# Bottom-up Object Detection by Grouping Extreme and Center Points - Supplementary Material

Xingyi Zhou
UT Austin
zhouxy@cs.utexas.edu

Jiacheng Zhuo
UT Austin
jzhuo@cs.utexas.edu

Philipp Krähenbühl
UT Austin
philkr@cs.utexas.edu

---

**Algorithm 1:** Center Grouping

---

**Input** : Center and Extremepoint heatmaps of an image for one category: $\hat{Y}^{(c)}, \hat{Y}^{(t)}, \hat{Y}^{(l)}, \hat{Y}^{(b)}, \hat{Y}^{(r)} \in (0,1)^{H \times W}$
Center and peak selection thresholds: $\tau_c$ and $\tau_p$

**Output**: Bounding box with score

// Convert heatmaps into coordinates of keypoints.

// $\mathcal{T}, \mathcal{L}, \mathcal{B}, \mathcal{R}$ are sets of points.

$\mathcal{T} \leftarrow$ ExtractPeak$(\hat{Y}^{(t)}, \tau_p)$
$\mathcal{L} \leftarrow$ ExtractPeak$(\hat{Y}^{(l)}, \tau_p)$
$\mathcal{B} \leftarrow$ ExtractPeak$(\hat{Y}^{(b)}, \tau_p)$
$\mathcal{R} \leftarrow$ ExtractPeak$(\hat{Y}^{(r)}, \tau_p)$
$B_x \leftarrow$ A list of emtpy lists of size $2W$
$B_y \leftarrow$ A list of emtpy lists of size $2H$
**for** $t \in \mathcal{T}, b \in \mathcal{B}$ **do**
  **if** $t_y \leq b_y$ **then**
  | Add $(t, b)$ to $B_y[t + b]$
  **end**
**end**
**for** $l \in \mathcal{L}, r \in \mathcal{R}$ **do**
  **if** $l_x \leq r_x$ **then**
  | Add $(l, r)$ to $B_x[l + r]$
  **end**
**end**
**for** $i \in [0, 2W - 1], j \in [0, 2H - 1]$ **do**
  $c_x = i/2, c_y = j/2$
  **if** $\hat{Y}^{(c)}_{c_x, c_y} \geq \tau_c$ **then**
    **for** $(l, r) \in B_x[i], (t, b) \in B_y[j]$ **do**
      **if** $t_y \leq l_y, r_y \leq b_y$ **and** $l_x \leq t_x, b_x \leq r_x$ **then**
        Add Bounding box $(l_x, t_y, r_x, b_y)$ with score
        $(\hat{Y}^{(t)}_{t_x, t_y} + \hat{Y}^{(l)}_{l_x, l_y} + \hat{Y}^{(b)}_{b_x, b_y} + \hat{Y}^{(r)}_{r_x, r_y} +$
        $\hat{Y}^{(c)}_{c_x, c_y})/5$.
      **end**
    **end**
  **end**
**end**

---

## 1. $O(n^2)$ grouping algorithm

Algorithm. 1 gives an $O(n^2)$ center grouping algorithm. The basic idea is to separate the enumeration between $x$ dimension and $y$ dimension and early reject invalid bounding boxes. However, the algorithm is hard to be accelerated by GPU. On the other hand, $n$ is less than 40 in our experiments and the $O(n^4)$ grouping algorithm runs effectively on GPU.

| | $AP$ | $AR_1$ | $AR_{10}$ | $AR_{100}$ | $AR_S$ | $AR_M$ | $AR_L$ |
|---|---|---|---|---|---|---|---|
| CornerNet (SS) | 40.5 | 35.3 | 54.3 | 59.1 | 37.4 | 61.9 | 76.9 |
| CornerNet (MS) | 42.1 | 36.4 | 55.7 | 60.0 | 38.5 | 62.7 | 77.4 |
| ExtremeNet (SS) | 39.8 | 31.6 | 49.9 | 53.0 | 30.0 | 56.5 | 69.5 |
| ExtremeNet (MS) | 43.2 | 34.7 | 55.6 | 60.0 | 37.0 | 63.3 | 78.1 |

Table 1: Average recall evaluation on COCO test-dev comparing to CornerNet [1] . SS/ MS are short for single-scale/ multi-scale tesing, respectively. $AR_k$ is the bounding box recall of the top $k$ predictions. $AR_S$, $AR_M$, $AR_L$ are the recall of the top 100 predictions for small, median, large objects, respectively.

