# Self-Tuning for Data-Efficient Deep Learning

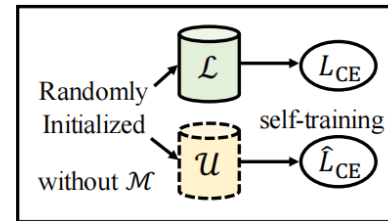**Ximei Wang** [*1]  **Jinghan Gao** [*1]  **Mingsheng Long** [1]  **Jianmin Wang** [1]

School of Software, BNRist, Tsinghua University, Beijing, China, 100084

ICML 2021

# SSL & TL

- SSL (pseudo labels) and Confirmation bias
  - Without a decent pre-trained model to provide an implicit regularization, will be easily misled by inaccurate pseudo-labels, especially in large-sized label space.
  - SSL need a well pre-trained model

- TL and Model shift
  - The fine-tuned model shifts towards the limited labeled data and leaves away from the original smooth model pre-trained on a large-scale datasets
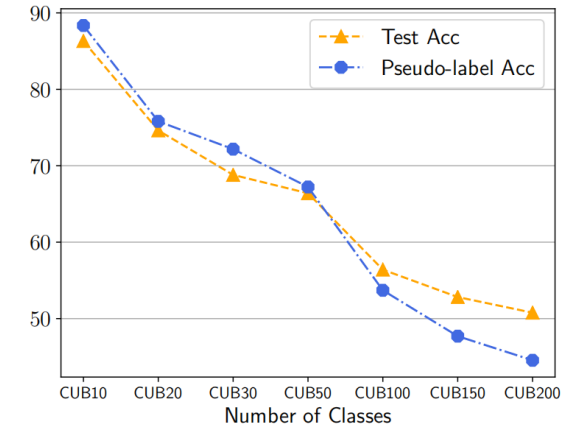  - Utilizing unlabeled data to alleviate the model shifting



(a) **Transfer Learning**        (b) **Semi-supervised Learning**
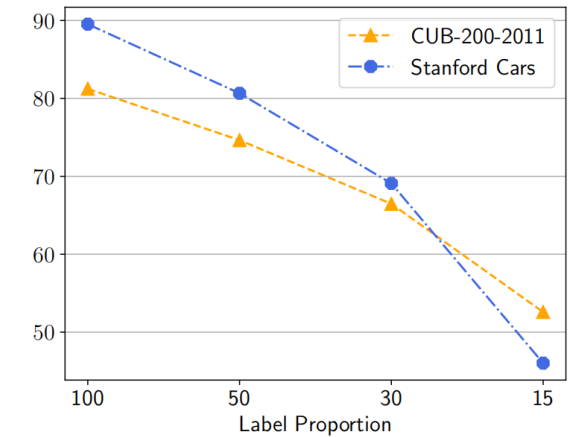
SSL and TL are complementary
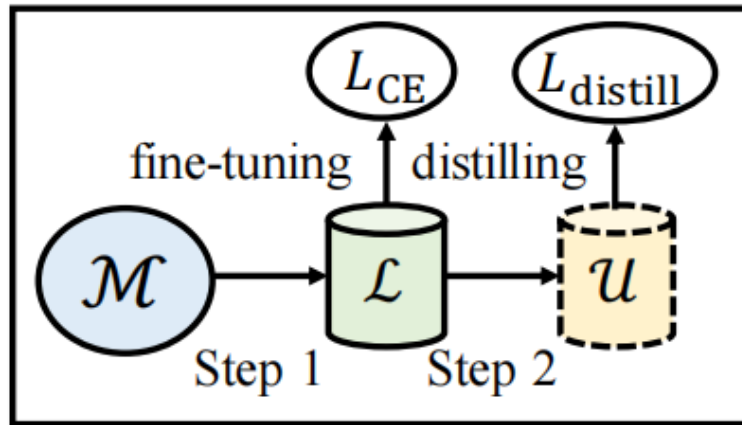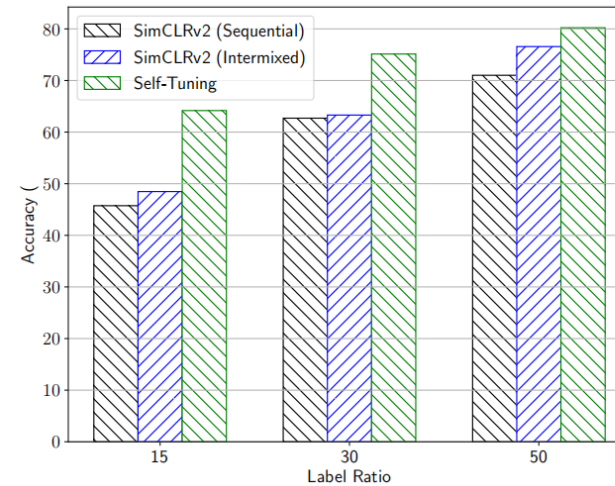
SSL



