

# RPTMask

Winner of VOT2021 short-term challenge

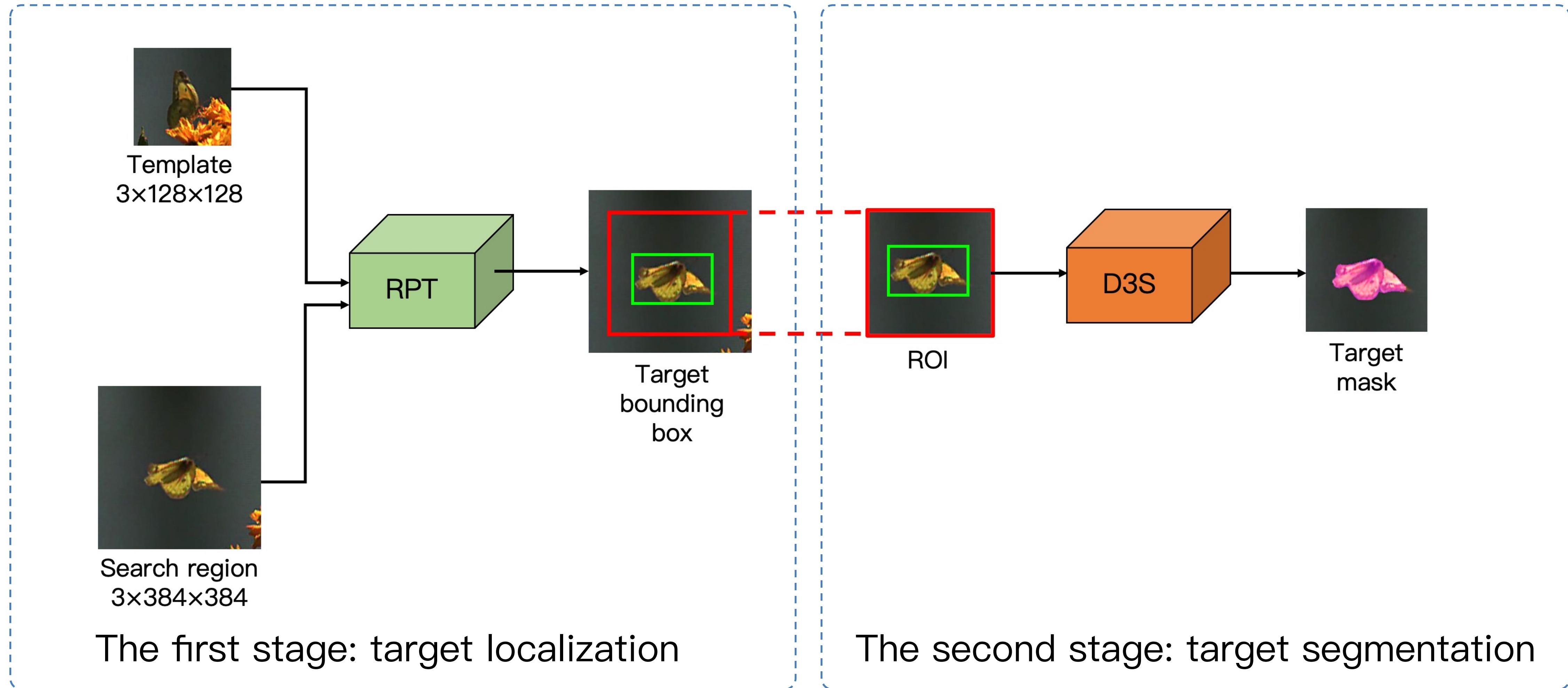
Zhihong Fu<sup>1</sup>, Liangliang Wang<sup>2</sup>, Qili Deng<sup>2</sup>, Kang Du<sup>2</sup>, Min Zheng<sup>2</sup>, Qingjie Liu<sup>1</sup>

