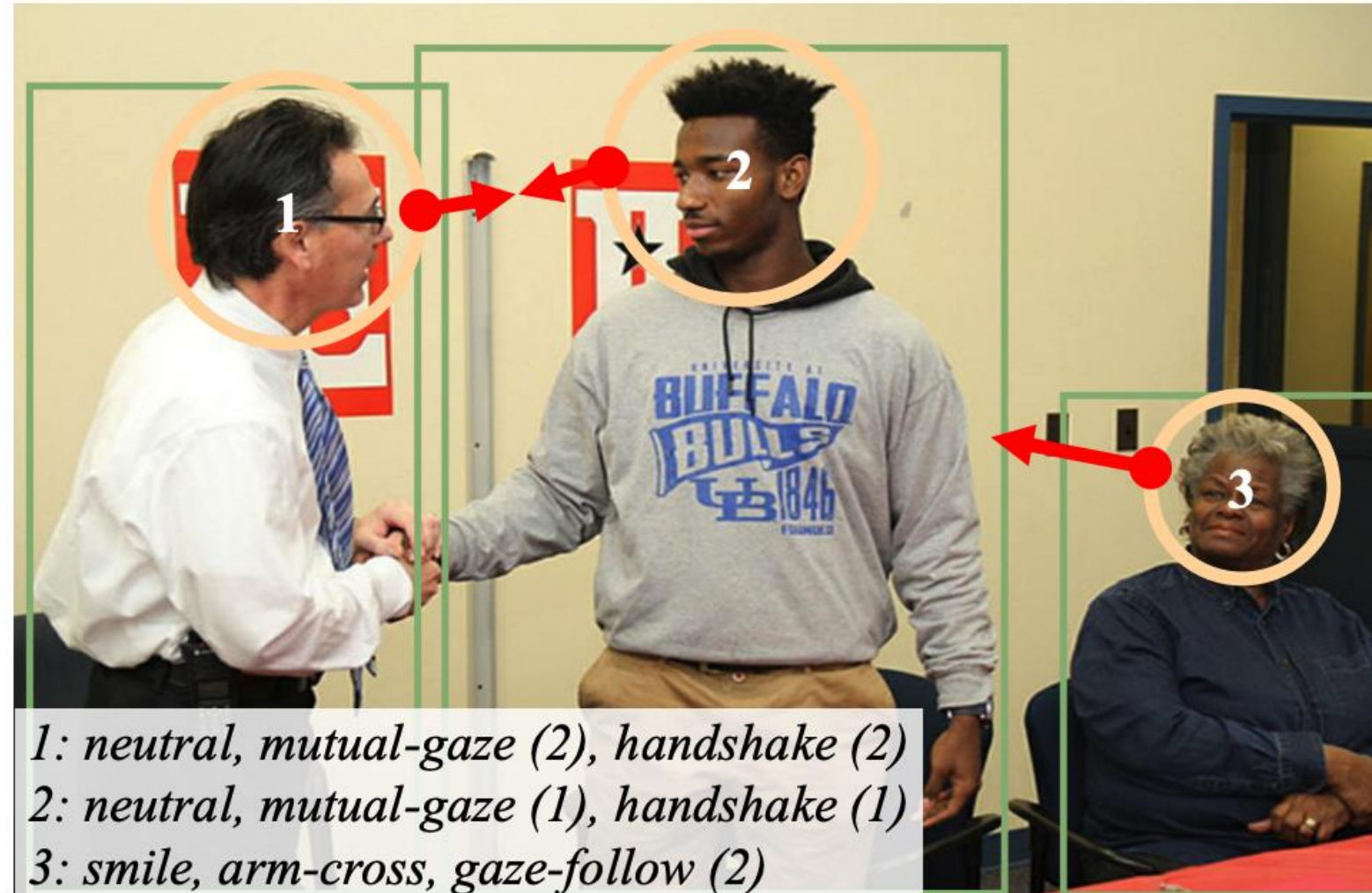
Perceiving subtle visual cues beyond global appearance