(a)  Acc of FixMatch on $CUB$

TL



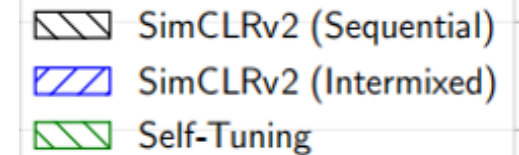(b)  Test accuracy of Co-Tuning

# SimCLRv2

- Both utilize labeled data and unlabeled data

- The fine-tuned model would easily shift towards the limited labeled data with sampling bias and leaves away from the original smooth model pre-trained on a large-scale dataset



(c) **SimCLRv2**



(b) Compare with SimCLRv2

# Confirmation bias: CE & CL loss

- **Cross Entropy**

    The model trained by CE loss will be easily confused by false pseudo-labels since it focuses on learning a hyperplane for discriminating each class from the other classes (CE loss overfitting easily)

- **Contrast Loss**

    While standard CL loss lacks a mechanism to tailor pseudo-labels into model training, leaving the useful discriminative information on the shelf. (CL doesn't take classes into account)



**CE**: Directly mislead a hyperplane   **CL**: No hyperplane is learnt

— Learnt Hyperplane   — True Hyperplane   ●▲■ Different Classes   ●▲■ Unlabeled Data   ● False Pseudo Labels   ····· Positive Key   ····· Negative Key

# Confirmation bias Solution-PGC

- Different from the standard CL which involves just a positive key in each contrast, PGC introduces a group of positive keys in the same pseudo-class to contrast with all negative keys from other pseudo-classes.



Queue size

Two different augments of one sample

$$\widehat{L}_{\text{PGC}} = -\frac{1}{D+1} \sum_{d=0}^{D} \log \frac{\exp(\mathbf{q} \cdot \mathbf{k}_d^{\widehat{y}}/\tau)}{\text{Pos} + \text{Neg}}$$

Unlabeled Data

$$\text{Pos} = \exp(\mathbf{q} \cdot \mathbf{k}_0^{\widehat{y}}/\tau) + \sum_{j=1}^{D} \exp(\mathbf{q} \cdot \mathbf{k}_j^{\widehat{y}}/\tau)$$

The second augments of queried sample

$$\text{Neg} = \sum_{c=1}^{\{1,2,\cdots,C\}\setminus\widehat{y}} \sum_{j=1}^{D} \exp(\mathbf{q} \cdot \mathbf{k}_j^{c}/\tau),$$

Labeled Data

$$L_{\text{PGC}} = -\frac{1}{D+1} \sum_{d=0}^{D} \log \frac{\exp(\mathbf{q} \cdot \mathbf{k}_d^{y}/\tau)}{\text{Pos} + \text{Neg}},$$

Loss

$$\mathbb{E}_{(\mathbf{x}_i, y_i)|\in\mathcal{L}} (L_{\text{CE}} + L_{\text{PGC}}) + \mathbb{E}_{(\mathbf{x}_i)\in\mathcal{U}} \widehat{L}_{\text{PGC}}.$$

queue list

**PGC**: Mitigate the reliance on pseudo-labels

— Learnt Hyperplane  — True Hyperplane  ●▲■ Different Classes
●▲■ Unlabeled Data  ● False Pseudo Labels  ····Positive Key  ····Negative Key

# Model Shift Solution-Unifying and Sharing

- By utilizing Unlabeled data at the same time in a unified form as shown the model shift challenge is expected to be alleviated.

