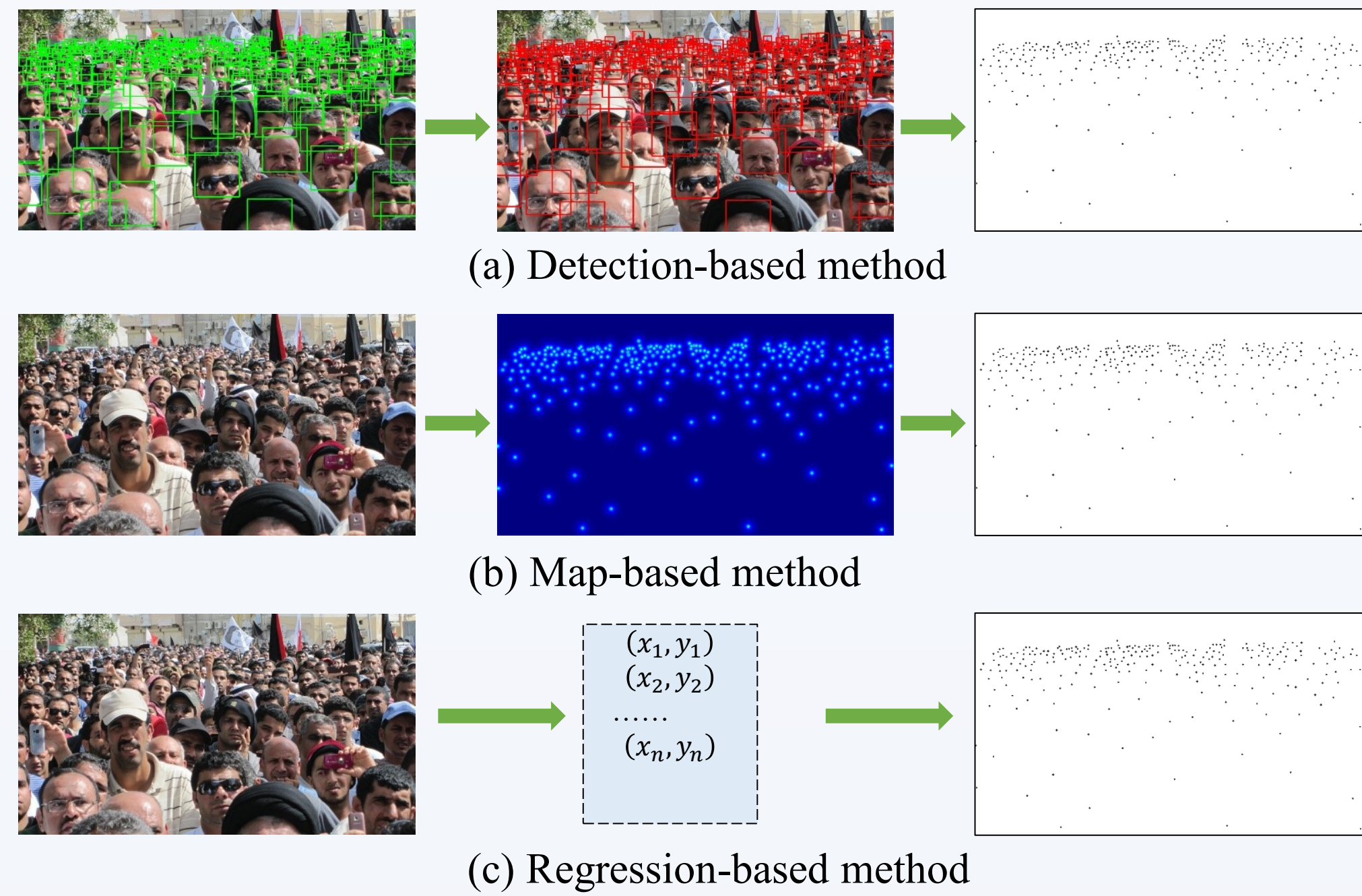


# An End-to-End Transformer Model for Crowd Localization

Dingkang Liang<sup>1</sup>, Wei Xu<sup>2</sup>, Xiang Bai<sup>1</sup>

<sup>1</sup> Huazhong University of Science and Technology; <sup>2</sup> Beijing University of Posts and Telecommunications

## MOTIVATION



- Crowd localization aims to provide the location of each instance.
- The regression-based methods, directly regressing the coordinates, are more straightforward than the detection-based and map-based methods.

