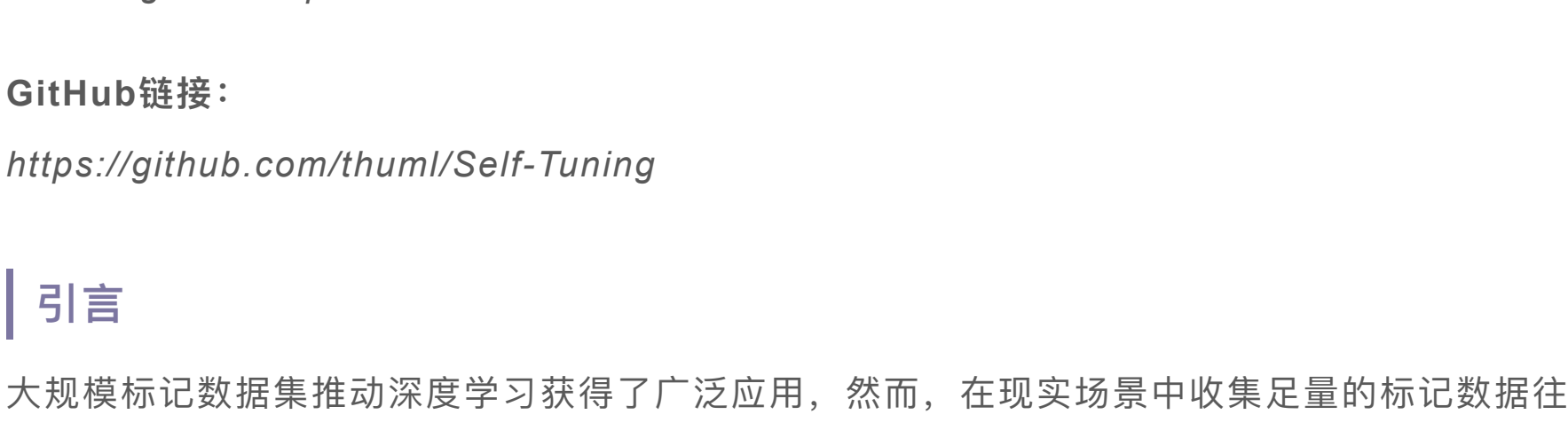


本文介绍ICML2021的中稿论文：Self-Tuning for Data-Efficient Deep Learning，就“如何减少对标记数据的需求”这一重要问题给出了我们的思考。



论文标题：

Self-Tuning for Data-Efficient Deep Learning

论文链接：

http://se.thss.tsinghua.edu.cn/~mlong/doc/Self-Tuning-for-Data-Efficient-Deep-Learning-icml21.pdf

GitHub链接：

https://github.com/thuml/Self-Tuning

## 引言

大规模标记数据集推动深度学习获得了广泛应用，然而，在现实场景中收集足量的标记数据往往耗时耗力。为了减少对标记数据的需求，半监督学习和迁移学习的研究者们从两个不同的视角给出了自己的思考：半监督学习(Semi-supervised Learning, SSL)侧重于同时探索标记数据和无标记数据，通过挖掘无标记数据的内在结构增强模型的泛化能力，而迁移学习(Transfer Learning, TL)旨在将预训练模型迁移到目标数据中，也就是我们耳熟能详的预训练-微调范式。