- Shared queue improves the accuracy keys for unlabeled queries than that of a separate queue for unlabeled data.

- Self-Tuning has a better starting point to provide an implicit regularization than the model trained from scratch on the target dataset.

# Experiments

| Dataset | Type | Method | Label Proportion | | | |
|---|---|---|---|---|---|---|
| | | | 15% | 30% | 50% | 100% |
| CUB-200-2011 | TL | Fine-Tuning (baseline) | $45.25_{\pm0.12}$ | $59.68_{\pm0.21}$ | $70.12_{\pm0.29}$ | $78.01_{\pm0.16}$ |
| | | $L^2$-SP (Li et al., 2018) | $45.08_{\pm0.19}$ | $57.78_{\pm0.24}$ | $69.47_{\pm0.29}$ | $78.44_{\pm0.17}$ |
| | | DELTA (Li et al., 2019) | $46.83_{\pm0.21}$ | $60.37_{\pm0.25}$ | $71.38_{\pm0.20}$ | $78.63_{\pm0.18}$ |
| | | BSS (Chen et al., 2019) | $47.74_{\pm0.23}$ | $63.38_{\pm0.29}$ | $72.56_{\pm0.17}$ | $78.85_{\pm0.31}$ |
| | | Co-Tuning (You et al., 2020) | $52.58_{\pm0.53}$ | $66.47_{\pm0.17}$ | $74.64_{\pm0.36}$ | $81.24_{\pm0.14}$ |
| | SSL | Π-model (Laine & Aila, 2017) | $45.20_{\pm0.23}$ | $56.20_{\pm0.29}$ | $64.07_{\pm0.32}$ | – |
| | | Pseudo-Labeling (Lee, 2013) | $45.33_{\pm0.24}$ | $62.02_{\pm0.31}$ | $72.30_{\pm0.29}$ | – |
| | | Mean Teacher (Tarvainen & Valpola, 2017) | $53.26_{\pm0.19}$ | $66.66_{\pm0.20}$ | $74.37_{\pm0.30}$ | – |
| | | UDA (Xie et al., 2020) | $46.90_{\pm0.31}$ | $61.16_{\pm0.35}$ | $71.86_{\pm0.43}$ | – |
| | | FixMatch (Sohn et al., 2020) | $44.06_{\pm0.23}$ | $63.54_{\pm0.18}$ | $75.96_{\pm0.29}$ | – |
| | | SimCLRv2 (Chen et al., 2020b) | $45.74_{\pm0.15}$ | $62.70_{\pm0.24}$ | $71.01_{\pm0.34}$ | – |
| | Combine | Co-Tuning + Pseudo-Labeling | $54.11_{\pm0.24}$ | $68.07_{\pm0.32}$ | $75.94_{\pm0.34}$ | – |
| | | Co-Tuning + Mean Teacher | $57.92_{\pm0.18}$ | $67.98_{\pm0.25}$ | $72.82_{\pm0.29}$ | – |
| | | Co-Tuning + FixMatch | $46.81_{\pm0.21}$ | $58.88_{\pm0.23}$ | $73.07_{\pm0.29}$ | – |
| | | **Self-Tuning (ours)** | $\mathbf{64.17_{\pm0.47}}$ | $\mathbf{75.13_{\pm0.35}}$ | $\mathbf{80.22_{\pm0.36}}$ | $\mathbf{83.95_{\pm0.18}}$ |

