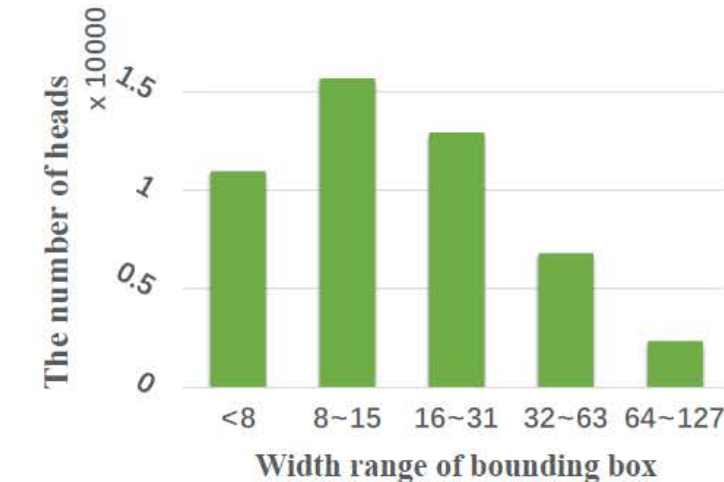




Introduction:

- **Crowd counting:** Estimate the number of persons in images or videos.
 - Recent methods: density map regression (cannot give the position of heads)
 - Ours: revisiting detection-based method.



(a) (b) (c)

Challenges:

- Underestimation: the number of detected heads is much smaller than the total number of heads [1] (especially for tiny heads);
- ground-truth annotation: heavy workload than point annotation.

Contributions:

- A regression guided detection network (RDNet) for RGB-D crowd counting and localization;
- Depth-adaptive kernel for density map generation, depth-anchor for anchor initialization, Depth-based bounding boxes ground-truth generation for training [figure (b)];
- A large-scale RGB-D crowd counting dataset (ShanghaiTechRGBD);
- Our method can be easily extended to RGB based counting and localization.

Our approach:

Density map regression module:

- Fixed-kernel density map:

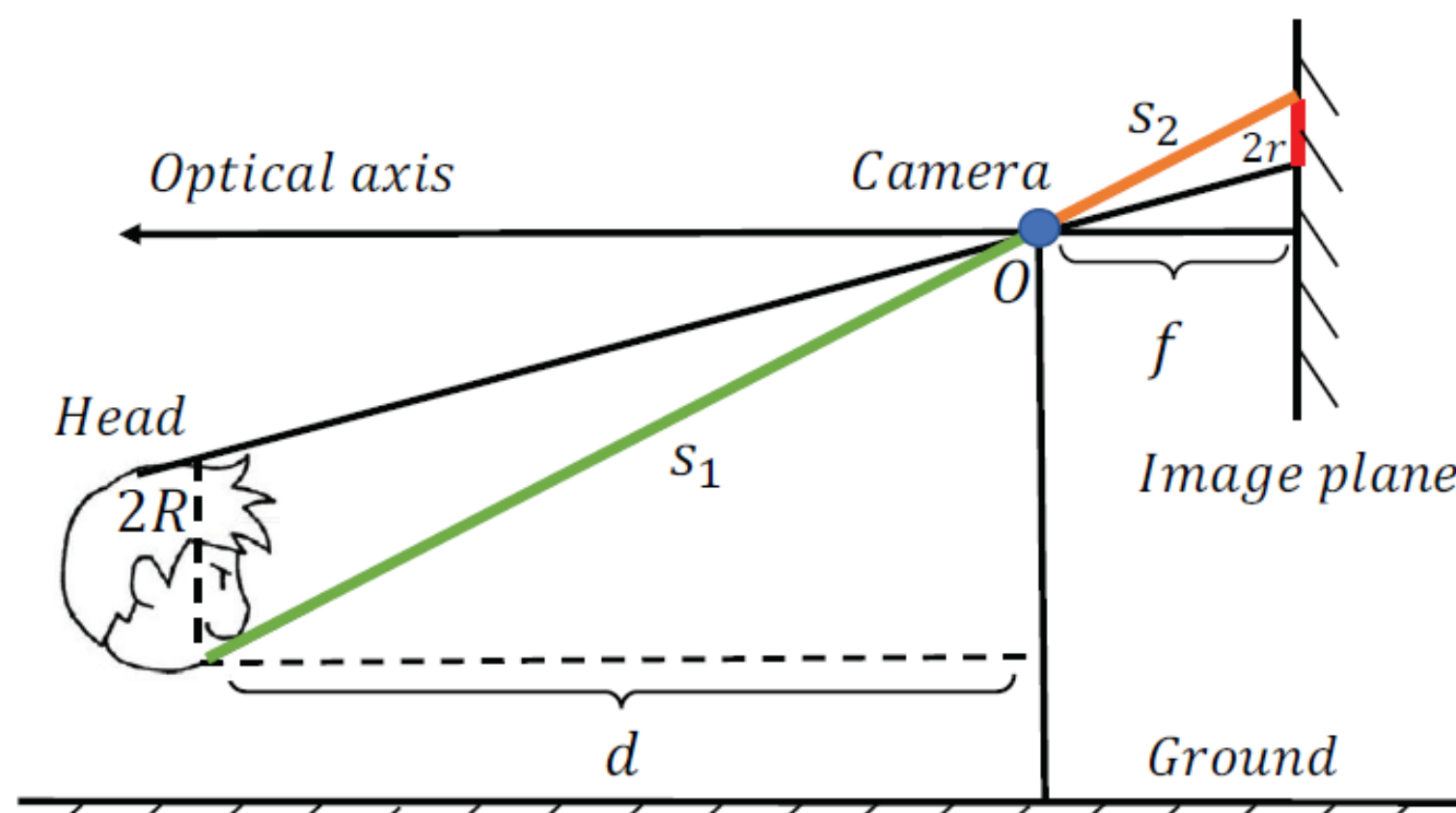
$$D(x) = \sum_{i=1}^N \delta(x - x_i) * G_{\sigma}(x).$$