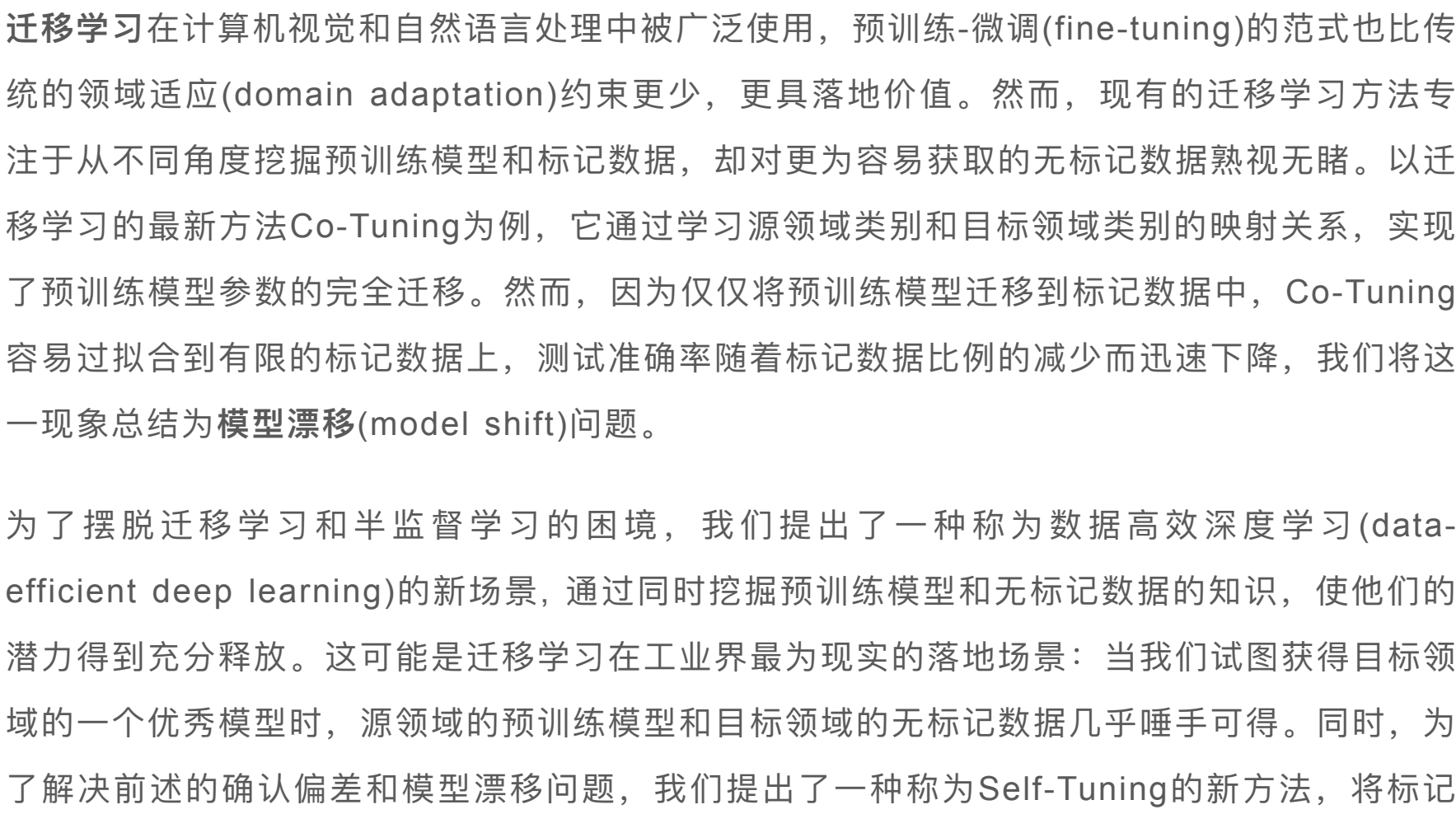


Figure 4. Test accuracy of a state-of-the-art SSL method and a TL method on various class numbers or label ratios respectively.

半监督学习的最新进展，例如UDA，FixMatch等方法，证明了自训练(Self-Tuning)的巨大潜力。通过蒸馏广样本生成伪标记(pseudo-label)，FixMatch就可以在Cifar10、SVHN、STL-10数据集上取得了令人耳目一新的效果。然而，细心的读者会发现，上述数据集都是类别数较少的简单数据集(都是10类)，当类别数增加到100时，FixMatch这种从头开始训练(train from scratch)的自训练方法的表现就强人意了。进一步地，我们在CUB200上将类别数从10逐渐增加到200时，发现FixMatch的准确率随着伪标签的准确率的下降而快速下降。这说明，随着类别数的增加，伪标签的质量逐渐下降，而自训练的模型也被错误的伪标签所误导，从而难以在测试数据集上取得可观的效果。这一现象，被前人总结为自训练的确认偏差(confirmiation bias)问题，说明Self-training虽然是很优秀，但确利刃双刃。

迁移学习在计算机视觉和自然语言处理中被广泛使用，预训练-微调(fine-tuning)的范式也比传统的领域适应(domain adaptation)约束更少，更具落地价值。然而，现有的迁移学习方法专注于从不同角度挖掘预训练模型和标记数据，却对更为容易获取的无标记数据熟视无睹。以迁移学习的最新方法Co-Tuning为例，它通过学习源领域类别和目标领域类别的映射关系，实现了预训练模型参数的完全迁移。然而，因为仅仅将预训练模型迁移到标记数据中，Co-Tuning容易过拟合到有限的标记数据上，测试准确率随着标记数据比例的减少而迅速下降，我们将这一现象总结为模型漂移(model shift)问题。

为了摆脱迁移学习和半监督学习的困境，我们提出了一种称为数据高效深度学习(data-efficient deep learning)的新场景，通过同时挖掘预训练模型和无标记数据的知识，使他们的潜力得到充分释放。这可能是迁移学习在工业界最为现实的落地场景：当我们试图获得目标领域的一个优秀模型时，源领域的预训练模型和目标领域的无标记数据几乎唾手可得。同时，为了解决前述的确认偏差和模型漂移问题，我们提出了一种称为Self-Tuning的新方法，将标记数据和无标记数据的探索与预训练模型的迁移融为一体，以及一种通用的伪标签组对比机制(Pseudo Group Contrast)，从而减轻对伪标签的依赖，提高模型的鲁棒性。在多个标准数据集的实验表明，Self-Tuning远远优于半监督学习和迁移学习的同类方法。例如，在标签比例为15%的Stanford-Cars数据集上，Self-Tuning的测试精度比fine-tuning几乎提高了一倍。