# Experiments

| Dataset | Type | Method | Label Proportion | | | |
|---|---|---|---|---|---|---|
| | | | 15% | 30% | 50% | 100% |
| *Stanford Cars* | TL | Fine-Tuning (baseline) | $36.77_{\pm0.12}$ | $60.63_{\pm0.18}$ | $75.10_{\pm0.21}$ | $87.20_{\pm0.19}$ |
| | | $L^2$-SP (Li et al., 2018) | $36.10_{\pm0.30}$ | $60.30_{\pm0.28}$ | $75.48_{\pm0.22}$ | $86.58_{\pm0.26}$ |
| | | DELTA (Li et al., 2019) | $39.37_{\pm0.34}$ | $63.28_{\pm0.27}$ | $76.53_{\pm0.24}$ | $86.32_{\pm0.20}$ |
| | | BSS (Chen et al., 2019) | $40.57_{\pm0.12}$ | $64.13_{\pm0.18}$ | $76.78_{\pm0.21}$ | $87.63_{\pm0.27}$ |
| | | Co-Tuning (You et al., 2020) | $46.02_{\pm0.18}$ | $69.09_{\pm0.10}$ | $80.66_{\pm0.25}$ | $89.53_{\pm0.09}$ |
| | SSL | Π-model (Laine & Aila, 2017) | $45.19_{\pm0.21}$ | $57.29_{\pm0.26}$ | $64.18_{\pm0.29}$ | – |
| | | Pseudo-Labeling (Lee, 2013) | $40.93_{\pm0.23}$ | $67.02_{\pm0.19}$ | $78.71_{\pm0.30}$ | – |
| | | Mean Teacher (Tarvainen & Valpola, 2017) | $54.28_{\pm0.14}$ | $66.02_{\pm0.21}$ | $74.24_{\pm0.23}$ | – |
| | | UDA (Xie et al., 2020) | $39.90_{\pm0.43}$ | $64.16_{\pm0.40}$ | $71.86_{\pm0.56}$ | – |
| | | FixMatch (Sohn et al., 2020) | $49.86_{\pm0.27}$ | $77.54_{\pm0.29}$ | $84.78_{\pm0.33}$ | – |
| | | SimCLRv2 (Chen et al., 2020b) | $45.74_{\pm0.16}$ | $61.70_{\pm0.18}$ | $77.49_{\pm0.24}$ | – |
| | Combine | Co-Tuning + Pseudo-Labeling | $50.16_{\pm0.23}$ | $73.76_{\pm0.26}$ | $83.33_{\pm0.34}$ | – |
| | | Co-Tuning + Mean Teacher | $52.98_{\pm0.19}$ | $71.42_{\pm0.24}$ | $75.38_{\pm0.29}$ | – |
| | | Co-Tuning + FixMatch | $42.34_{\pm0.19}$ | $73.24_{\pm0.25}$ | $83.13_{\pm0.34}$ | – |
| | | **Self-Tuning (ours)** | $\mathbf{72.50}_{\pm0.45}$ | $\mathbf{83.58}_{\pm0.28}$ | $\mathbf{88.11}_{\pm0.29}$ | $\mathbf{90.67}_{\pm0.23}$ |

# Experiments

| Dataset | Type | Method | Label Proportion | | | |
|---|---|---|---|---|---|---|
| | | | 15% | 30% | 50% | 100% |
| *FGVC Aircraft* | TL | Fine-tuning (baseline) | $39.57_{\pm0.20}$ | $57.46_{\pm0.12}$ | $67.93_{\pm0.28}$ | $81.13_{\pm0.21}$ |
| | | $L^2$-SP (Li et al., 2018) | $39.27_{\pm0.24}$ | $57.12_{\pm0.27}$ | $67.46_{\pm0.26}$ | $80.98_{\pm0.29}$ |
| | | DELTA (Li et al., 2019) | $42.16_{\pm0.21}$ | $58.60_{\pm0.29}$ | $68.51_{\pm0.25}$ | $80.44_{\pm0.20}$ |
| | | BSS (Chen et al., 2019) | $40.41_{\pm0.12}$ | $59.23_{\pm0.31}$ | $69.19_{\pm0.13}$ | $81.48_{\pm0.18}$ |
| | | Co-Tuning (You et al., 2020) | $44.09_{\pm0.67}$ | $61.65_{\pm0.32}$ | $72.73_{\pm0.08}$ | $83.87_{\pm0.09}$ |
| | SSL | $\Pi$-model (Laine & Aila, 2017) | $37.32_{\pm0.25}$ | $58.49_{\pm0.26}$ | $65.63_{\pm0.36}$ | – |
| | | Pseudo-Labeling (Lee, 2013) | $46.83_{\pm0.30}$ | $62.77_{\pm0.31}$ | $73.21_{\pm0.39}$ | – |
| | | Mean Teacher (Tarvainen & Valpola, 2017) | $51.59_{\pm0.23}$ | $71.62_{\pm0.29}$ | $80.31_{\pm0.32}$ | – |
| | | UDA (Xie et al., 2020) | $43.96_{\pm0.45}$ | $64.17_{\pm0.49}$ | $67.42_{\pm0.53}$ | – |
| | | FixMatch (Sohn et al., 2020) | $55.53_{\pm0.26}$ | $71.35_{\pm0.35}$ | $78.34_{\pm0.43}$ | – |
| | | SimCLRv2 (Chen et al., 2020b) | $40.78_{\pm0.21}$ | $59.03_{\pm0.29}$ | $68.54_{\pm0.30}$ | – |
| | Combine | Co-Tuning + Pseudo-Labeling | $49.15_{\pm0.32}$ | $65.62_{\pm0.34}$ | $74.57_{\pm0.40}$ | – |
| | | Co-Tuning + Mean Teacher | $51.46_{\pm0.25}$ | $64.30_{\pm0.28}$ | $70.85_{\pm0.35}$ | – |
| | | Co-Tuning + FixMatch | $53.74_{\pm0.23}$ | $69.91_{\pm0.26}$ | $80.02_{\pm0.32}$ | – |
| | | **Self-Tuning (ours)** | $\mathbf{64.11}_{\pm0.32}$ | $\mathbf{76.03}_{\pm0.25}$ | $\mathbf{81.22}_{\pm0.29}$ | $\mathbf{84.28}_{\pm0.14}$ |