<sup>1</sup> Beihang University, China

<sup>2</sup> ByteDance Inc. China



# Pipeline of VOT2020-ST winner



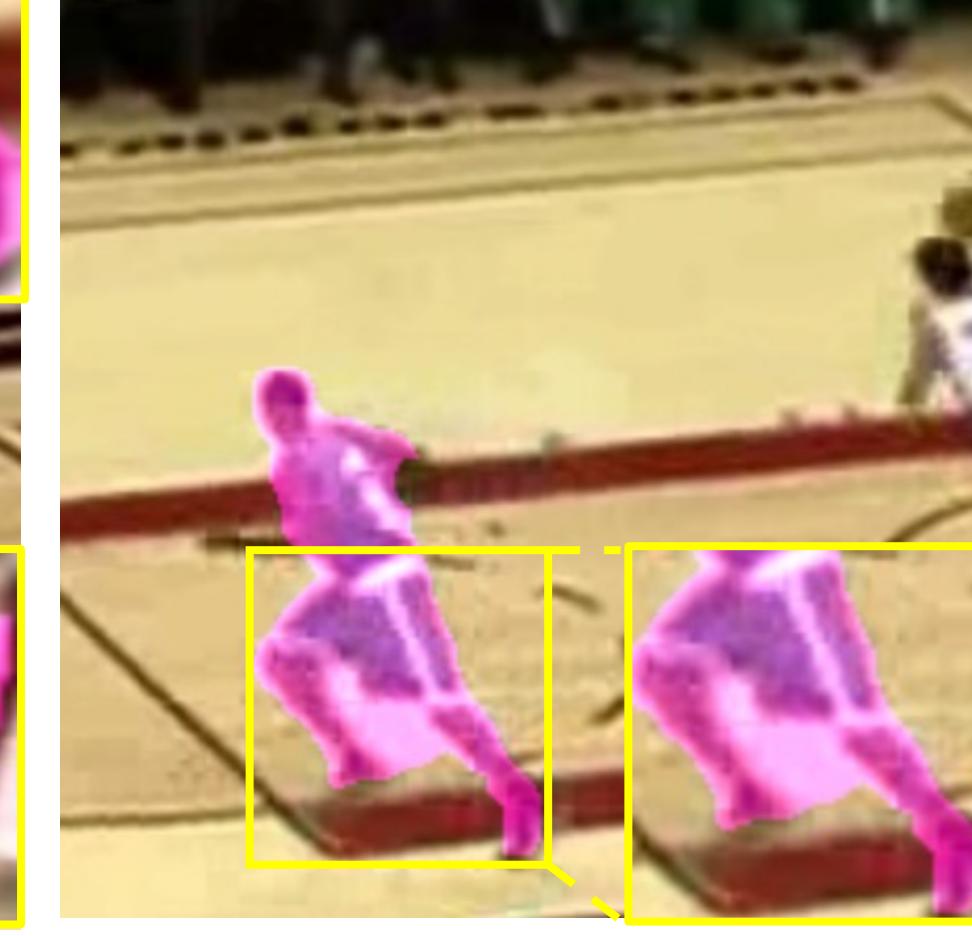
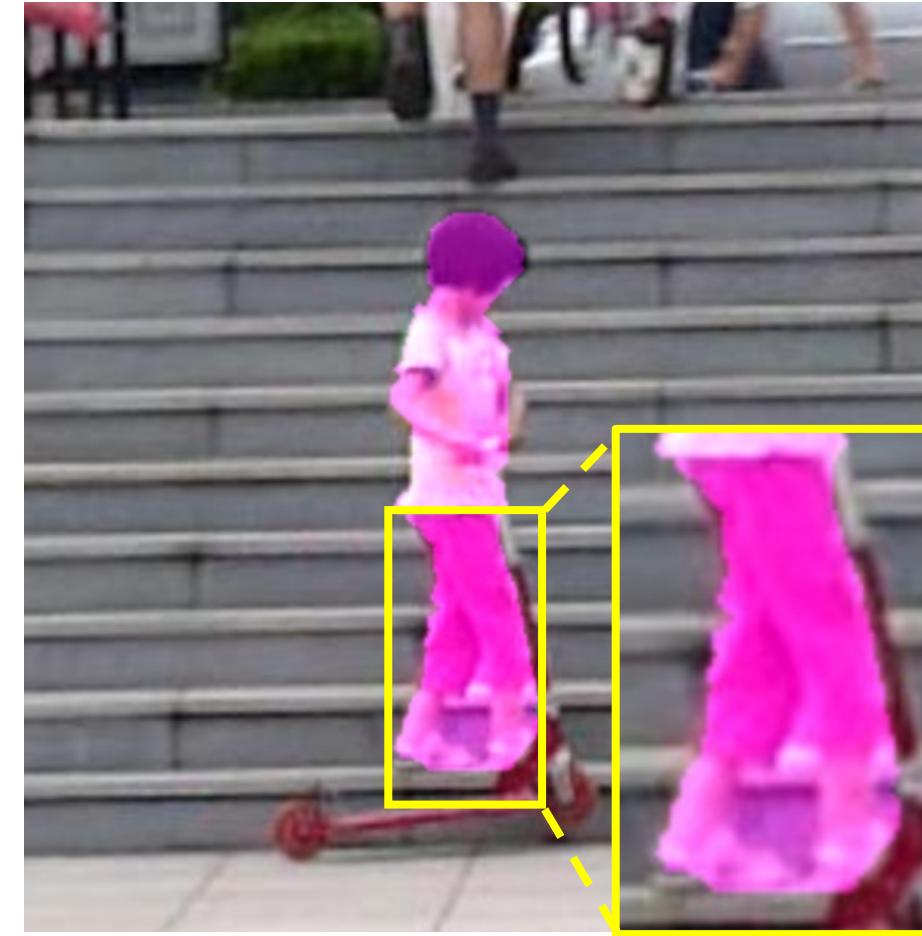
RPT: Learning Point Set Representation for Siamese Visual Tracking (ECCV 2020 Workshops)