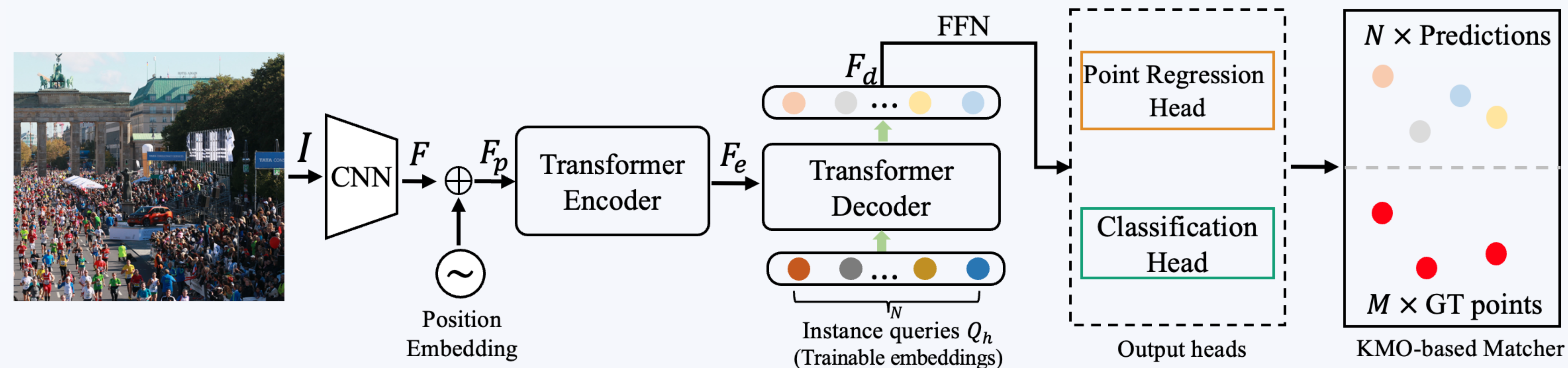


- DETR shows terrible performance in the crowd localization task, attributed to the intrinsic limitation of the matcher. Due to lack of context, the  $L1$  distance easily causes the ambiguous match pair.

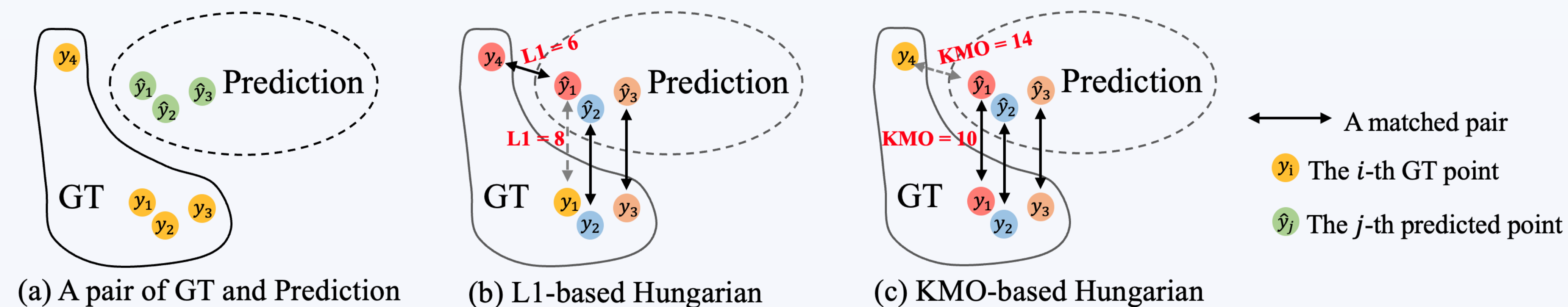
## CONTRIBUTION

- We propose an end-to-end Crowd Localization TRansformer framework named CLTR, which formulates the crowd localization as a point set prediction task.
- We introduce the KMO-based Hungarian bipartite matching, which takes the context from nearby heads as an auxiliary matching cost. As a result, the matcher can effectively reduce the ambiguous points and generate more reasonable matching results.

## METHOD



- The overview of our CLTR. First, the input image  $I$  is fed to the CNN-based backbone to extract the features  $F$ . Second, the features  $F$  are added position embedding, resulting in  $F_p$ , fed to the transformer-encoder layers, outputting  $F_e$ . Third, we define  $N \times$  trainable embeddings  $Q_h$  as query,  $F_e$  as key, and transformer decoder takes the  $Q_h$  and  $F_e$  as input to generate the decoded feature  $F_d$ . Finally, the  $F_d$  can be decoupled to the point coordinate and corresponding confidence score.



- (a) A pair of GT and predictions. (b) The  $L1$ -based Hungarian generate unsatisfactory matching results. (c) The proposed KMO-based Hungarian models the context as the matching cost, generating more reasonable matching results.

$$L_m(y_i, \hat{y}_j) = \|y_i^p - \hat{y}_j^p\|_1 - \hat{c}_j, i \in M, j \in N,$$

$$L_m^k(y_i, \hat{y}_j) = \|y_i^p - \hat{y}_j^p\|_1 + \|y_i^k - \hat{y}_j^k\|_1 - \hat{c}_j,$$

$$y_i^k = \frac{1}{k} \sum_{k=1}^k d_i^k, \quad \hat{y}_j^k = \frac{1}{k} \sum_{k=1}^k \hat{d}_j^k,$$

- Merely taking the  $L1$  with confidence will generate unsatisfactory matching results on specific cases.

- The proposed KMO-based matcher, revisiting the label assignment from a context view, turns to find the whole optimum.