# Fine-Grained Social Interaction Detection

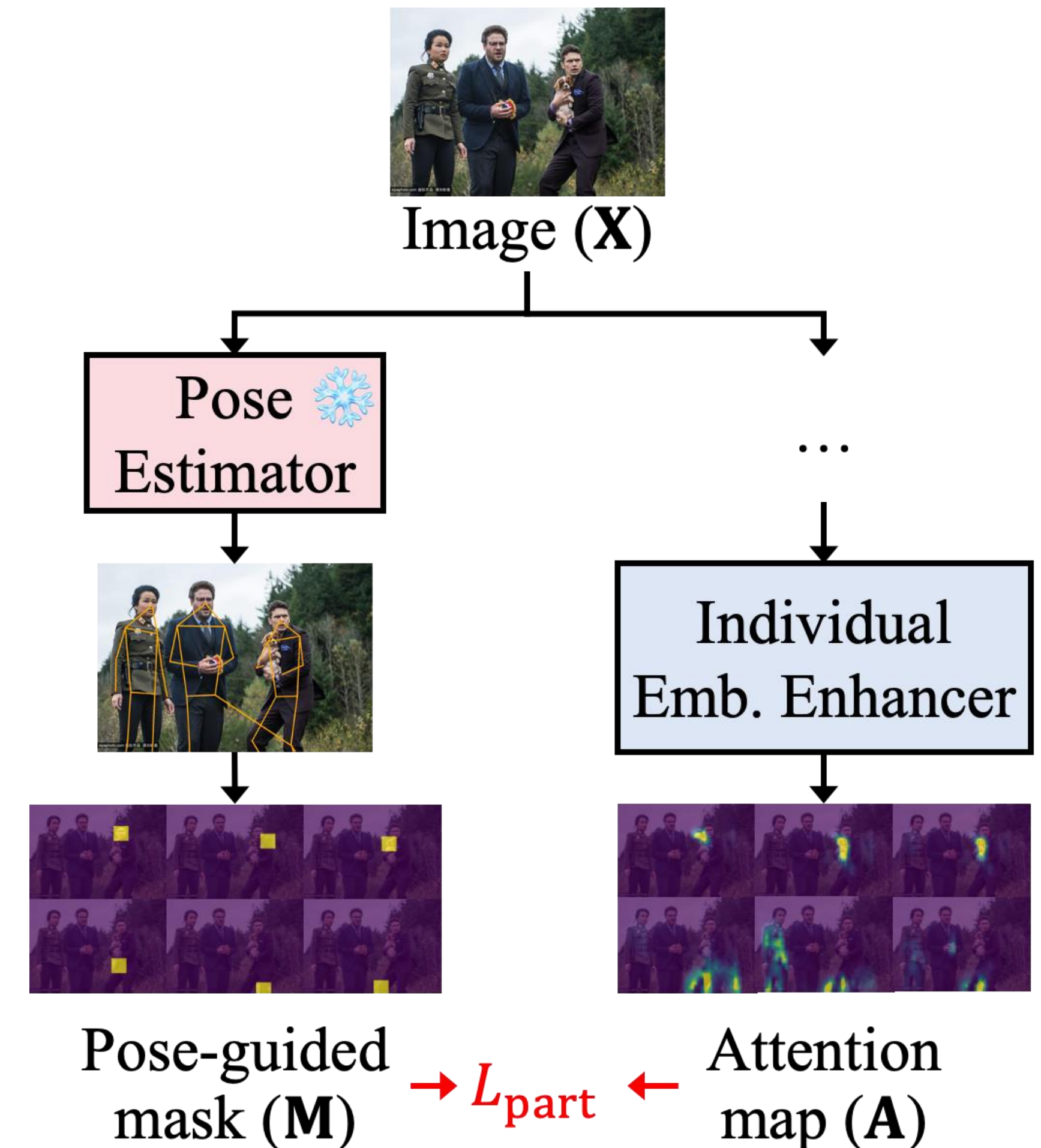
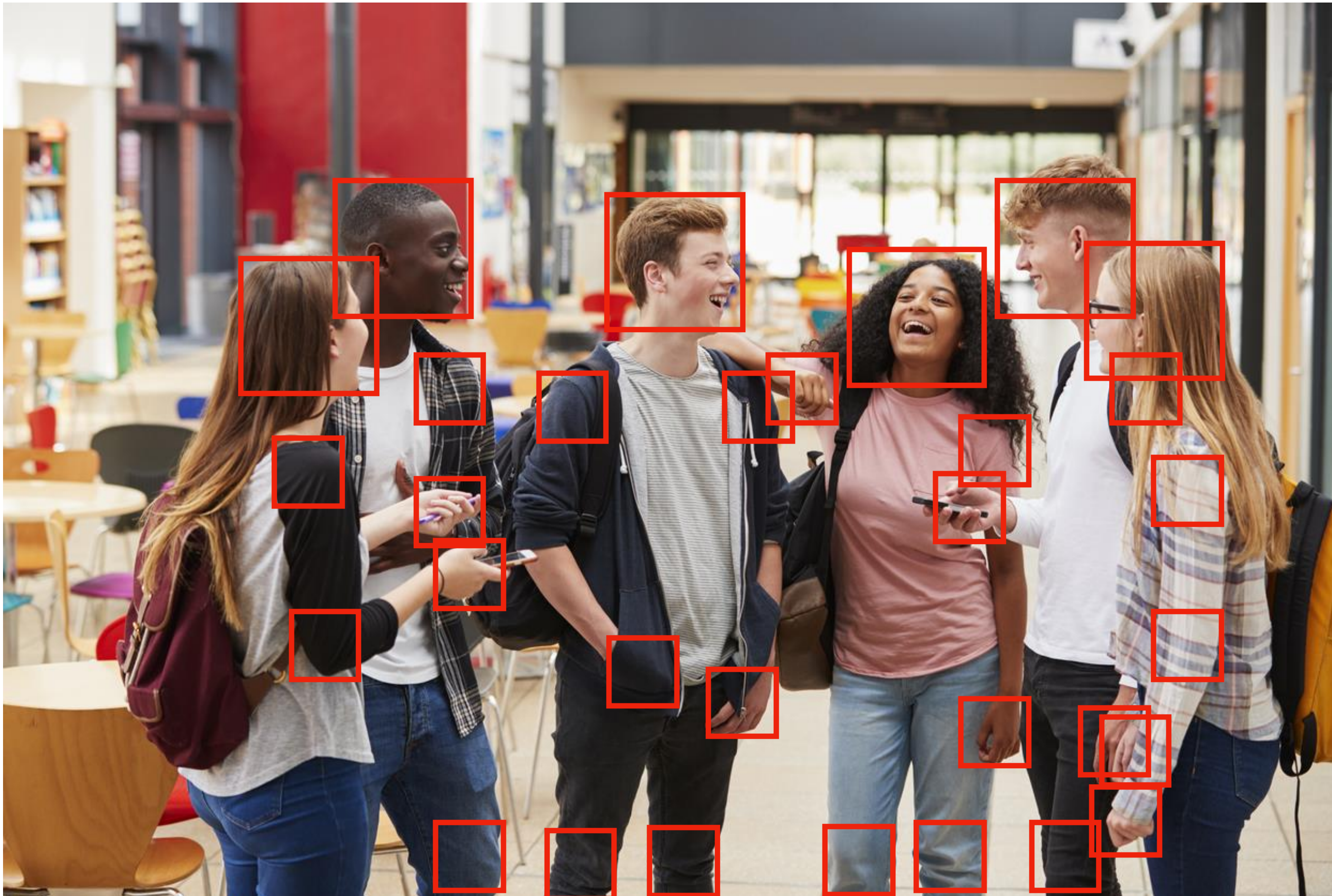
Predict <individual, group, interaction> triplets