# Motivation

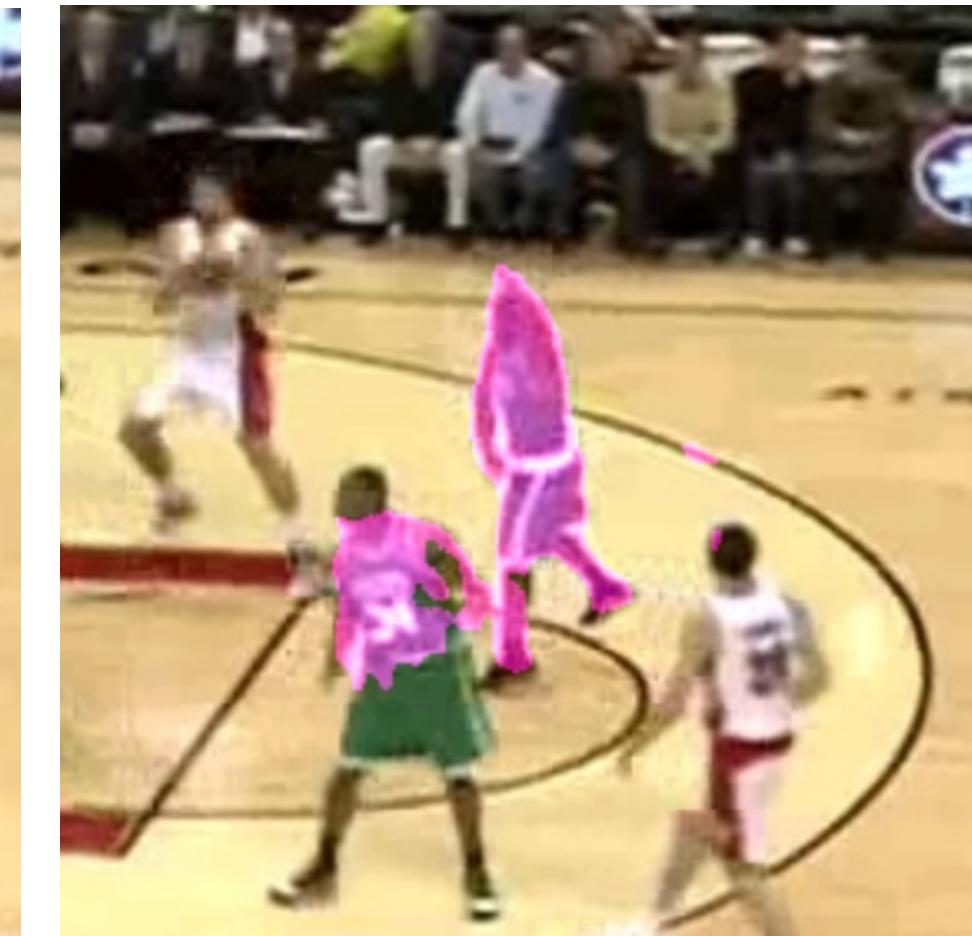
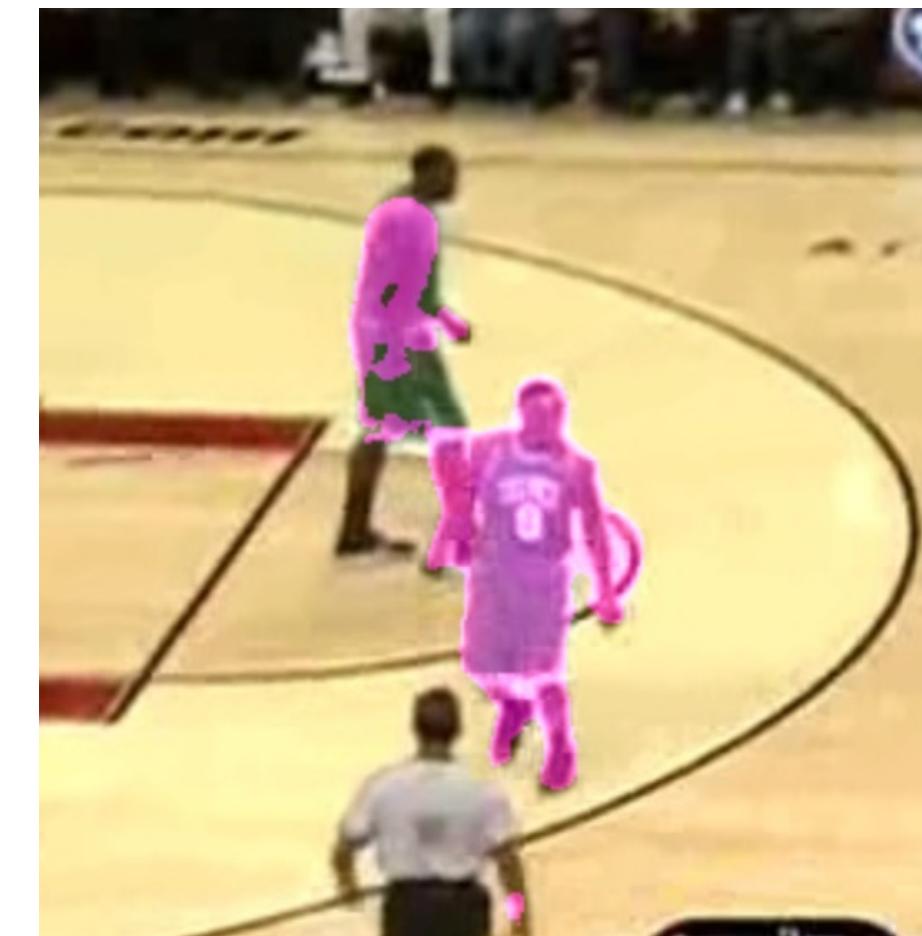
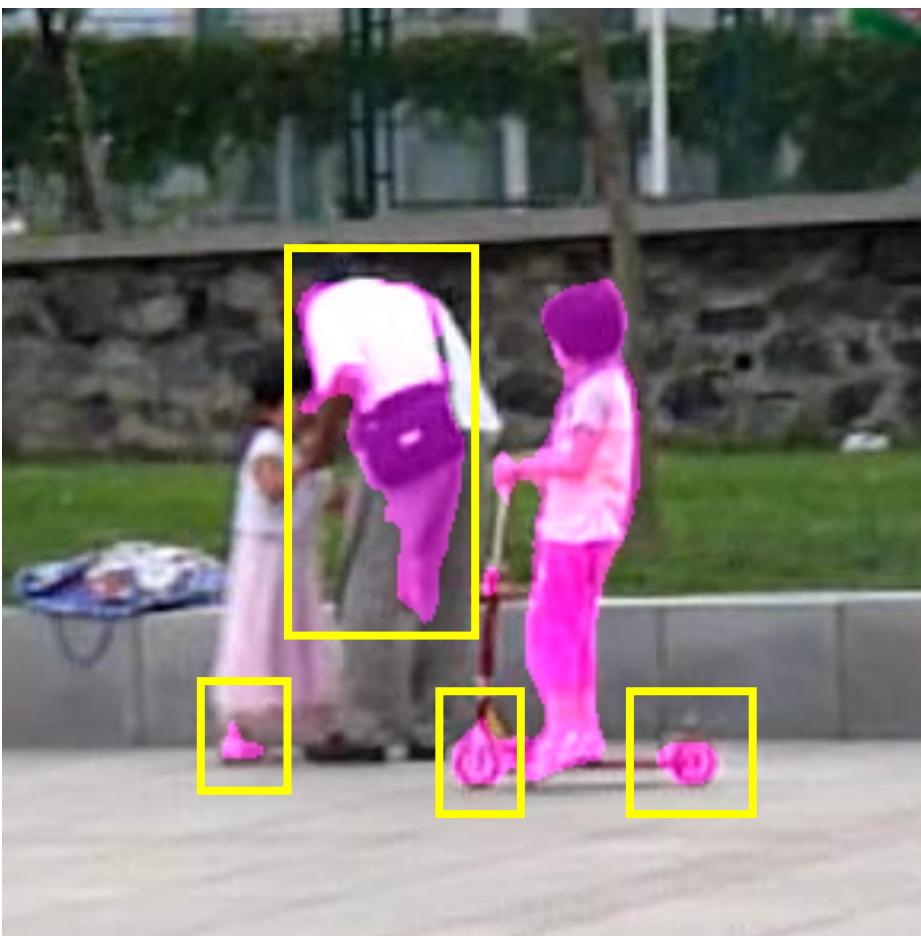


