

# Deep Structured Instance Graph for Distilling Object Detectors

## – Supplementary Material –

Yixin Chen<sup>1</sup>, Pengguang Chen<sup>1</sup>, Shu Liu<sup>2</sup>, Liwei Wang<sup>1</sup>, Jiaya Jia<sup>1,2</sup>

The Chinese University of Hong Kong<sup>1</sup> SmartMore<sup>2</sup>

Here we provide some additional details about what was left in the main text. First we show the algorithm of our method involved in the whole distillation procedure. Next, we elaborate the evaluation metrics of the COCO [2] benchmark and the student training configurations and settings. We also give definition of Kullback-Leibler Divergence Loss that was not mentioned in the submitted paper. At last, visualization results are provided to afford more intuitive understanding.

### 1. Algorithm

We show the pseudo algorithm of our method involved in detection models, including structured graph building and graph distillation loss calculation:

---

**Algorithm 1** Graph Distillation Loss

---

**Input:** image  $x$ , target  $t$ , teacher  $T$ , student  $S$

**Output:** Graph Distillation Loss

- 1:  $\mathcal{F}_T \leftarrow T.\text{backbone}(x)$
  - 2:  $\mathcal{F}_S \leftarrow S.\text{backbone}(x)$
  - 3:  $\text{proposals} \leftarrow S.\text{rpn}(x, t, \mathcal{F}_S)$
  - 4:  $\text{Node}_T^{fg}, \text{Node}_T^{bg} \leftarrow \text{RoIPool}(\mathcal{F}_T, \text{proposals})$
  - 5:  $\text{Node}_S^{fg}, \text{Node}_S^{bg} \leftarrow \text{RoIPool}(\mathcal{F}_S, \text{proposals})$
  - 6:  $\text{Node}_T^{mine}, \text{Node}_S^{mine} \leftarrow \text{BackGroundMining}(\text{Node}_T^{bg}, \text{Node}_S^{bg}, T)$
  - 7:  $\text{Node}_T \leftarrow \text{cat}(\text{Node}_T^{fg}, \text{Node}_T^{mine})$
  - 8:  $\text{Node}_S \leftarrow \text{cat}(\text{Node}_S^{fg}, \text{Node}_S^{mine})$
  - 9:  $\text{Edge}_T \leftarrow \text{similarity\_function}(\text{Node}_T)$
  - 10:  $\text{Edge}_S \leftarrow \text{similarity\_function}(\text{Node}_S)$
  - 11:  $L_{Node} \leftarrow \text{mseloss}(\text{Node}_T, \text{Node}_S)$
  - 12:  $L_{Edge} \leftarrow \text{mseloss}(\text{Edge}_T, \text{Edge}_S)$
  - 13: **return**  $L_{Node} + L_{Edge}$
- 

where the algorithm of function BackGroundMining is detailed as follows:

---

**Algorithm 2** BackGroundMining

---

**Input:**  $\text{Node}_T^{bg}, \text{Node}_S^{bg}$ , Teacher  $T$

**Output:** Mined Background Nodes  $\text{Node}_T^{mine}, \text{Node}_S^{mine}$

- 1: **Initialize**  $\text{thrs} = t$

- 2:  $\text{RoILoss}_T^{bg} \leftarrow T.\text{RoIclsLoss}(\text{Node}_T^{bg})$
  - 3:  $\text{NodeIndex} \leftarrow \text{SelectAboveThrs}(\text{RoILoss}_T^{bg}, \text{thrs})$
  - 4:  $\text{Node}_T^{mine} \leftarrow \text{SelectIndex}(\text{Node}_T^{bg}, \text{NodeIndex})$
  - 5:  $\text{Node}_S^{mine} \leftarrow \text{SelectIndex}(\text{Node}_S^{bg}, \text{NodeIndex})$
  - 6: **return**  $\text{Node}_T^{mine}, \text{Node}_S^{mine}$
- 

### 2. Experimental Details

We skip some details in the experiment of the main paper due to pages limitation. Here we present a brief description about the COCO evaluation metrics and the training settings.