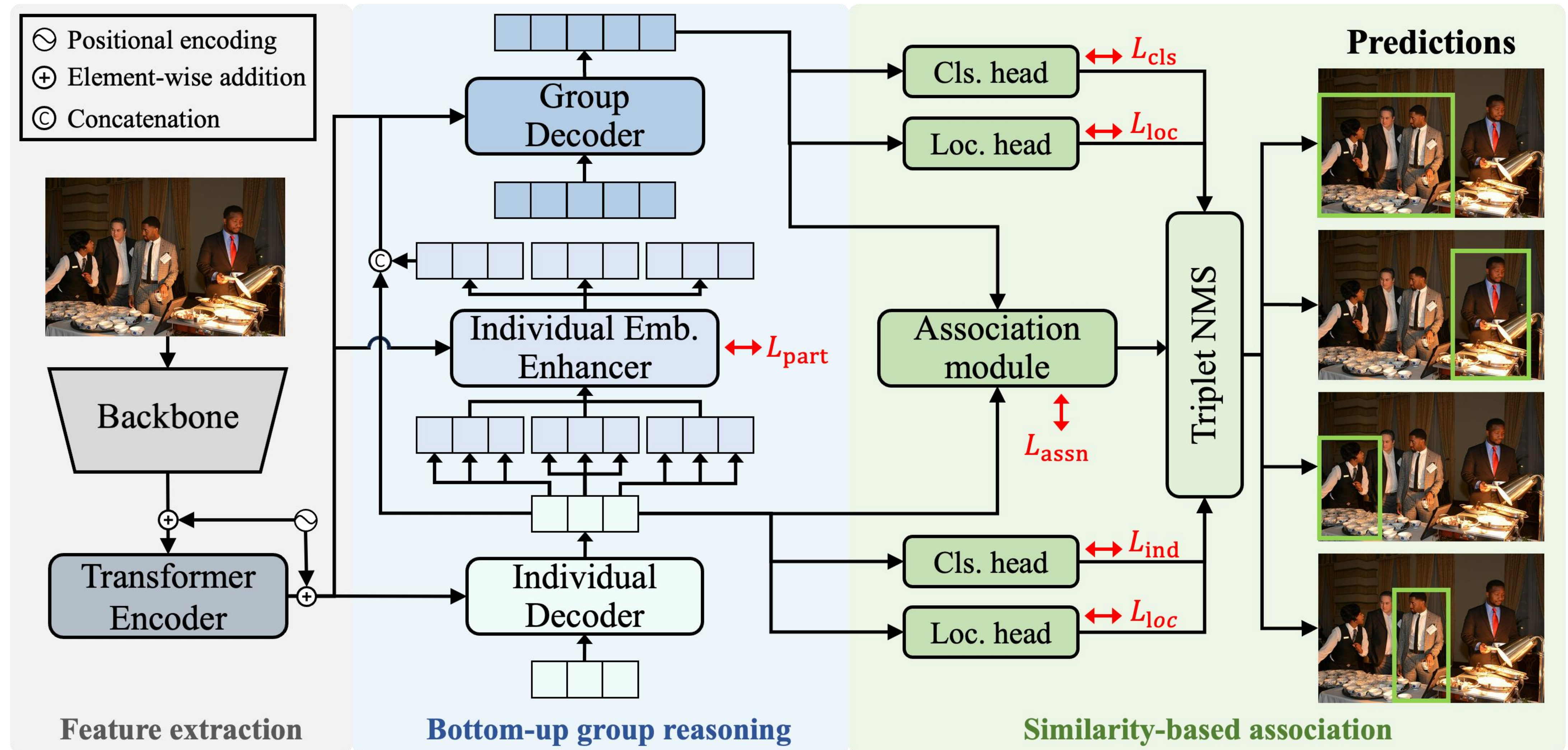
# Motivation

Leveraging pose as privileged information to capture part cues