baseline's results

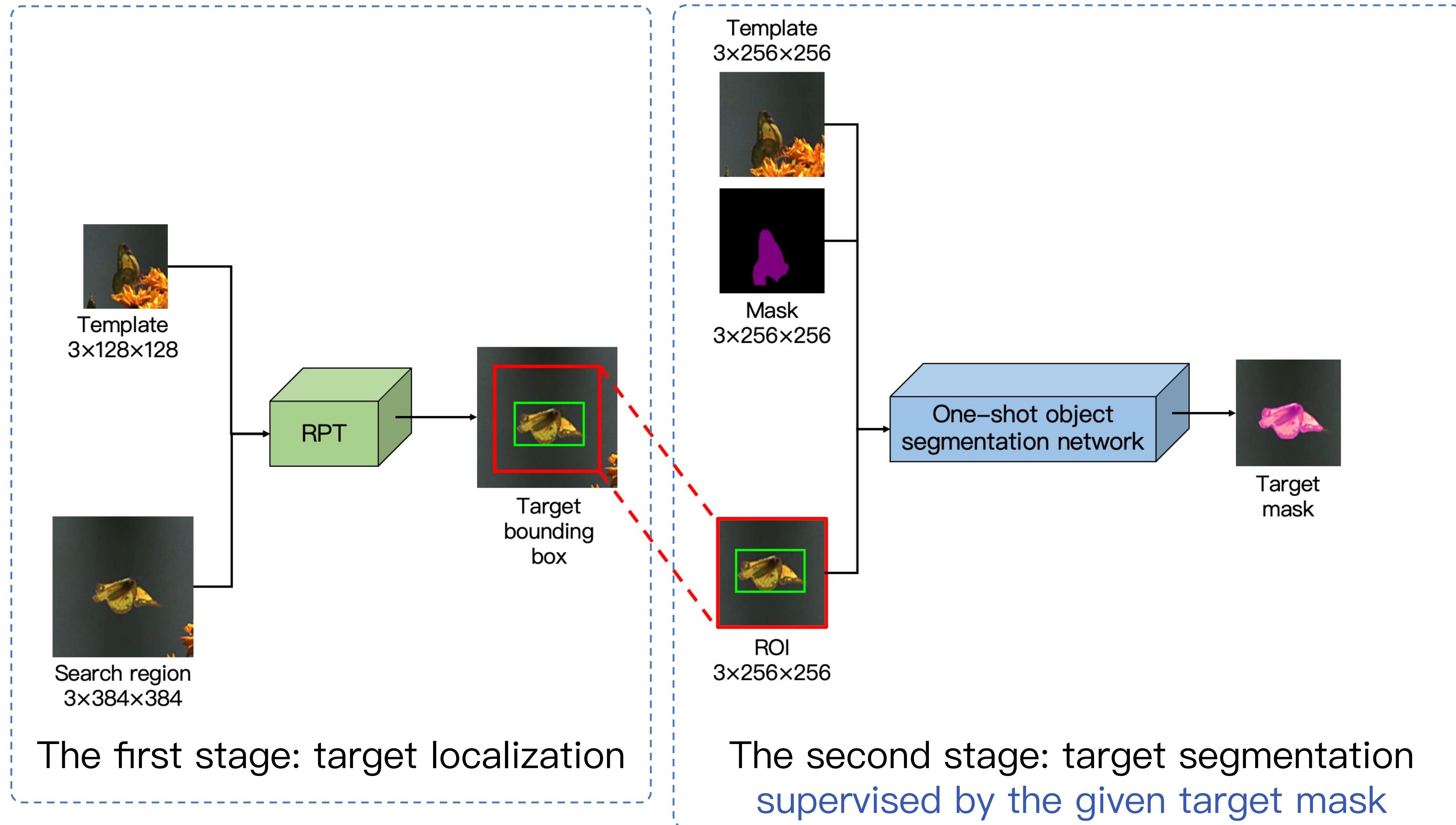
Poor detail segmentation ability



Poor discriminative ability



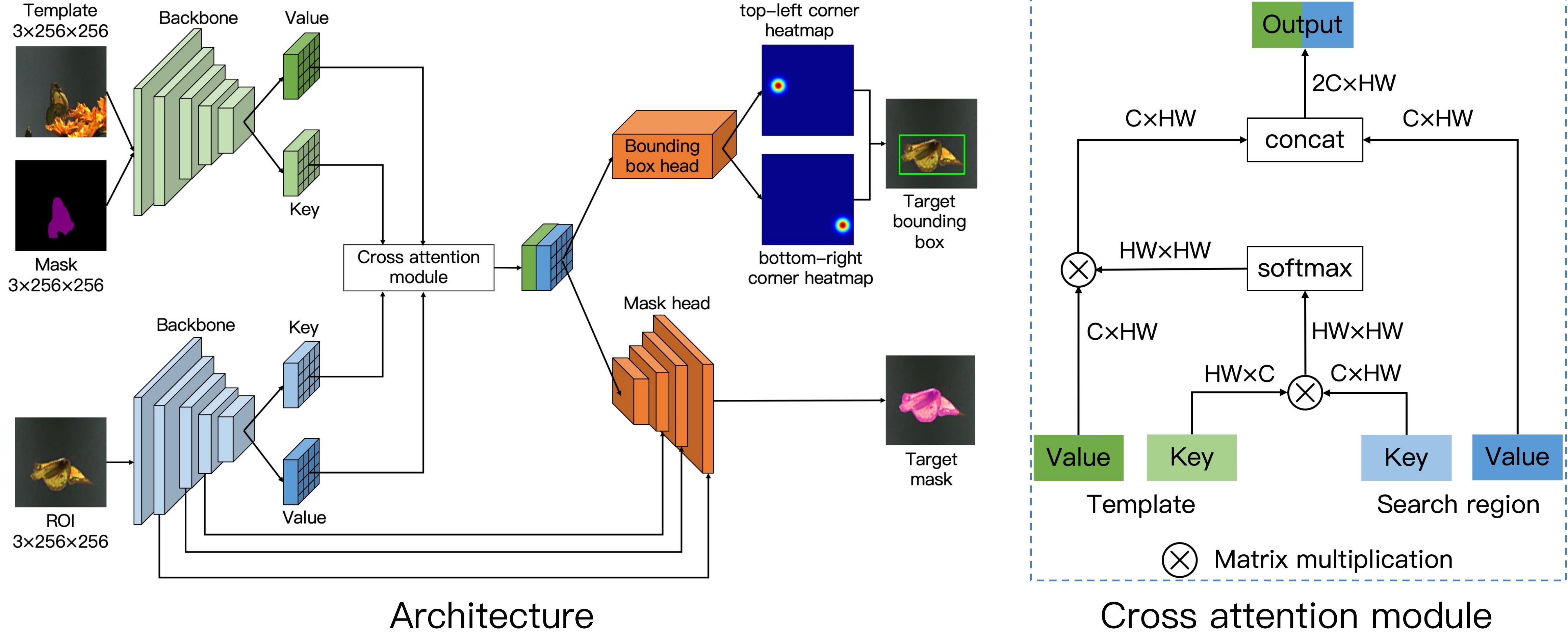
# Pipeline of RPTMask (Ours)