#### 2.1. Evaluation Metrics

Here we adopt the evaluation metric Average Precision(AP) in our experiments. AP is a quite popular evaluation metric in measuring object detector’s performance, here we only refer to the COCO style AP. Before that, we first introduce IoU(Intersection over union) to determine whether one prediction is correct w.r.t. a ground truth object or not. Specifically, IoU defines the intersection between the predicted bounding box and ground truth bounding box:

$$IOU = \frac{\text{area of overlap}}{\text{area of union}} \quad (1)$$

We also set a IoU *threshold* to classify whether a prediction is true positive or false positive. Precision actually means the percentage of true positive predictions in all predictions while recall is the percentage of true positive predictions in all positive samples:

$$\begin{aligned} \text{Precision} &= \frac{TP}{TP + FP} \\ \text{Recall} &= \frac{TP}{TP + FN} \end{aligned} \quad (2)$$

where  $TP, FP, FN$  represents true positive, false positive, and false negative respectively.

We rank all predictions in a descending order according to their confidences. With every rank, we can calculate the

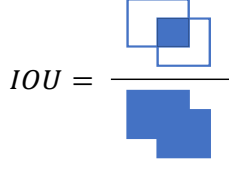


Figure 1. IoU

Average Precision (AP):	
AP	% AP at IoU=.50:.95 (primary challenge metric)
AP <sub>IoU=.50</sub>	% AP at IoU=.50 (PASCAL VOC metric)
AP <sub>IoU=.75</sub>	% AP at IoU=.75 (strict metric)
AP Across Scales:	
AP <sub>small</sub>	% AP for small objects: area < 32 <sup>2</sup>
AP <sub>medium</sub>	% AP for medium objects: 32 <sup>2</sup> < area < 96 <sup>2</sup>
AP <sub>large</sub>	% AP for large objects: area > 96 <sup>2</sup>

Figure 2. AP series. <https://cocodataset.org/#detection-eval>.

precision and recall, which will finally form a Precision-Recall curve. In the COCO style, we sample 101 points recall@[0:0.01:1] and average all these corresponding precision to calculate mean AP. With different IoU *threshold*, true positive predictions can be different. In our experiments, we adopt six different kinds of AP according to various IoU *thresholds* and object scales, descriptions are detailed in Figure 2.

## 2.2. Training Settings

For training, all our experiments are performed on 4 Nvidia RTX 2080Ti GPUs, and all the student models are sticking to the 1x/3x(around 12/37 epochs) COCO training schedule. Specifically, we set batch size as 8 with 2 images allocated to each GPU and double the training iterations to 180000/540000 for 1x/3x training. Our student models are optimized with Stochastic Gradient Descent(SGD) with the base learning rate 0.01 and we reduce it by gamma 0.1 at iteration (120000, 160000)/(420000, 500000) for 1x/3x training. All our backbone models are initialized with weights pre-trained on ImageNet.

## 3. Kullback-Leibler(KL) Divergence Loss

The definition of KLD loss is as follows:

$$D_{KL}(p_t|p_s) = - \sum_i p(x_{t,i}) \log p(x_{s,i})$$

$$= \sum_{x_t, x_s} p(x_{t,i}) \log \frac{p(x_{t,i})}{p(x_{s,i})} \quad (3)$$

where  $x_t, x_s$  represents the logits of teacher and student respectively, and probability

$$p(x_{.,i}) = \frac{\exp(x_{.,i})}{\sum_j \exp(x_{.,j})} \quad (4)$$

Actually, in our experiment, following [1], we use the logits softened by temperature  $T$ , so the probability is:

$$\hat{p}(x_{.,i}) = \frac{\exp(x_{.,i}/T)}{\sum_j \exp(x_{.,j}/T)} \quad (5)$$

since the gradients produced by soft logits are also scaled as  $1/T^2$ , we should multiply them by  $T^2$ , so the KLD loss is transformed to:

$$D_{KL}(p_t|p_s) = - \sum_i \hat{p}(x_{t,i}) \log \hat{p}(x_{s,i}) * T^2$$

$$= \sum_{x_t, x_s} \hat{p}(x_{t,i}) \log \frac{\hat{p}(x_{t,i})}{\hat{p}(x_{s,i})} * T^2 \quad (6)$$

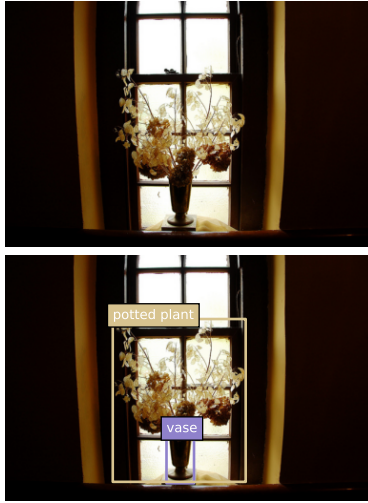
In our experiments, we simply set  $T$  to 1.