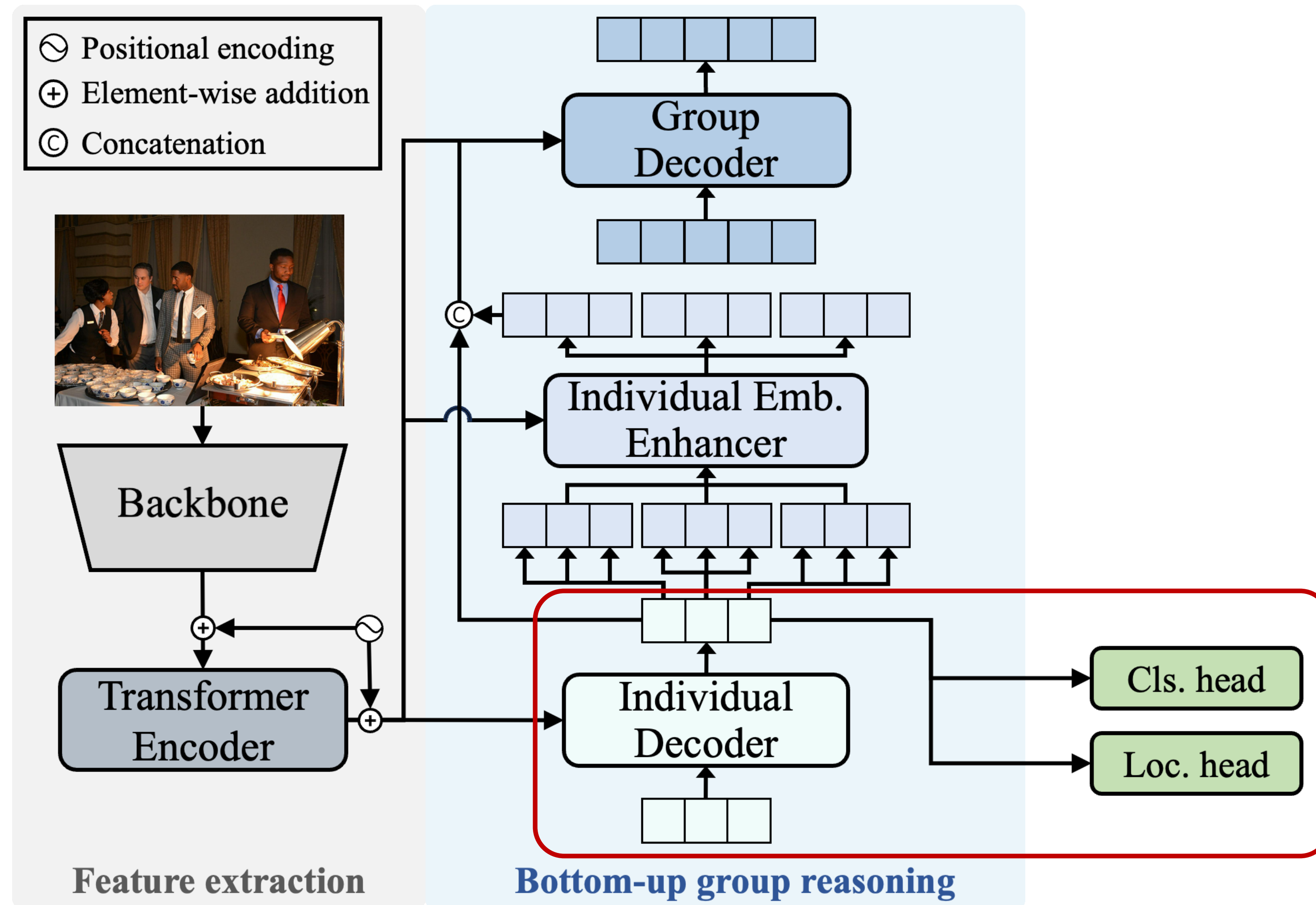


# Part-Aware Bottom-Up Group Reasoning

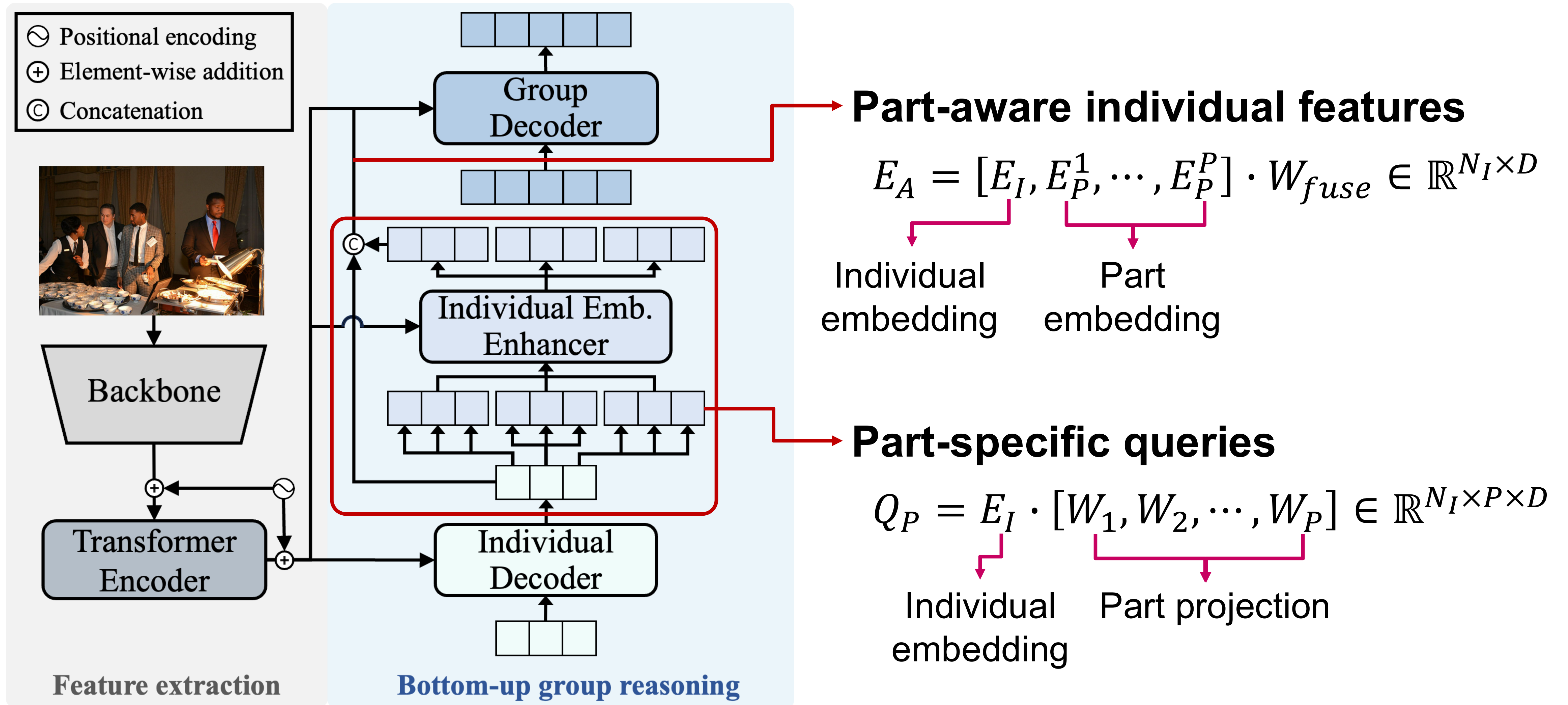