We find the $O(n^2)$ algorithm runs slower in practice.

## 2. Average recall evaluation

Table. 1 shows the average recall on COCO test-dev set. We observe that our single-scale model yields considerably lower recall than single scale CornerNet, although their overall AP are close. While in the multi-scale setting, the AR of ExtremeNet and CornerNet are comparable. This shows multi-scale testing benefits ExtremeNet by increasing its recall. It can be understood that our center grouping requires strict pixel-level accuracy for the matching between the computed geometric center and the center heatmap. Multi-scale testing implicitly gives an error tolerance and gives more detections at a good precision.
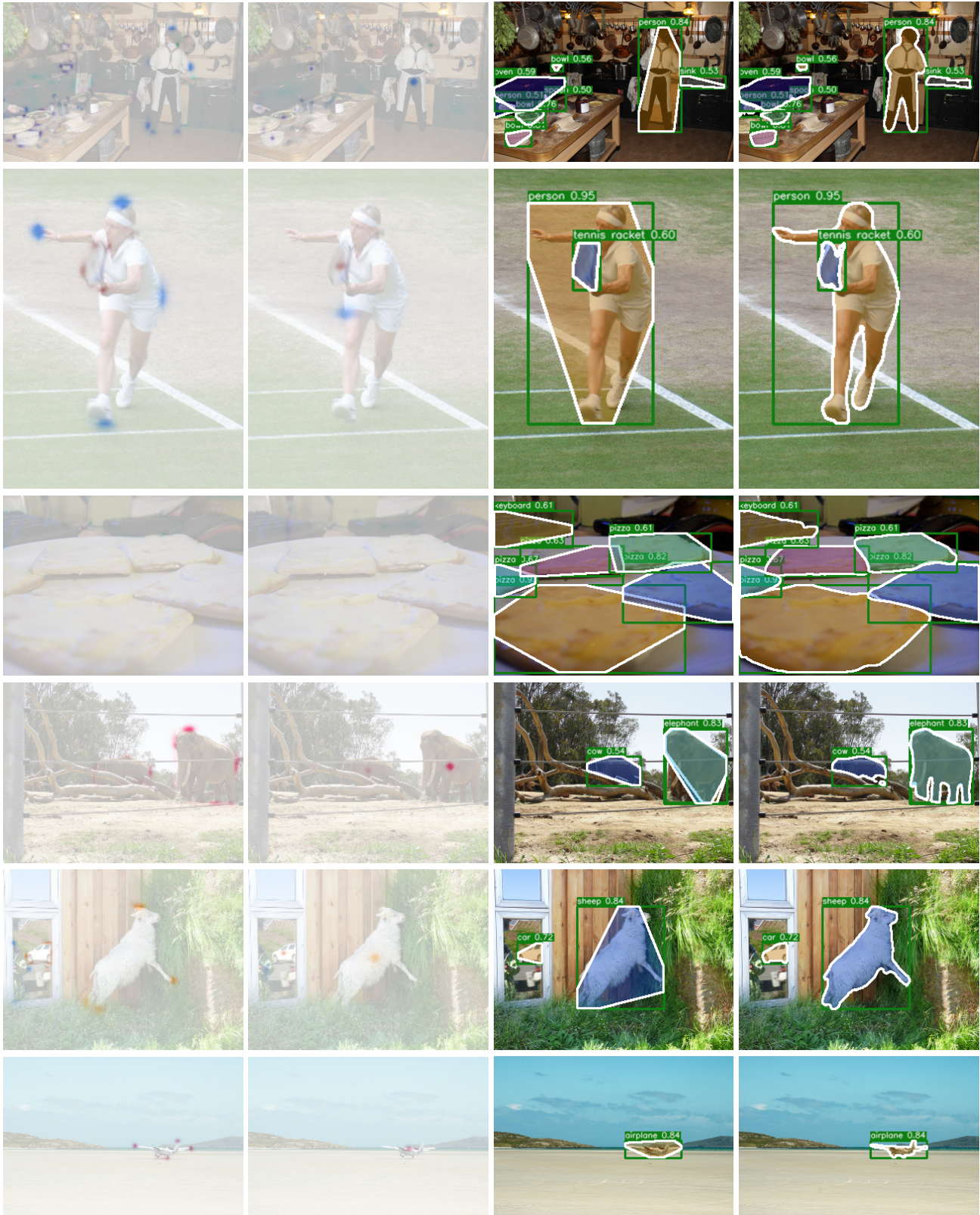
## 3. More qualitative results