- Depth-adaptive density map:

$$\frac{r}{R} = \frac{s_2}{s_1} = \frac{f}{d}$$

$$\sigma = \beta r = \beta \frac{Rf}{d} = \beta \frac{\gamma}{d}.$$

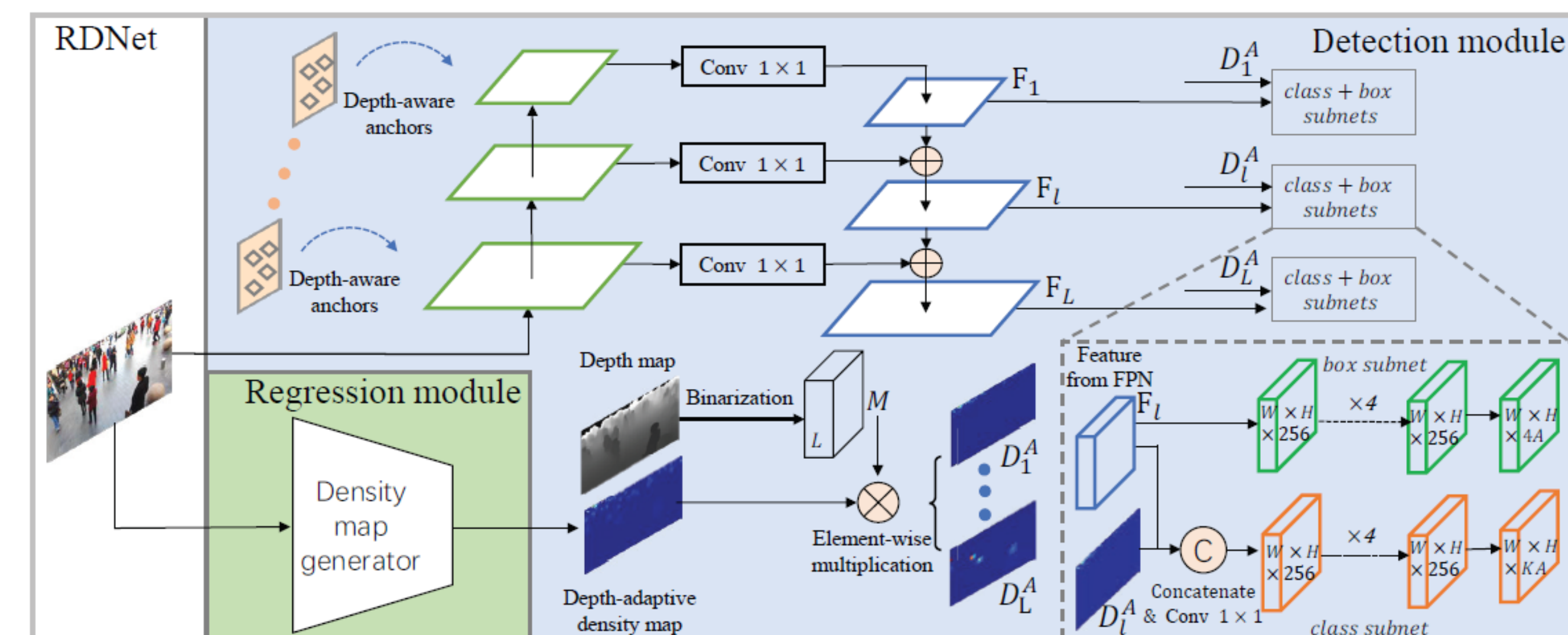
$$D^A(x) = \sum_{i=1}^N \delta(x - x_i) * G_{\sigma(d_i)}(x).$$