# Part-Aware Bottom-Up Group Reasoning

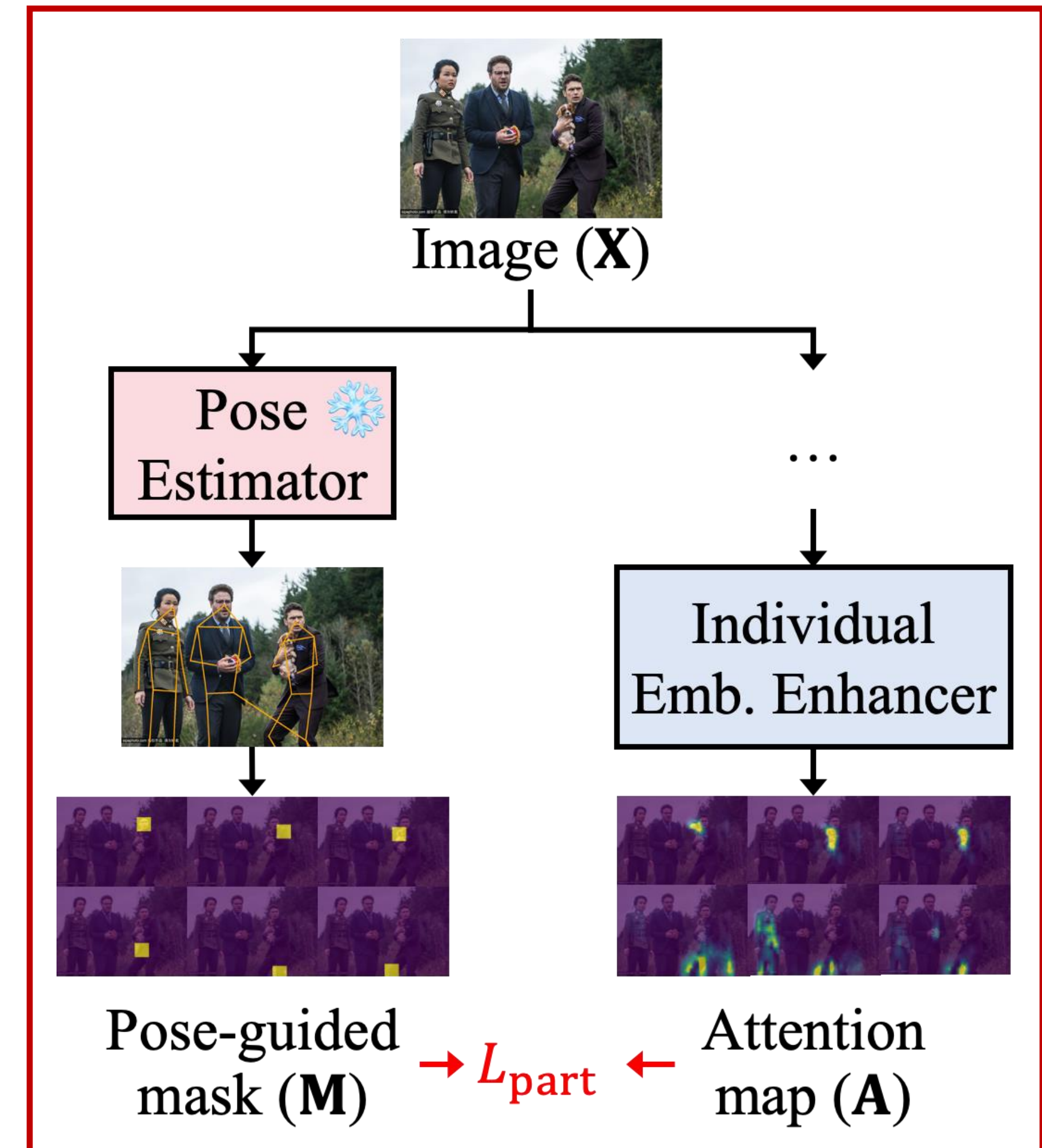
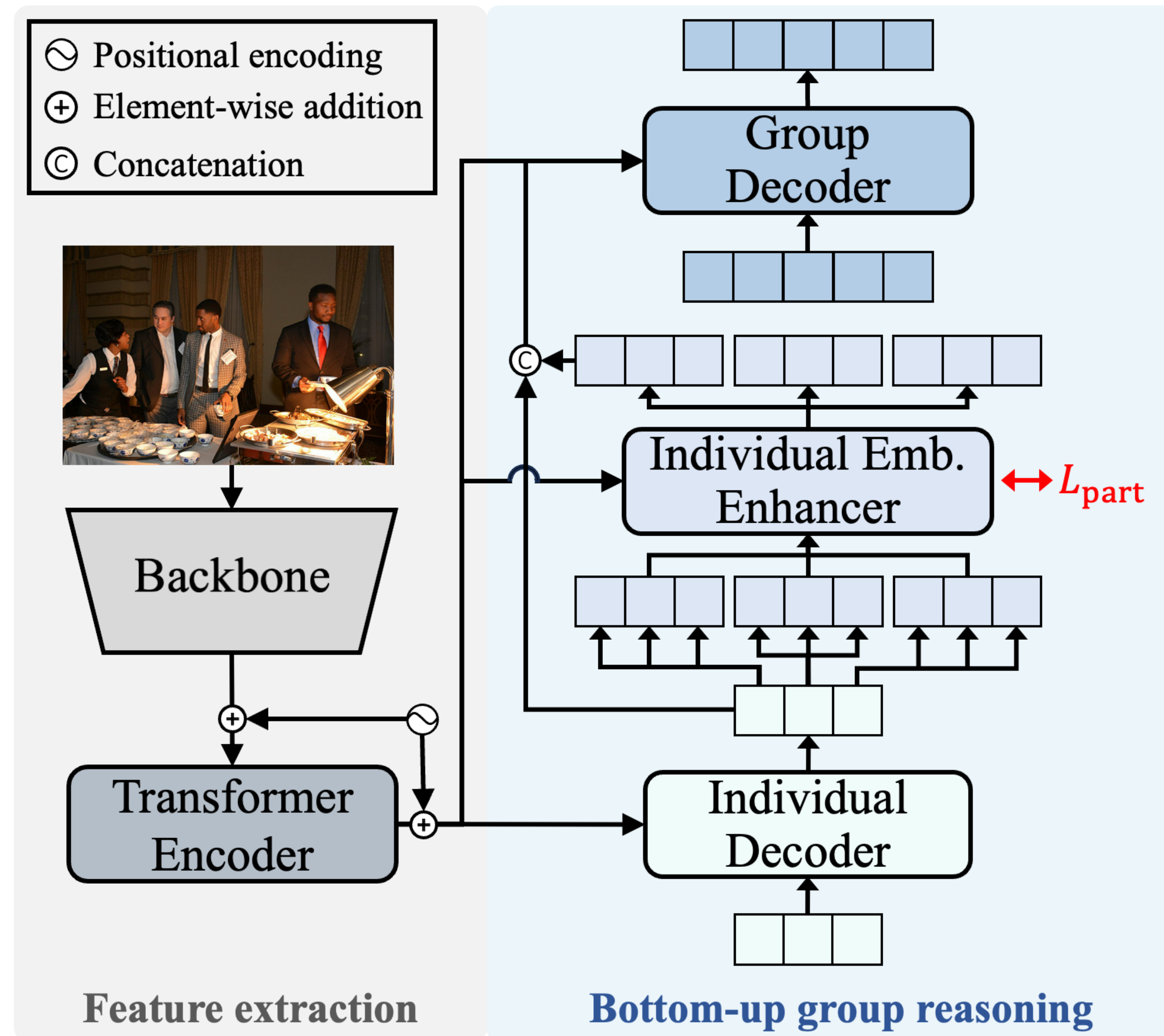