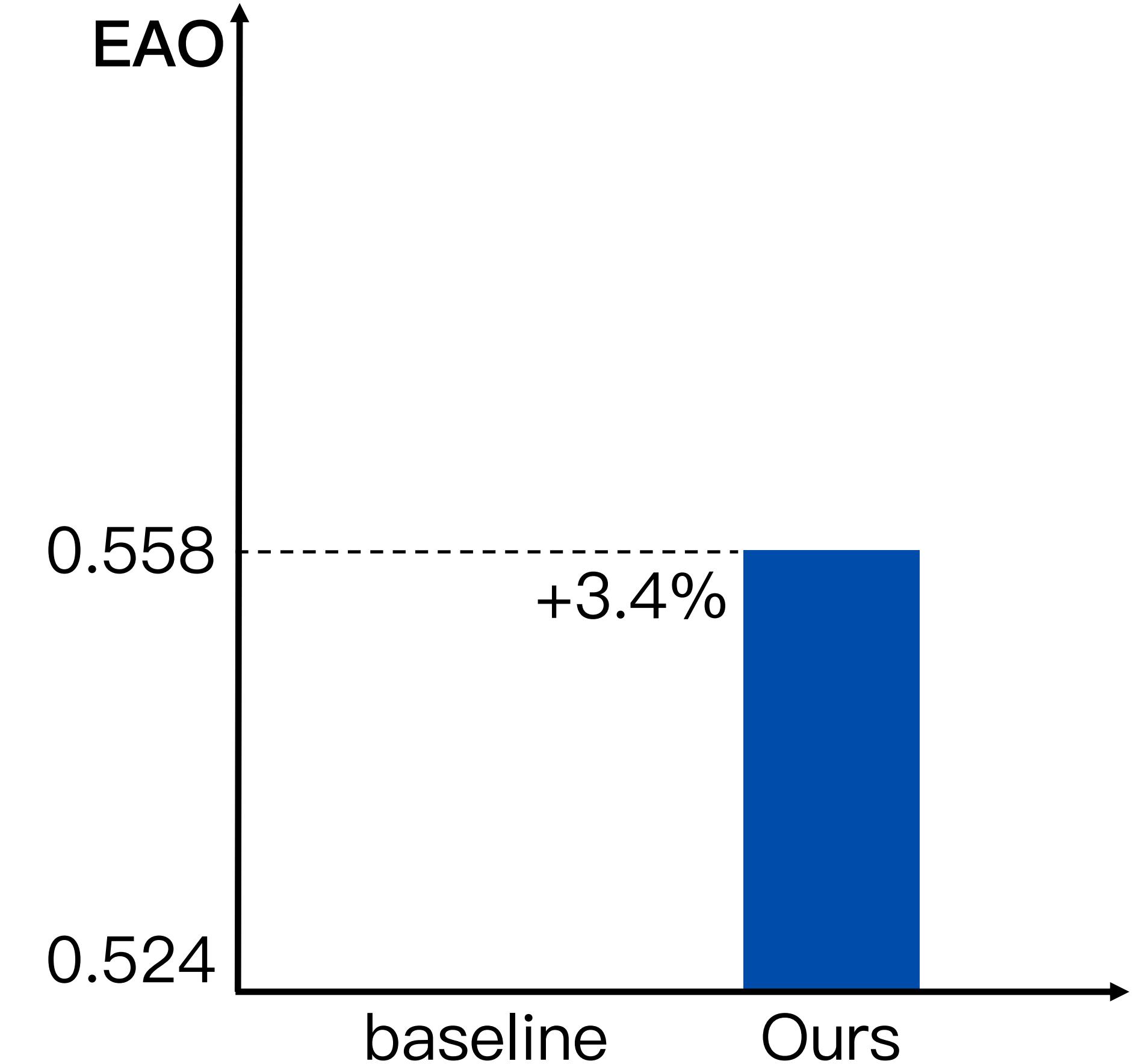
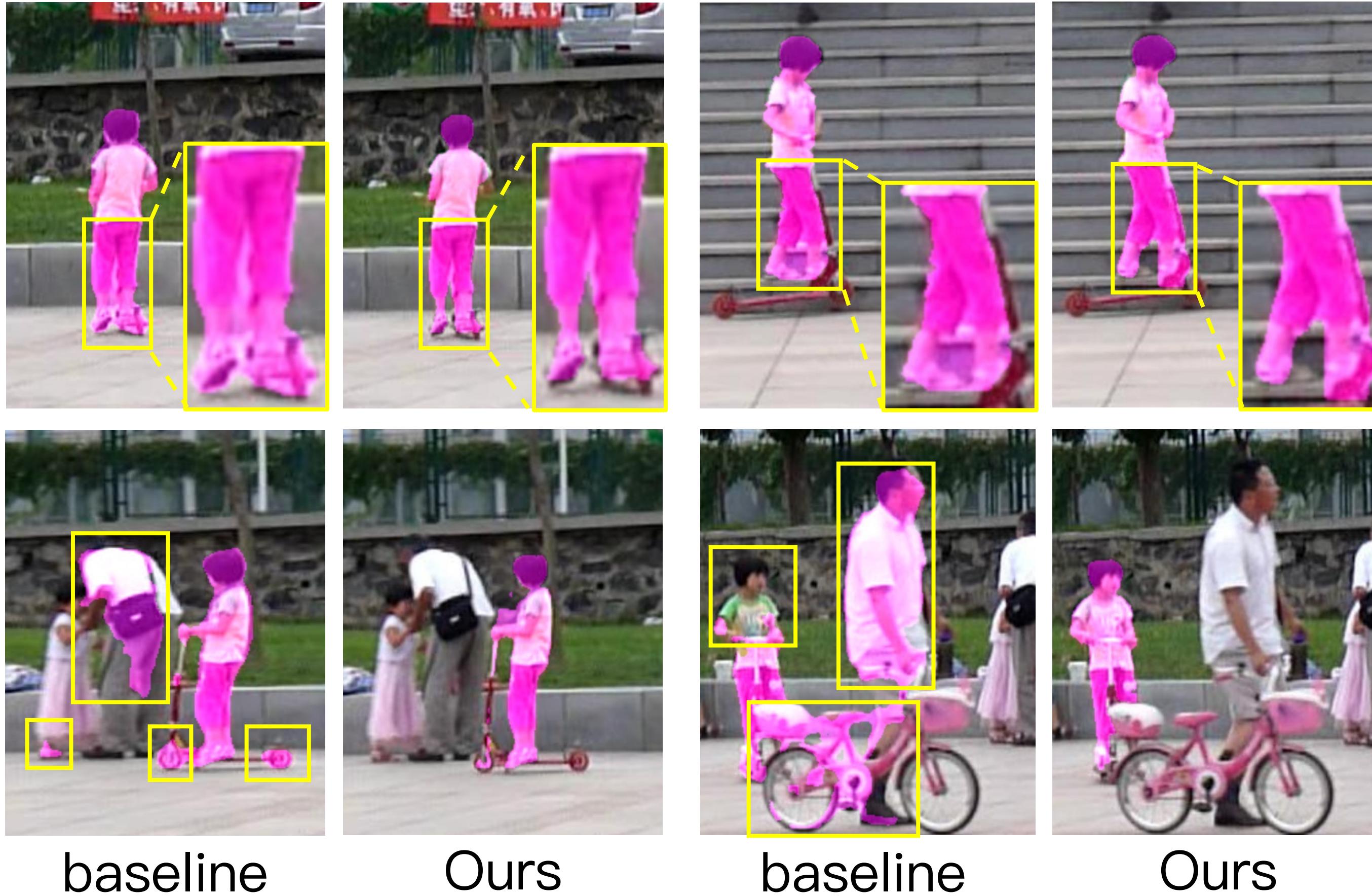
The first stage: target localization