Detection module:

- Density map guided classification:
 - Decoding layer l ($l=1, \dots, L$), head size $[r_1, r_2]$, corresponding depth $d \in [\frac{\gamma}{r_2}, \frac{\gamma}{r_1}]$
 - Masked density map: $D_l^A = D^A \odot M_l$
- Depth-aware anchor: $H(m, n) = \frac{\gamma}{d(m, n)}$ for anchor initialization.
- Generation of bounding box for training:
 - Bounding boxes: $\mathcal{B} = \{b_1, \dots, b_N\}$ for N heads. Head width: $w_i = \frac{\gamma}{d_i}$,

Network architecture:



Dataset:

Our ShanghaiTechRGBD dataset:

- 2,193 images with 144,512 annotated head counts;
- 1,193 images for training.

Table 1. Comparisons of ShanghaiTechRGBD with some existing datasets: Num is the number of images; Max is the maximal crowd count within one image; Min is the minimal crowd count; Ave is the average crowd count; Total is total number of labeled heads.

Dataset	Resolution	Num	Max	Min	Ave	Total	Modality
CBSR [36]	Dataset 1	240 × 320	2834	7	0	1.6	4,541
	Dataset 2	240 × 320	1500	7	0	1	1,553
MICC [11]	480 × 640	3358	11	0	5.32	17,630	RGB + depth
ShanghaiTechRGBD	1080 × 1920	2193	234	6	65.9	144,512	RGB + depth



Experiments:

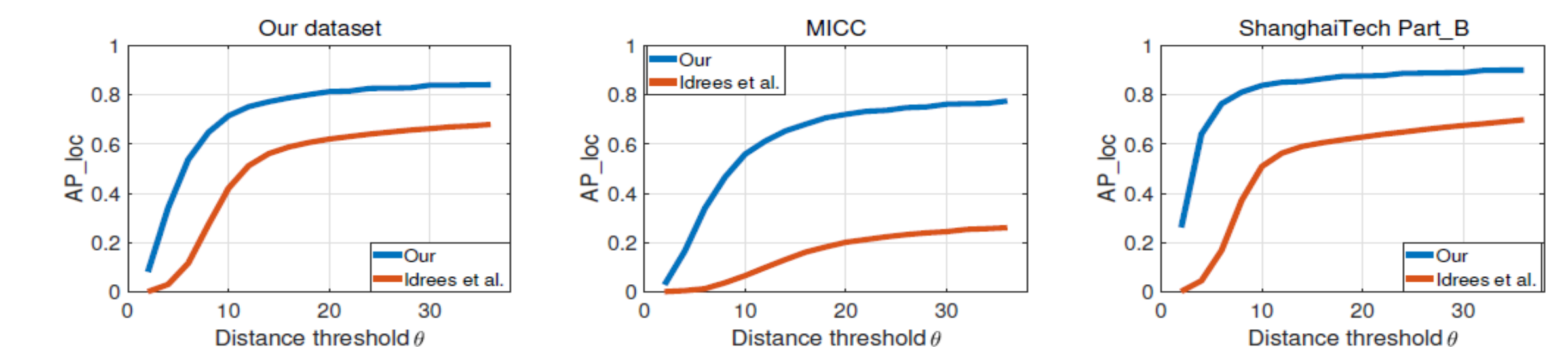
Counting and localization:

$$\text{MAE: } \text{MAE} = \frac{1}{M} \sum_{j=1}^M |N_j - \tilde{N}_j|,$$

$$\text{MSE: } \text{MSE} = \sqrt{\frac{1}{M} \sum_{j=1}^M (N_j - \tilde{N}_j)^2}$$

- AP_det for head detection evaluation [2];
- AP_loc for localization evaluation [3].

$$\text{Prediction} = \begin{cases} \text{Positive,} & \text{if dist} \leq \theta \\ \text{Negative,} & \text{otherwise} \end{cases}$$



Visualization:



Reference:

- [1] Jiang Liu, Chenqiang Gao, Deyu Meng, and Alexander G Hauptmann. Decidenet: Counting varying density crowds through attention guided detection and density estimation, in CVPR 2018.
- [2] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollar. Focal Loss for Dense Object Detection. In ICCV, 2017.
- [3] Haroon Idrees, Muhammad Tayyab, Kishan Athreya, Dong Zhang, Somaya Al-Maadeed, Nasir Rajpoot, and Mubarak Shah. Composition loss for counting, density map estimation and localization in dense crowds. In ECCV 2018.

Code & dataset