# Part-Aware Bottom-Up Group Reasoning



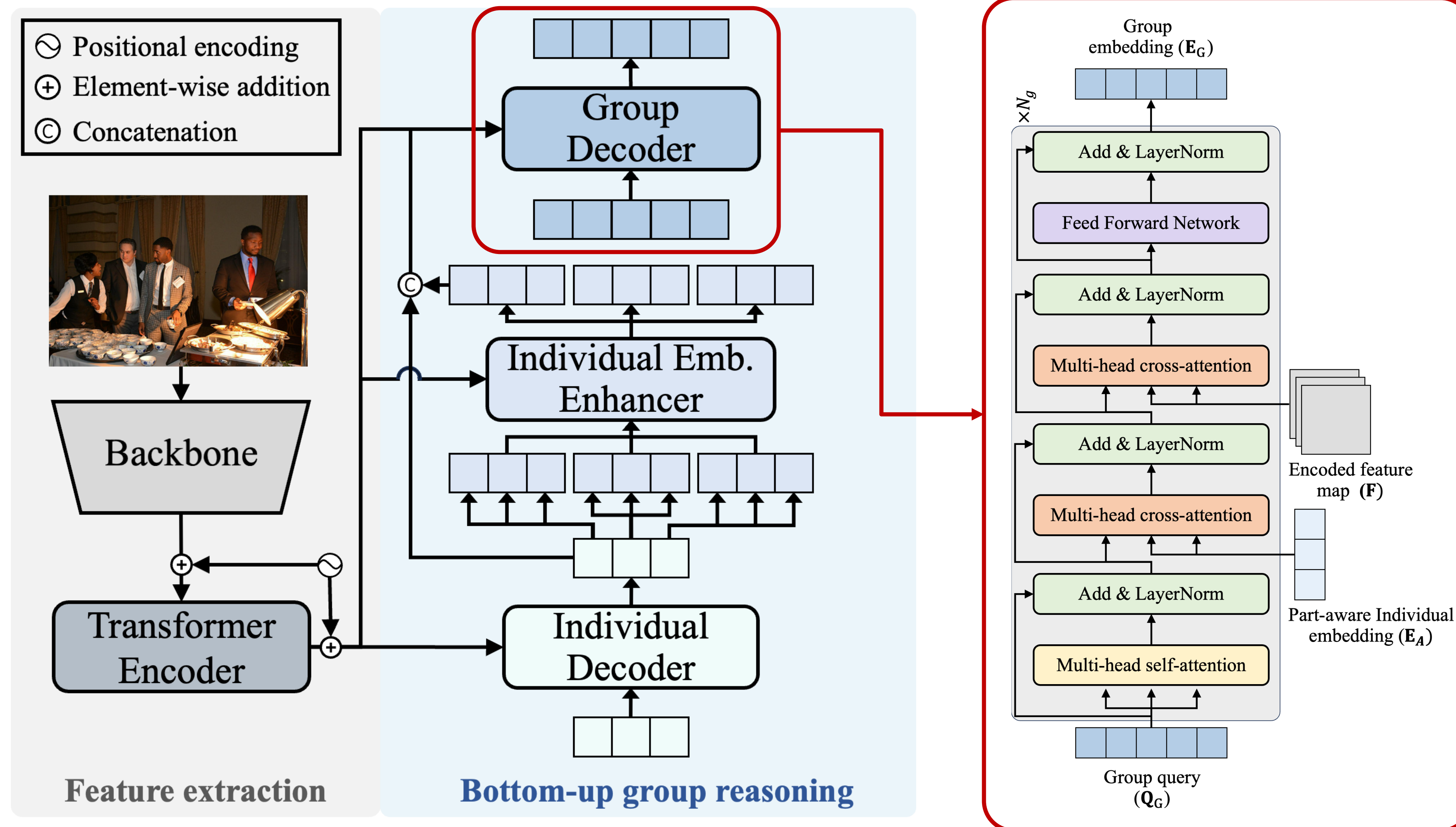


# Part-Aware Bottom-Up Group Reasoning



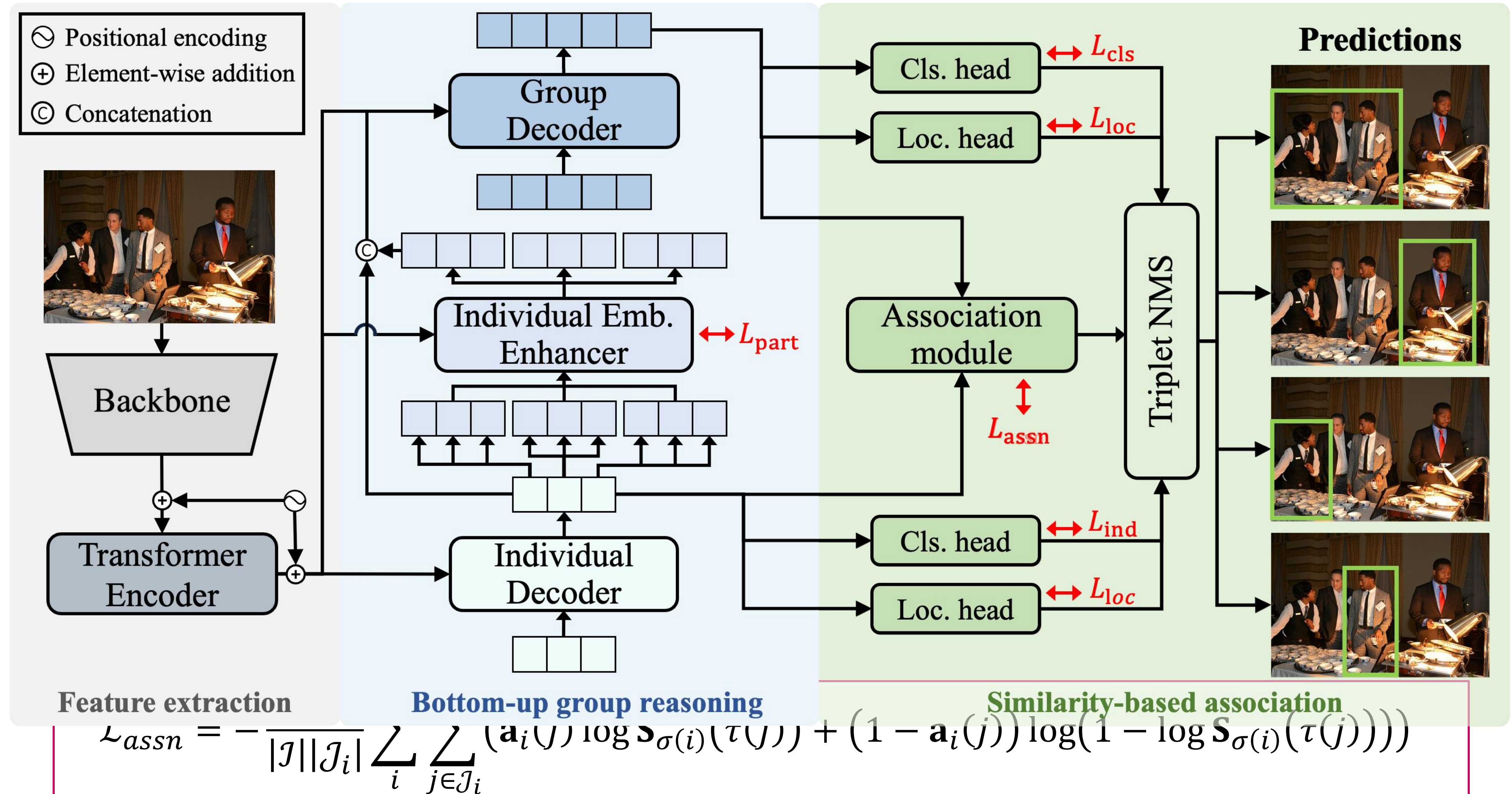


# Part-Aware Bottom-Up Group Reasoning





# Part-Aware Bottom-Up Group Reasoning





# Quantitative Result

| Method             | val          |              |              |              | test         |              |              |              |
|--------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|                    | mR@25        | mR@50        | mR@100       | AR           | mR@25        | mR@50        | mR@100       | AR           |
| <i>m</i> -QPIC     | 56.89        | 69.52        | 78.36        | 68.26        | 59.44        | 71.46        | 80.07        | 70.32        |
| <i>m</i> -CDN      | 55.57        | 71.06        | 78.81        | 68.48        | 59.01        | 72.94        | 82.61        | 71.52        |
| <i>m</i> -GEN-VLKT | 50.59        | 70.87        | 80.08        | 67.18        | 56.68        | 74.32        | 84.18        | 71.72        |
| NVI-DEHR           | 54.85        | 73.42        | 85.33        | 71.20        | 59.46        | 76.01        | 88.52        | 74.67        |
| Ours               | <b>59.43</b> | <b>76.62</b> | <b>87.43</b> | <b>74.49</b> | <b>63.59</b> | <b>80.62</b> | <b>91.34</b> | <b>78.52</b> |