The second stage: target segmentation  
supervised by the given target mask

# One-shot object segmentation network



# Results on the VOT2021 public dataset



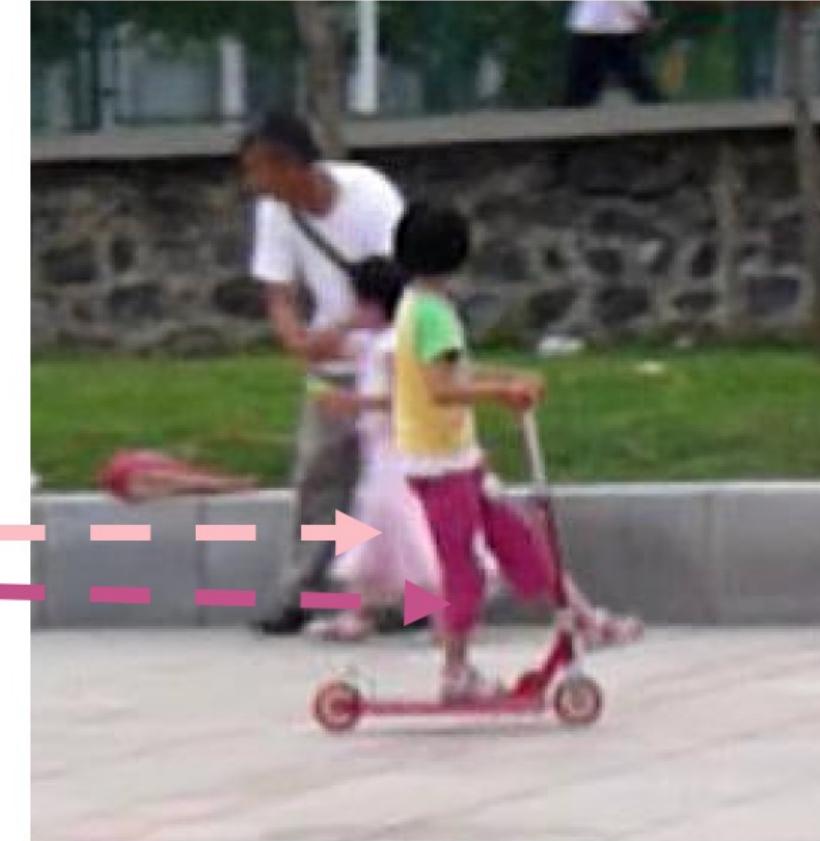
# Motivation #2



Template



ROI



Template



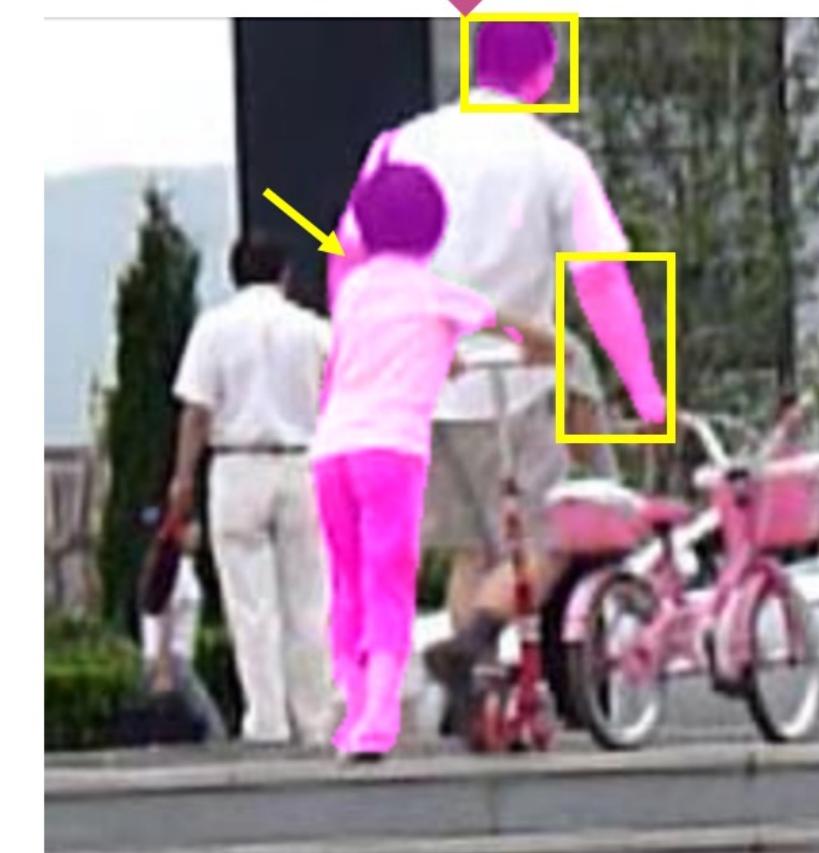
ROI



Output

— Good matching 😊

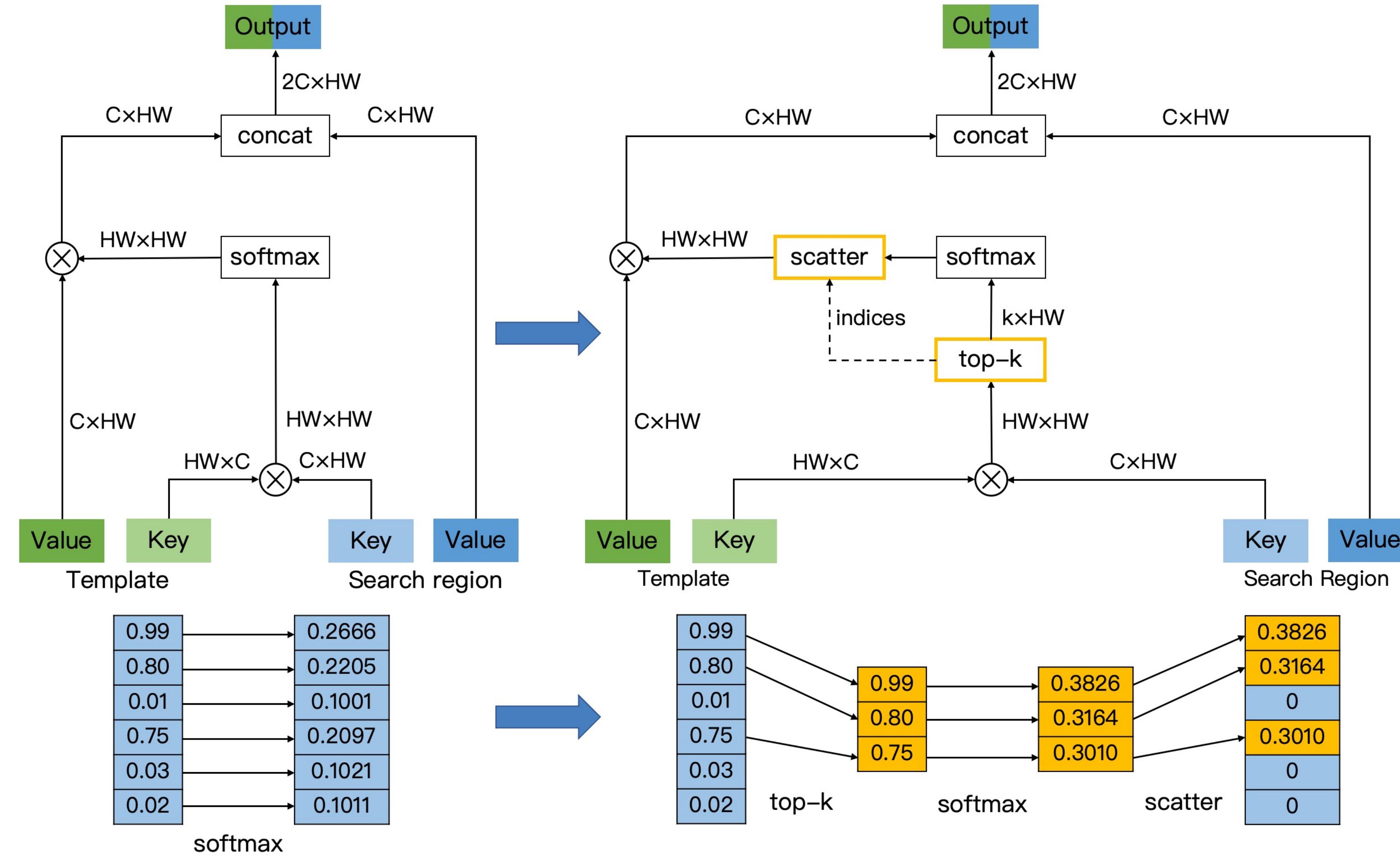
