Figure A: The failure examples produced by our proposed method.

1  We thank reviewers for their comments which are very helpful for improving the paper. We address reviewers' comments
2  in the rebuttal and will revise the paper accordingly. All the results are evaluated by $\text{mAP}^r_{0.5}$ unless specified otherwise.
3  **R1:** *About the tightness assumption. How sensitive/robust is the proposed method?* As suggested, we evaluate the
4  proposed method by expanding/contracting the bounding boxes and show its performance under different expan-
5  sion/contraction ratios in following table. The tightness prior is helpful to construct the training bags, and the proposed
6  method is quite robust to the small ratios from $-5\%$ to $+5\%$. Moreover, by slightly contracting the bounding boxes,
7  e.g. $-5\%$, we can reduce noisy pixels in each positive bag, resulting in even better performance. However, excessively
8  expanding or contracting ratios leads to unreliable positive and negative bags, and thus produces sub-optimal results.

9

| Ratio | $+15\%$ | $+10\%$ | $+5\%$ | $0\%$ | $-5\%$ | $-10\%$ | $-15\%$ |
|---|---|---|---|---|---|---|---|
| Ours | 54.0 | 55.0 | 58.0 | 58.9 | 59.5 | 46.0 | 28.1 |

10  **R1:** *About the CVPR'18 paper "Learning ..."* We will cite it, and discuss its similarities and differences from our paper.
11  **R2:** *About average pooling.* Only with the unary term, the result of average pooling is 36.8, which falls behind the
12  result of max pooling, 43.9, because average pooling considers all pixels in the positive bags as the foreground and
13  overestimates the object region. Therefore, using max pooling can better diminish false alarms.
14  **R2:** *Result on COCO and small objects.* The result on *coco minival* is shown in following table (BoxMask is the
15  baseline method in the paper). Our method outperforms the baseline and reaches $78.3\%$ ($= 45.5/58.1$) and $68.3\%$
16  ($= 11.2/16.4$) of the performance of fully-supervised Mask R-CNN in $\text{AP}^r_{0.5}$ and $\text{AP}^r_S$, respectively. The results show
17  that our method is effective in segmenting diverse objects with varied sizes.

18

| method | AP | $\text{AP}^r_{0.5}$ | $\text{AP}^r_{0.75}$ | $\text{AP}^r_S$ | $\text{AP}^r_M$ | $\text{AP}^r_L$ |
|---|---|---|---|---|---|---|
| BoxMask | 11.1 | 31.1 | 6.0 | 5.3 | 11.6 | 15.9 |
| Ours | 21.1 | 45.5 | 17.2 | 11.2 | 22.0 | 29.8 |
| Mask R-CNN [7] | 36.3 | 58.1 | 38.5 | 16.4 | 38.9 | 53.5 |

19  **R2:** *Failure cases.* Figure A shows failure cases of our method. In (a) and (b), segments of small objects are incomplete
20  due to inaccurate boundary on low-resolution regions. In (c) and (d), different instances are wrongly merged. In (e) and
21  (f), inaccurate object contours are segmented due to inter-instance similarity and cluttered scenes.
22  **R3:** *Performance versus annotation effort.* Based on [Bellver et al, Budget-aware Semi-Supervised Semantic and
23  Instance Segmentation, CVPR'19 workshop], the instance-level and box-level annotation costs of Pascal VOC are 239.7
24  and 38.1 seconds per image, respectively. We train Mask R-CNN by instance-level annotation, and limit the amount of
25  annotation so that the annotation budget is comparable with $1.0\times, 1.5\times, 2.0\times$ of the box-level annotation. The results
26  of the three different settings are 48.3, 53.5 and 59.9, respectively. The first two results fall behind our method by the
27  margins 10.6 and 5.4, while the last one surpasses ours by 1.0 but with $2\times$ annotation cost. Therefore, with the same
28  annotation cost, our method outperforms Mask R-CNN because less training data lead to overfitting for Mask R-CNN.
29  **R3:** *The four questions about Eq.4.* (1) Epsilon is an 8-neighbor set. (2) Because neighboring pixels are connected in
30  the pairwise term, we can propagate the segment scores, and thus enlarge the segment. (3) We assume "patches" in the
31  review refers to "bags" in our method. If all instances in positive bags are treated as positive samples, many background
32  pixels are mistaken as positive samples. Therefore, more false alarms could be predicted. (4) With the tightness prior,
33  each positive bag meets the MIL assumption, so this task could be solved with the MIL formulation.
34  **R3:** *Comparing efficiency with [16].* The method [16] requires pre-generated proposals from Selective Search, obtains
35  the detected boxes by Fast R-CNN, and finally takes the detected boxes and the RGB images as the input to generate
36  the instance result. In [16], the proposal generation step alone takes $\approx 10$ seconds per image. In contrast, our method
37  requires only one forward pass which takes only $\leq 0.1$ seconds per image or $\approx 1$ second per image if DenseCRF is
38  applied. Therefore, our method is more efficient than [16].
39  **R3:** *About L85.* We agree with this comment, and will make it more clear in the revision. However, what we want to
40  claim is that in general, the applicability of the fully supervised methods "may" be limited in the real world because of
41  the high annotation cost.
42  **R3:** *Comparing performance with fully supervised Mask-RCNN.* Our method adopting the MIL formulation tends to
43  highlight the discriminative parts of objects, while Mask R-CNN with mask-level annotation emphasizes the whole
44  objects. The IOUs between the ground truth and the discriminative regions are often larger than 0.25 but less than
45  0.5. It is why our method slightly outperforms Mask R-CNN in $\text{mAP}^r_{0.25}$, but falls behind it in $\text{mAP}^r_{0.5}$, $\text{mAP}^r_{0.7}$, and
46  $\text{mAP}^r_{0.75}$, as reported in Table 1 of the submitted manuscript.