## 4. Visualization Results

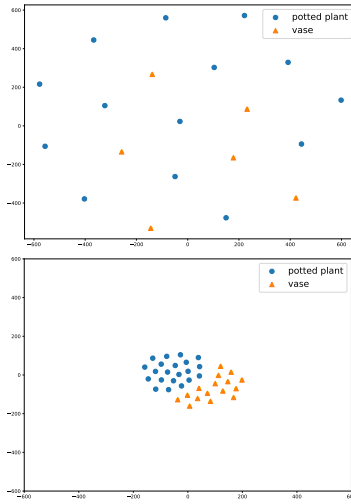
**Graph Visualization:** We show more visualization results on COCO test images. As is shown in Figure 3. It can be seen that from node visualization, the scattered nodes from several categories are mixing with each other, which means student’s feature space is quite different from that of teacher. Also, in edge visualization, some edges in student still reveal strong intensities among nodes in different categories while edges in teacher almost show strong links among nodes in the same categories. These appearing visualization results enlighten us that there are great gaps between the knowledge space of student and teacher, which propels us to conduct graph transferring to close that.

**Object Detection Visualization:** To intuitively illustrate how our distillation method makes a difference to improve student detector’s performance, we visualize student detection results before distillation and the counterparts distilled after. See Figure 4. Before distillation, student detectors with backbones in small capacity show poor results including missing boxes, duplicate boxes, and inaccurate boxes *etc.* However, refined with our method, students with substantial promotions are able to bridge the gaps and behave much more accurate in classification and localization.

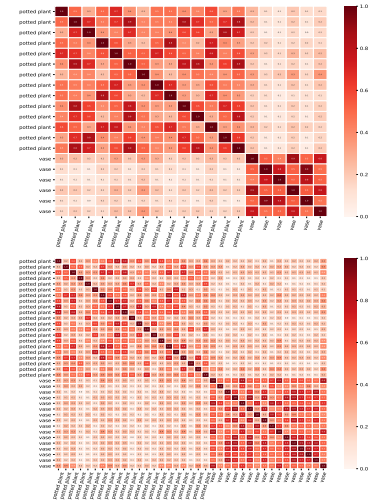
**Instance Segmentation Visualization:** We also show the visualization of distillation performance in instance segmentation. As is shown in Figure 5, the students are not behaving well in classification and localization, and even show the poor segmentation results in each predicted mask. As for teachers, the segmentation results seem better, yet some wrongly predicted pixels still exist. Since our method has added more regularization than normal baseline models, the masks in our method have exhibited better shapes and finer borders than the student baseline, even beyond the teacher in some test cases.