More qualitative results. We *uniformly* sample 50 images from the 5000 images of COCO val2017. We show detection results with confidence $\geq 0.5$. First and second column: our predicted extreme point heatmap and center heatmap. Third column: our predicted bounding box and the octagon mask formed by extreme points. Fourth column: resulting masks of feeding our extreme point predictions to DeepExtremeCut [2].

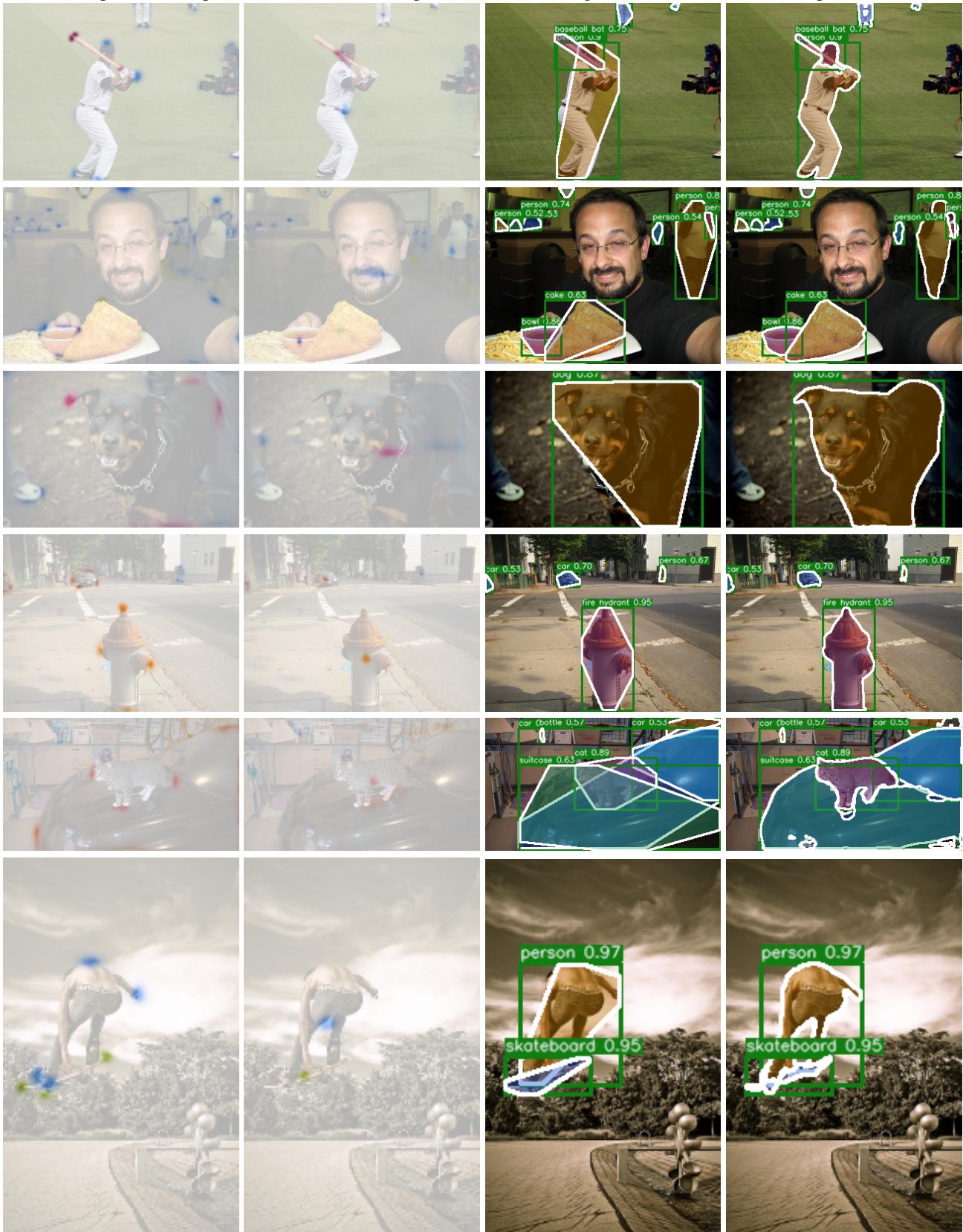| Extreme point heatmap | Center heatmap | Octagon mask | Extreme points+DEXTR |
| --- | --- | --- | --- |

