Cross-Dataset Collaborative Learning for Semantic Segmentation

Li Wang¹, Dong Li¹, Yousong Zhu², Lu Tian¹, Yi Shan¹ ¹ Xilinx Inc., Beijing, China. ² Institute of Automation, Chinese Academy of Sciences, Beijing, China. {liwa, dongl, lutian, yishan}@xilinx.com, yousong.zhu@nlpr.ia.ac.cn

CVPR 2021



Cross-Dataset Collaborative Learning for Semantic Segmentation

• Proposed a simple, flexible, and general cross-dataset collaborative learning algorithm.



(c) Cross-dataset training with label remapping

(d) Cross-dataset collaborative training (ours)

Analysis

- (a) Train and evaluate a network for each dataset separately
- (b) Finetuning
 - Time-consuming
 - Not applicable for joint optimization
 - Manually adjust hyper-parameters
- (c) Generate a hybrid dataset
 - Simple label concatenation
 - Label mapping



(c) Cross-dataset training with label remapping

(d) Cross-dataset collaborative training (ours)

Label mapping

- Prior knowledge needed
 - Class duplication and conflict
- Poor performance
 - Discrepancy of camera viewpoints, scenes, etc.



Figure 3: A simple example of label mapping. Two datasets with labels l_1, l_2, l_3, l_4, l_5 (labels are represented by circles) and m_1, m_2, m_3 (represented by rectangles) are merged into a new hybrid dataset with labels n_1, n_2, \dots, n_7 (represented by round rectangles). The solid shapes in blue indicate that the labels belong to the same class that can be merged into one class in label mapping operation. Other labels stay unchanged during mapping. https://blog.csdn.net/stezio

Cross-Dataset Collaborative Learning

- Different segmentation datasets vary greatly
 - (e.g., scenes and illumination)



(a) Visualization for samples of different datasets

Rethinking BatchNorm and Convolution

- Analyze the **parameter** distributions of conv and bn layers.
- Conv,
 - weight, bias from conv kernel
- BN,
 - running_var, running_mean, weight, bias

$$y = rac{x - \mathrm{E}[x]}{\sqrt{\mathrm{Var}[x] + \epsilon}} * \gamma + eta \qquad \hat{x}_{\mathrm{new}} = (1 - \mathrm{momentum}) imes \hat{x} + \mathrm{momentum} imes x_t$$

where x^ is the estimated statistic and x_t is the new observed value

Rethinking BatchNorm and Convolution



Rethinking BatchNorm and Convolution



- For Conv, weights hold the same distributions
- For BN, distributions of both running mean and running var have different shapes for different datasets

Network



Dataset-Aware Block

- Dataset-invariant conv layer
- Dataset-specific bn layers
- A switch is automatically to determine which BN should be activated based on the data source (if statement?)
- Dataset-aware classifiers (conv with different out channel?)

$$DSBN\{D_i\}(X_i;\gamma_i,\beta_i) = \gamma_i \hat{X}_i + \beta_i$$
(1)

where,

$$\hat{X}_{i} = \frac{X_{i} - \mu_{i}}{\sqrt{\sigma_{i}^{2} + \epsilon}}$$

$$\mu_{i} = \frac{1}{B} \sum_{j=1}^{B} X_{i}^{j}, \ \sigma_{i}^{2} = \frac{1}{B} \sum_{j=1}^{B} (X_{i}^{j} - \mu_{i})^{2}$$
(2)



Dataset-Aware Block (DAB)

Dataset Alternation Training

- One iteration
 - Construct the batch with samples from a single dataset, compute loss
- Next iteration
 - Samples from another dataset, compute loss
- Dataset-number iteration
 - •
 - Accumulate all losses and backpropagate



Dataset Alternation Training

- Backpropagate the loss of each dataset in each iteration will incur training instability.
- Efficient way to train the samples from multiple datasets



Result

| Method | Cityscapes (%) | | BDD100K(%) |
|-----------------|----------------|-------|------------|
| | Val. | Test | Val. |
| Single-dataset | 67.52 | 67.75 | 53.88 |
| Finetuning | 67.79 | 66.52 | 58.30 |
| Label remapping | 66.23 | 66.39 | 58.74 |
| CDCL (Ours) | 72.63 | 71.55 | 60.47 |

(a) Cityscapes + BDD100K

| Method | Cityscapes (%) | | CamVid (%) | | |
|-------------------------|-------------------------|-------|------------|-------|--|
| | Val. | Test | Val. | Test | |
| Single-dataset | 67.52 | 67.75 | 73.05 | 70.41 | |
| Finetuning | 67.35 | 67.87 | 74.83 | 71.16 | |
| Label remapping | 67.13 | 68.22 | 78.03 | 76.86 | |
| CDCL (Ours) | 69.77 | 68.56 | 78.45 | 77.34 | |
| (b) Cityscapes + CamVid | | | | | |
| Method | Cityscapes (%) COCO (%) | | | 0 (%) | |
| | Val. Test Val. | | | al. | |
| Single-dataset | 67.52 | 67.75 | 32 | 32.86 | |
| Finetuning | 70.44 | 70.23 | 32 | .10 | |
| Label remapping | 50.35 | 52.39 | 32 | 32.67 | |
| CDCL (Ours) | 72.63 | 72.52 | 32 | .87 | |

(c) Cityscapes + COCO

Result

| Method | Cityscapes (%) | | CamVid (%) | | BDD100K (%) |
|-------------------------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|
| | Validation | Test | Validation | Test | Validation |
| Single-dataset CDCL (Ours) | 67.52 73.17 (+5.65) | 67.75 70.98 (+3.23) | 73.05 78.84 (+5.79) | 70.41 75.52 (+5.11) | 53.88 60.45 (+6.57) |

Result

| Method | Norm | DAT | Cityscapes | BDD100K |
|----------------|------|--------------|------------|---------|
| Single-dataset | BN | | 67.75% | 53.88% |
| Cross-dataset | BN | | 62.10% | 53.34% |
| Cross-dataset | BN | \checkmark | 64.69% | 56.33% |
| Cross-dataset | DSBN | | 68.55% | 58.93% |
| CDCL (Ours) | DSBN | \checkmark | 72.63% | 60.47% |

Table 3: Ablation studies on DSBN and DAT with the ResNet-18 backbone on the Cityscapes and BDD100K validation sets.