(a) Labeled Image



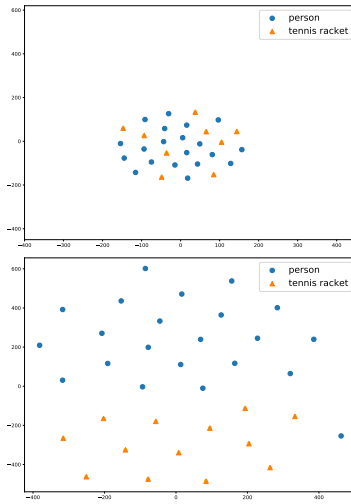
(b) Node Visualization



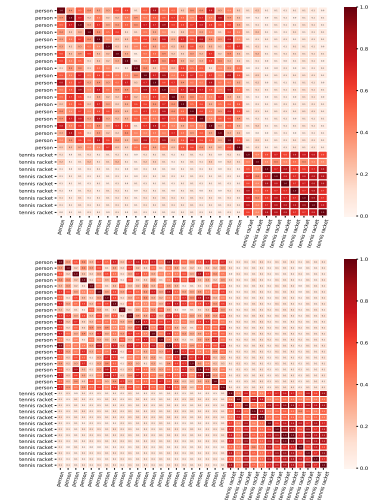
(c) Edge Visualization



(d) Labeled Image



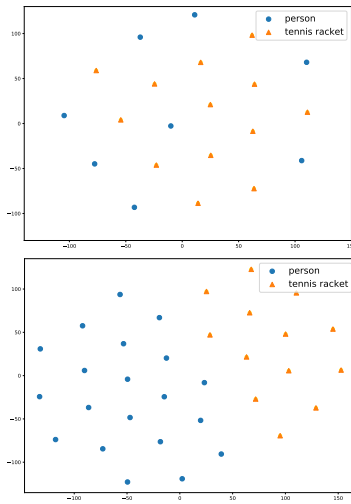
(e) Node Visualization



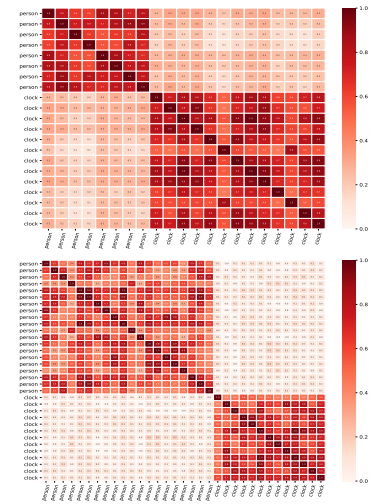
(f) Edge Visualization



(g) Labeled Image



(h) Node Visualization



(i) Edge Visualization

Figure 3. Graph visualization on test image.



Figure 4. Object detection distillation visualization on COCO2017 test. We adopt (a) Student Faster R-CNN with ResNet18, (b) Student distilled using our method, and (c) Teacher ResNet50 as our trained models for visualization.

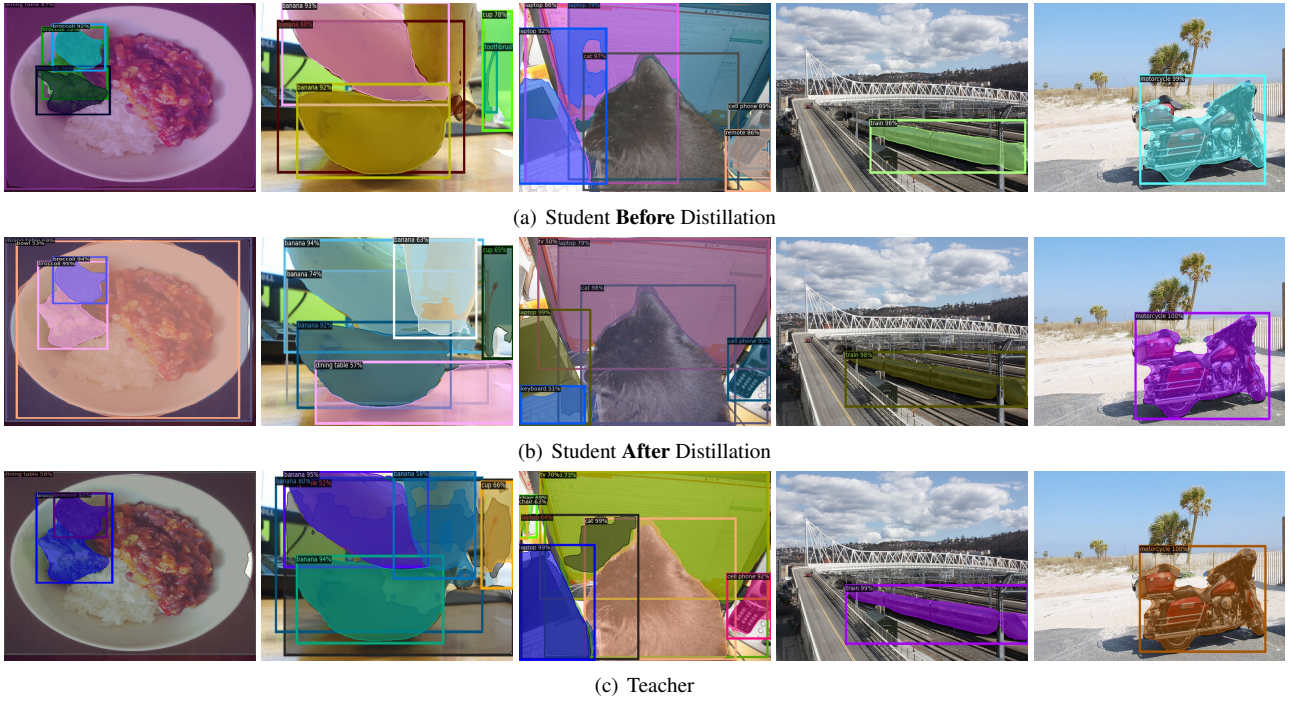


Figure 5. Instance segmentation distillation visualization on COCO2017 test. We adopt (a) Student Faster R-CNN with ResNet18, (b) Student distilled using our method, and (c) Teacher ResNet50 as our trained models for visualization.

## References

- [1] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. *ArXiv*, abs/1503.02531, 2015. [2](#)
- [2] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *ECCV*, 2014. [1](#)