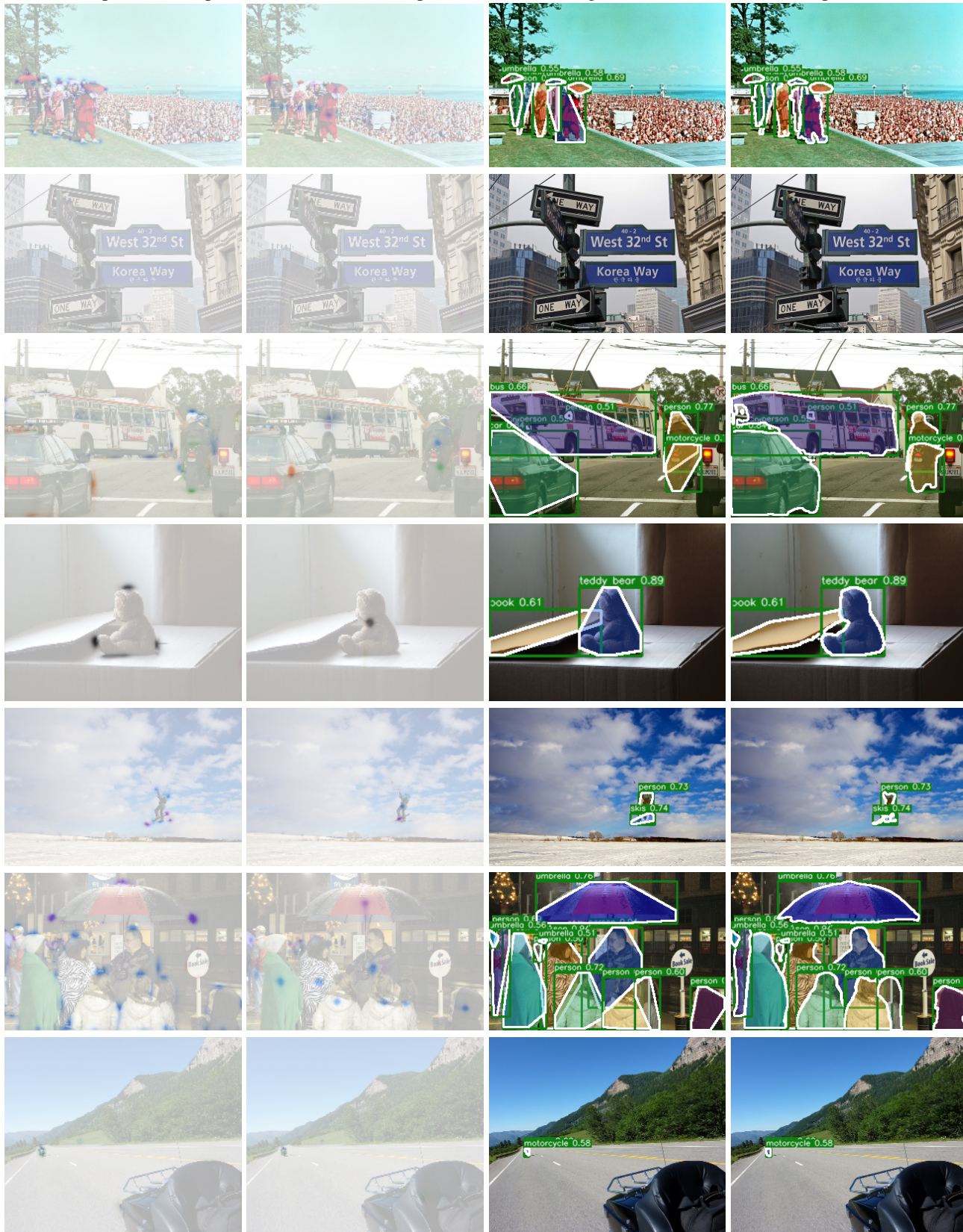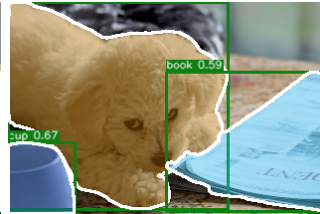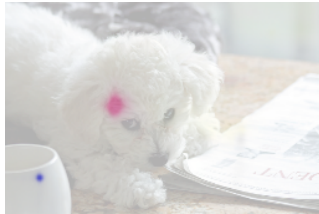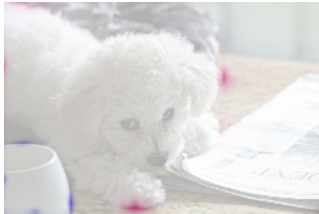