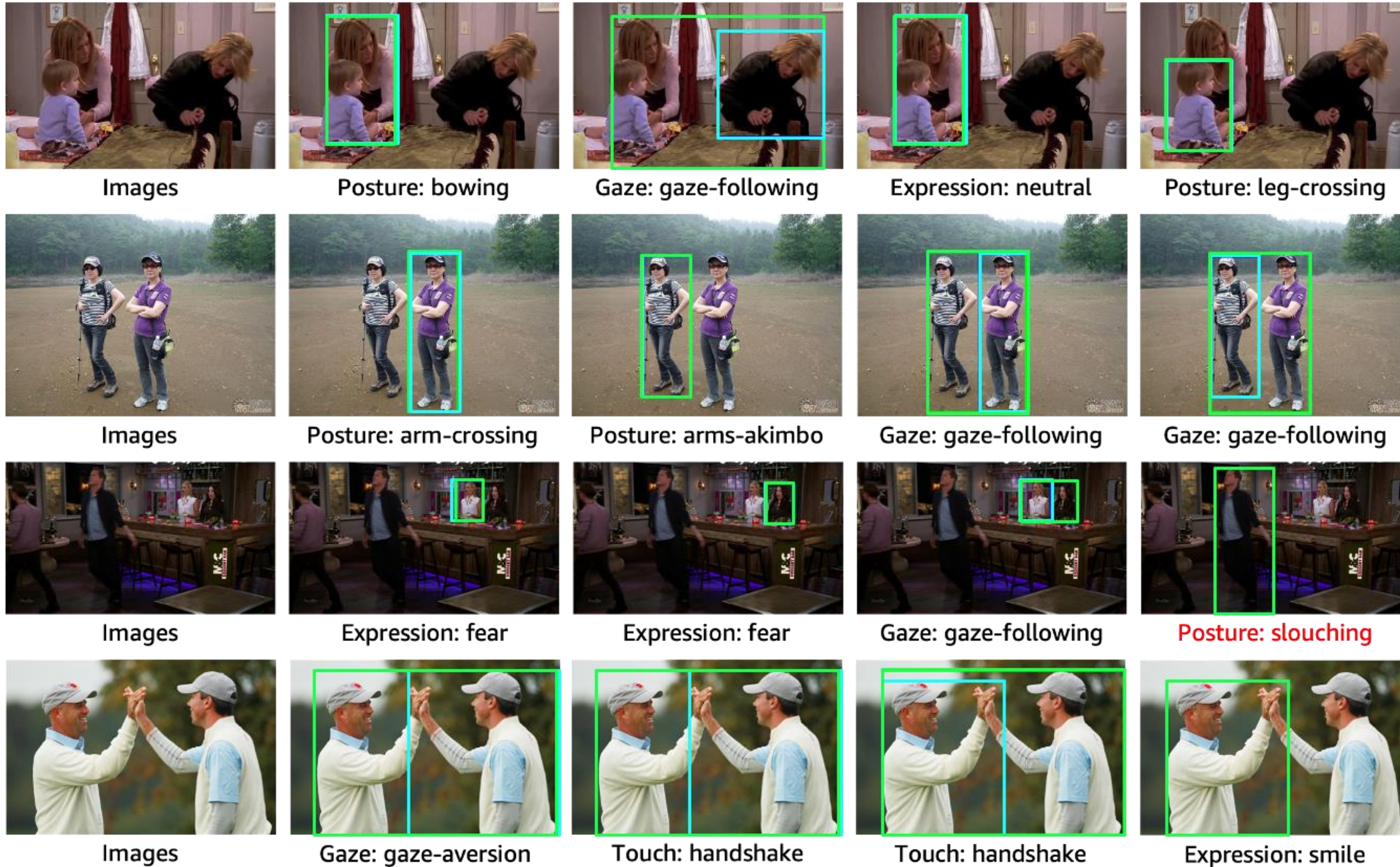
## Comparison on NVI dataset

| Method    | Split by view    |                  |                 | Split by place   |                  |                 |
|-----------|------------------|------------------|-----------------|------------------|------------------|-----------------|
|           | Group<br>mAP 1.0 | Group<br>mAP 0.5 | Outlier<br>mIoU | Group<br>mAP 1.0 | Group<br>mAP 0.5 | Outlier<br>mIoU |
| Joint     | 9.14             | 31.83            | 42.93           | 6.08             | 18.43            | 2.83            |
| JRDB-base | 12.63            | 35.53            | 31.85           | 8.15             | 22.68            | 33.03           |
| HGC       | 6.77             | 31.08            | 57.65           | 4.27             | 24.97            | 57.70           |
| Café-base | 14.36            | 37.52            | 63.70           | 8.29             | 28.72            | 59.60           |
| Ours      | <b>18.23</b>     | <b>46.88</b>     | <b>67.62</b>    | <b>10.65</b>     | <b>39.03</b>     | <b>63.60</b>    |

## Comparison on Café dataset



# Qualitative Result





# Conclusion

- Leverage human pose as privileged information to obtain part-aware representation.
- Propose a bottom-up group reasoning framework.
- Achieve superior performance on NVI and Café datasets.

Poster Session #1

December 3 (Wednesday)

11:00 ~ 14:00

In San Diego



Project Page



Paper