# Experiments (Unsupervised Pretrained Model)

Table 4. Classification accuracy (%) ↑ with a typical unsupervised pre-trained model MoCov2 on *CUB-200-2011*.

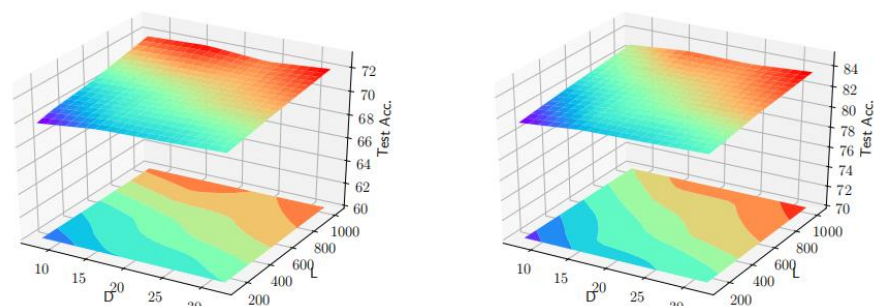| Type | Method | 800 labels | 5k labels |
|------|--------|-----------|-----------|
| TL | Fine-Tuning (baseline) | 20.04 | 71.50 |
| | Co-Tuning | 20.99 | 71.61 |
| SSL | Mean Teacher | 28.13 | 71.26 |
| | FixMatch | 21.18 | 71.28 |
| Combine | Co-Tuning + Mean Teacher | 28.43 | 72.21 |
| | Co-Tuning + FixMatch | 21.08 | 71.40 |
| | **Self-Tuning (ours)** | **36.80** | **74.56** |

# Experiments (Ablation studies)

Table 5. Ablation studies of Self-Tuning on *Stanford Cars*.

| Perspective | Method | 15% | 30% |
|---|---|---|---|
| Loss Function | w/ CE loss | 40.93 | 67.02 |
| | w/ CL loss | 46.29 | 68.82 |
| | w/ PGC loss | **72.50** | **83.58** |
| Info. Exploration | w/o $\widehat{L}_{\mathrm{PGC}}$ | 58.82 | 81.71 |
| | w/o $L_{\mathrm{PGC}}$ | 58.85 | 77.52 |
| | separate queue | 70.43 | 80.78 |
| | unified exploration | **72.50** | **83.58** |

# Experiments (Sensitivity Analysis & Others)



(a) Acc on *Car* with 15% labels (b) Acc on *Car* with 30% labels

*Figure 6.* Sensitivity analysis for embedded size $L$ of the projector and queue size $D$ of each class on *Stanford Cars*. (Warmer colors indicate higher values)



(a) Training Process on *CUB30* (b) $\mathrm{Acc}_{test} - \mathrm{Acc}_{pseudo\_labels}$

*Figure 7.* Comparisons between Self-Tuning with FixMatch on pseudo label accuracy and test accuracy.