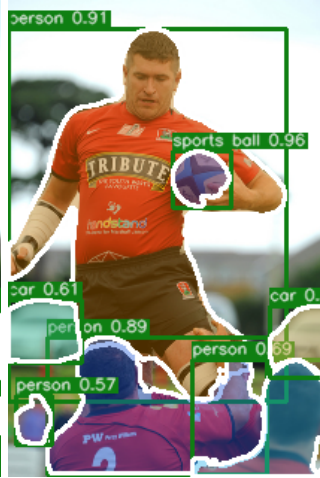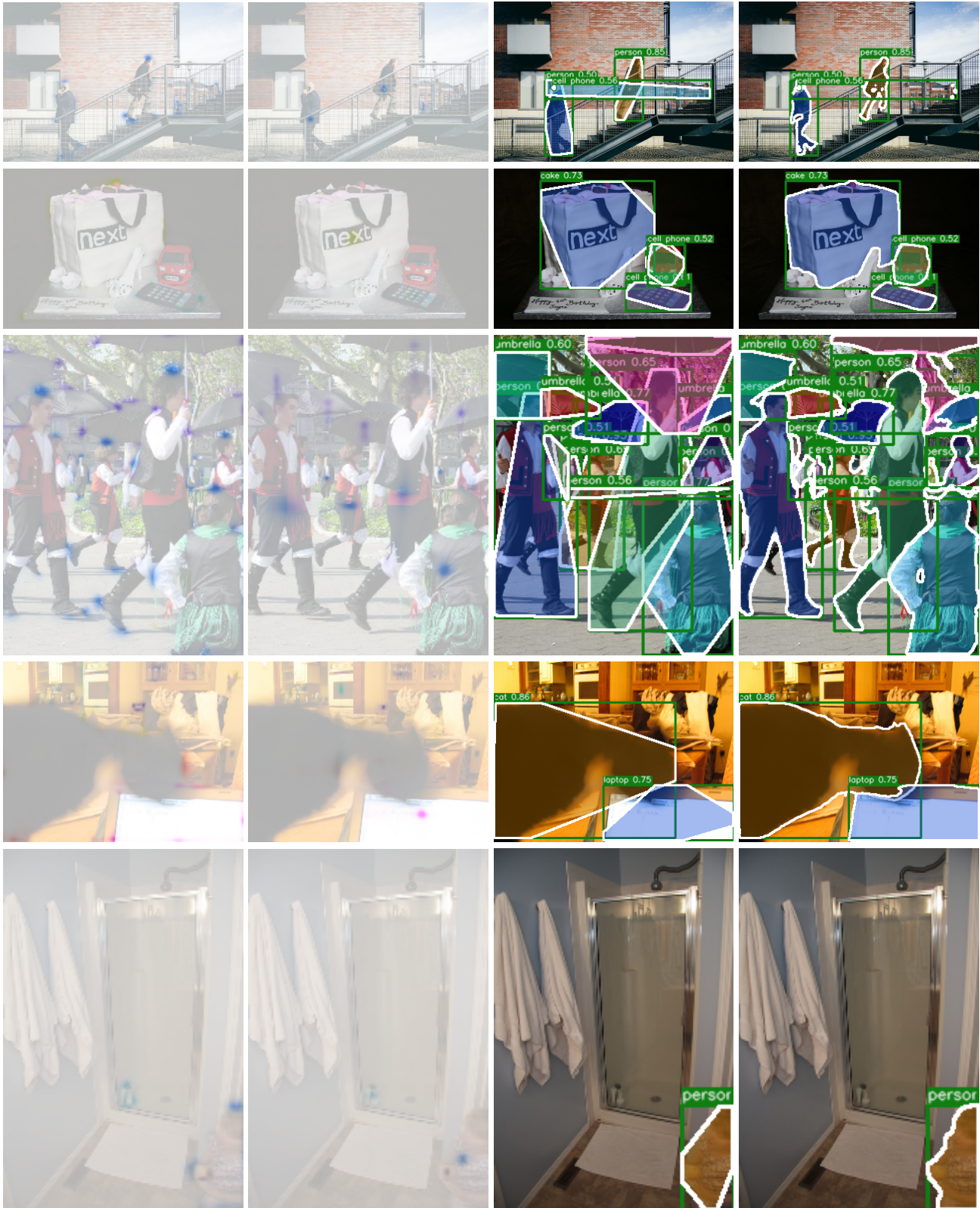| Extreme point heatmap | Center heatmap | Octagon mask | Extreme points+DEXTR |
|---|---|---|---|