— Bad matching 😞



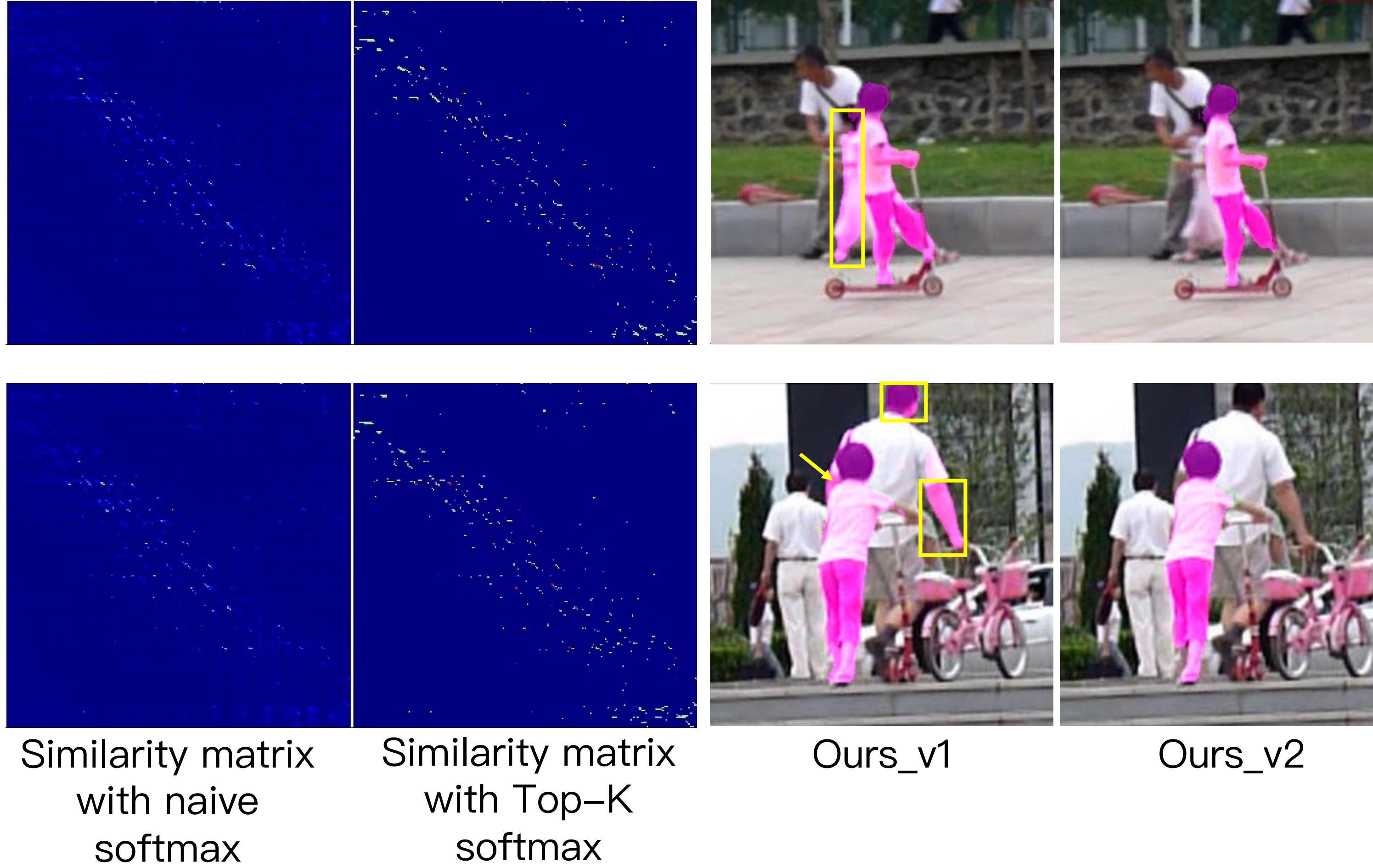
Output

— Undesired response 😞

# Improved cross-attention module



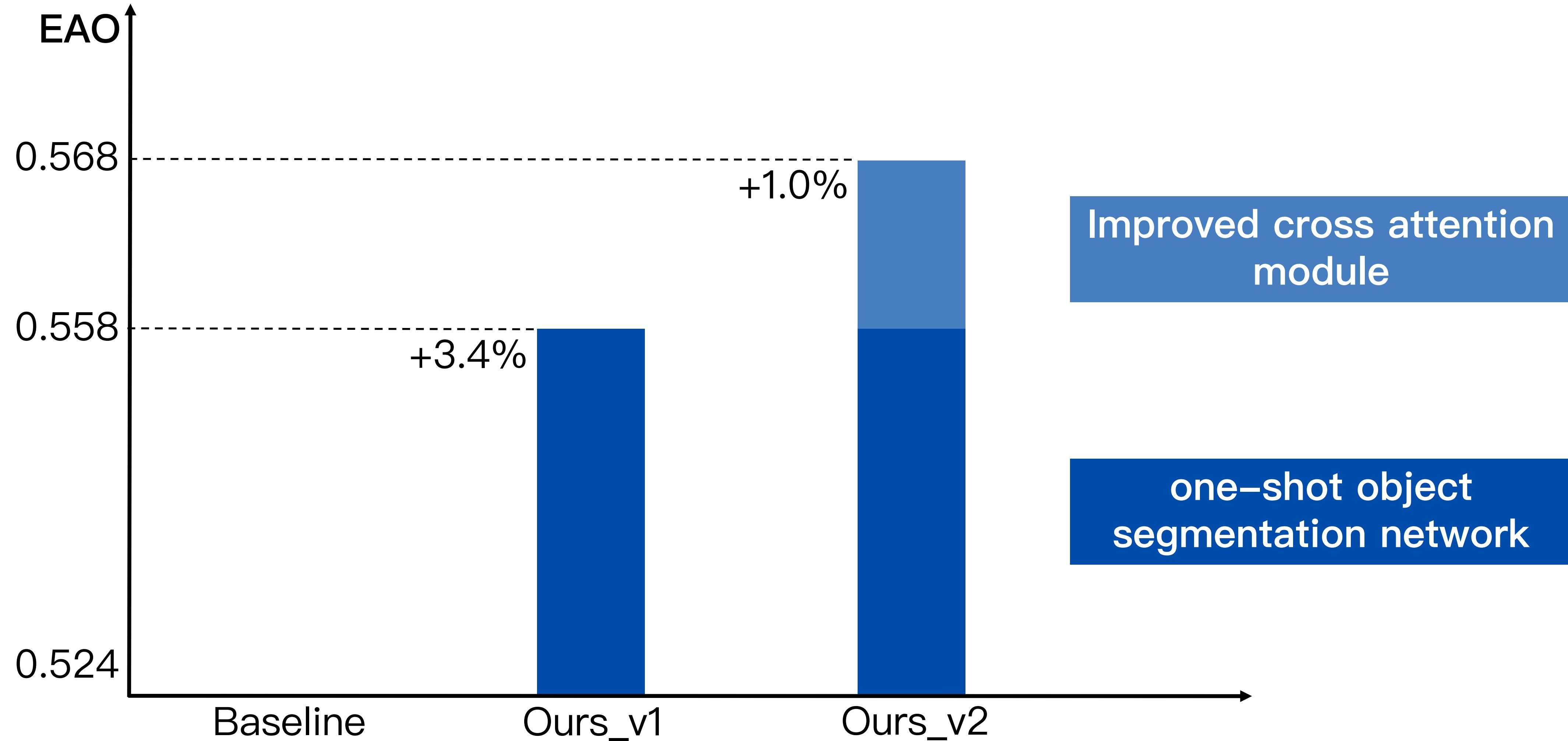
# Results on the VOT2021 public dataset



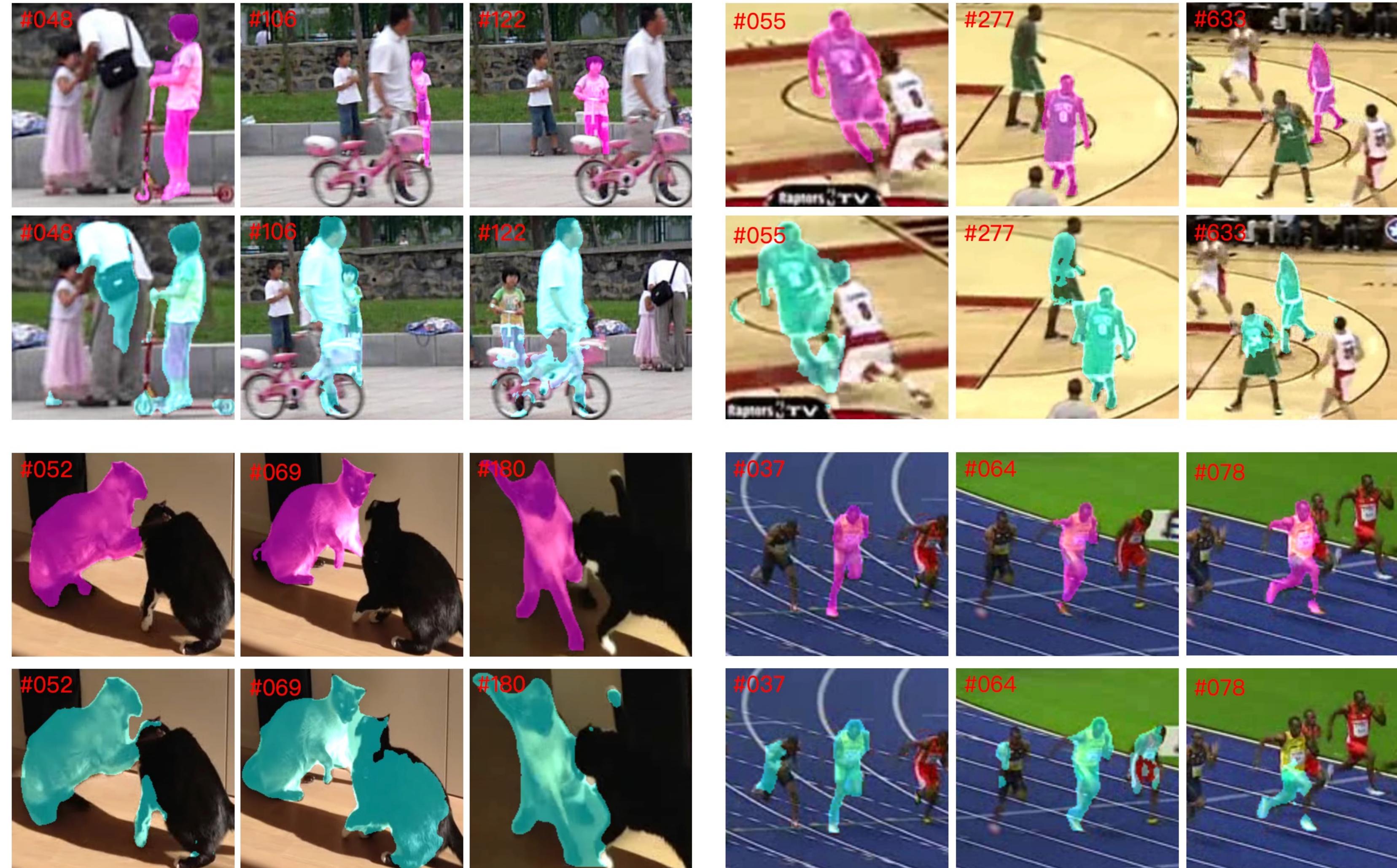
# Quantitative results



VOT2021 public dataset



# Qualitative results



Pink: RPTMask's results

green: baseline's results



# Training dataset

- RPT (we directly use the trained model provided by the VOT2020-ST winner)

- Bounding box

- YouTube-BB
    - COCO
    - ImageNet VID
    - ImageNet DET

- One-shot object segmentation network

- Mask

- YouTube-VOS
    - Saliency

- Bounding box

- GOT-10k
    - LaSOT
    - COCO
    - ImageNet VID
    - ImageNet DET

Thanks