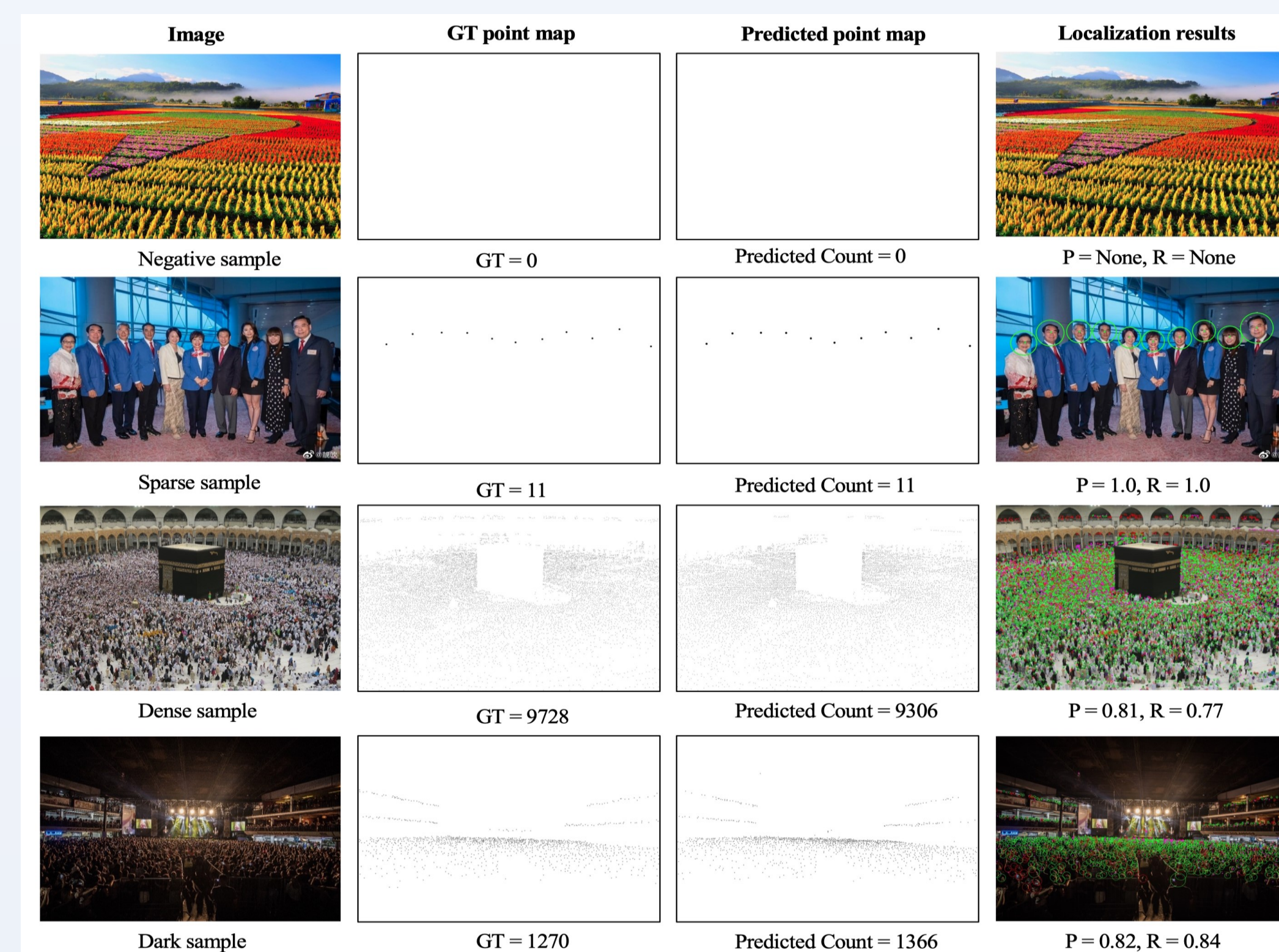
## RESULTS

### Localization performance on NWPU dataset.

Method	Validation set			Test set	
	P(%)	R(%)	F(%)	P(%)	F(%)
Faster RCNN* [29]	96.4%	3.8%	7.3%	95.8%	6.7%
TinyFaces* [11]	54.3%	66.6%	59.8%	52.9%	56.7%
TopoCount* [1]	-	-	-	69.5%	69.1%
GPR [7]	61.0%	52.2%	56.3%	55.8%	52.5%
RAZ_Loc [19]	69.2%	56.9%	62.5%	66.6%	59.8%
AutoScale_Loc [46]	70.1%	63.8%	66.8%	67.3%	62.0%
Crowd-SDNet [44]	-	-	-	65.1%	63.7%
GL [39]	-	-	-	80.0%	66.0%
CLTR (ours)	73.9%	71.3%	72.6%	69.4%	68.5%

### Counting performance on NWPU dataset.

## VISUALIZATIONS



Some examples from the NWPU dataset. From left to right, there are images, GT points, predicted points, and localization results.

## ACKNOWLEDGEMENT

This work was supported by National Key R&D Program of China (Grant No. 2018YFB1004602).