| Extreme point heatmap | Center heatmap | Octagon mask | Extreme points+DEXTR |
|---|---|---|---|

| Extreme point heatmap | Center heatmap | Octagon mask | Extreme points+DEXTR |
|---|---|---|---|

| Extreme point heatmap | Center heatmap | Octagon mask | Extreme points+DEXTR |
| --- | --- | --- | --- |

| Extreme point heatmap | Center heatmap | Octagon mask | Extreme points+DEXTR |
| --- | --- | --- | --- |

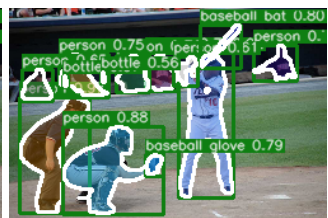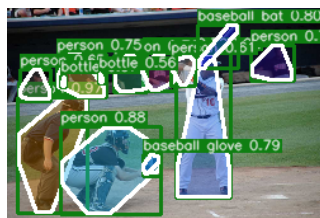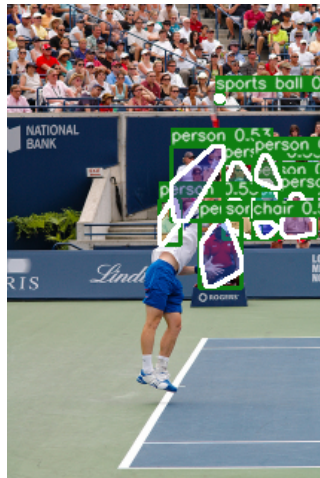| Extreme point heatmap | Center heatmap | Octagon mask | Extreme points+DEXTR |

| Extreme point heatmap | Center heatmap | Octagon mask | Extreme points+DEXTR |

| Extreme point heatmap | Center heatmap | Octagon mask | Extreme points+DEXTR |
|---|---|---|---|

## References

[1] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *ECCV*, 2018. 1

[2] K.K. Maninis, S. Caelles, J. Pont-Tuset, and L. Van Gool. Deep extreme cut: From extreme points to object segmentation. In *CVPR*, 2018. 1