## 如何解决确认偏差问题？

为了找出自训练的确认偏差(confirmiation bias)问题的根源，我们首先分析了伪标签(pseudo-label)广泛采用的交叉熵损失函数(Cross-Entropy, CE):

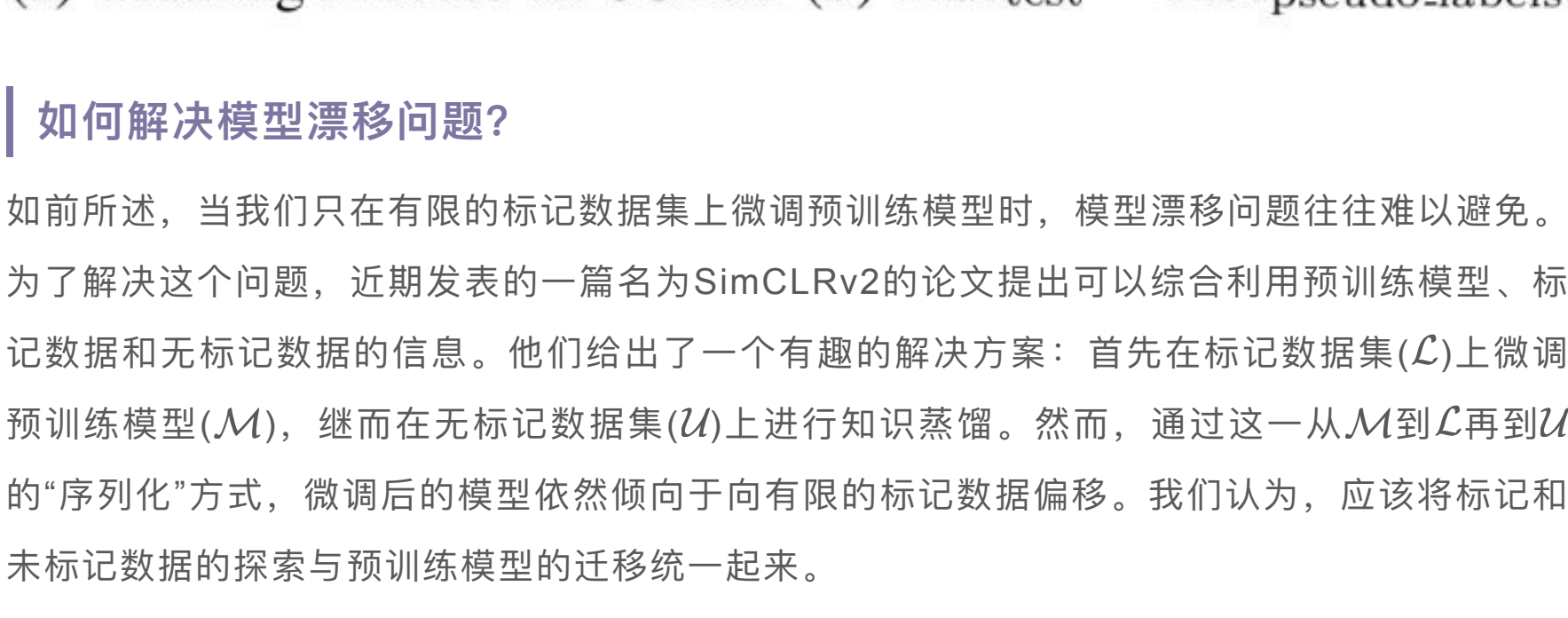
$$\hat{L}_{CE} = - \sum_{c=1}^C \mathbf{1}(\hat{y}_i = c) \mathbf{1}(z_i > t) \log \mathbf{p}_i^c$$

其中， $\hat{y}_i$ 是输入 $x_i$ 生成的伪标签，而 $z_i = \max_c(\mathbf{p}_i^c)$ 是模型对于样本 $x_i$ 。通常地，大多数自训练方法都会针对confidence做一个阈值过滤，只有大于阈值 $t$  (比如FixMatch中设置了0.95的阈值)的样本的预测标签才会被假设为合格的伪标签加入模型训练。然而，如图2所示，由于交叉熵损失函数专注于学习不同类别的分类面，如果某些伪标签存在错误，通过交叉熵损失函数训练的模型就会轻易地被错误的伪标签所误导。

为了解决交叉熵损失函数的类别鉴别(class discrimination)特性对自训练带来的挑战，最近取得突破进展的基于样本鉴别(sample discrimination)思想的对比学习损失函数吸引了我们的注意。给定由输入 $x_i$ 生成的查询样本 $q_i$ ， $x_i$ 在不同数据源下生成的副本 $k_j$ ，以及 $D$ 个不同输入生成的负样本 $\{k_1, k_2, \dots, k_D\}$ ，则通过内积度量相似性的对比学习(Contrastive Learning, CL)损失函数可以定义为

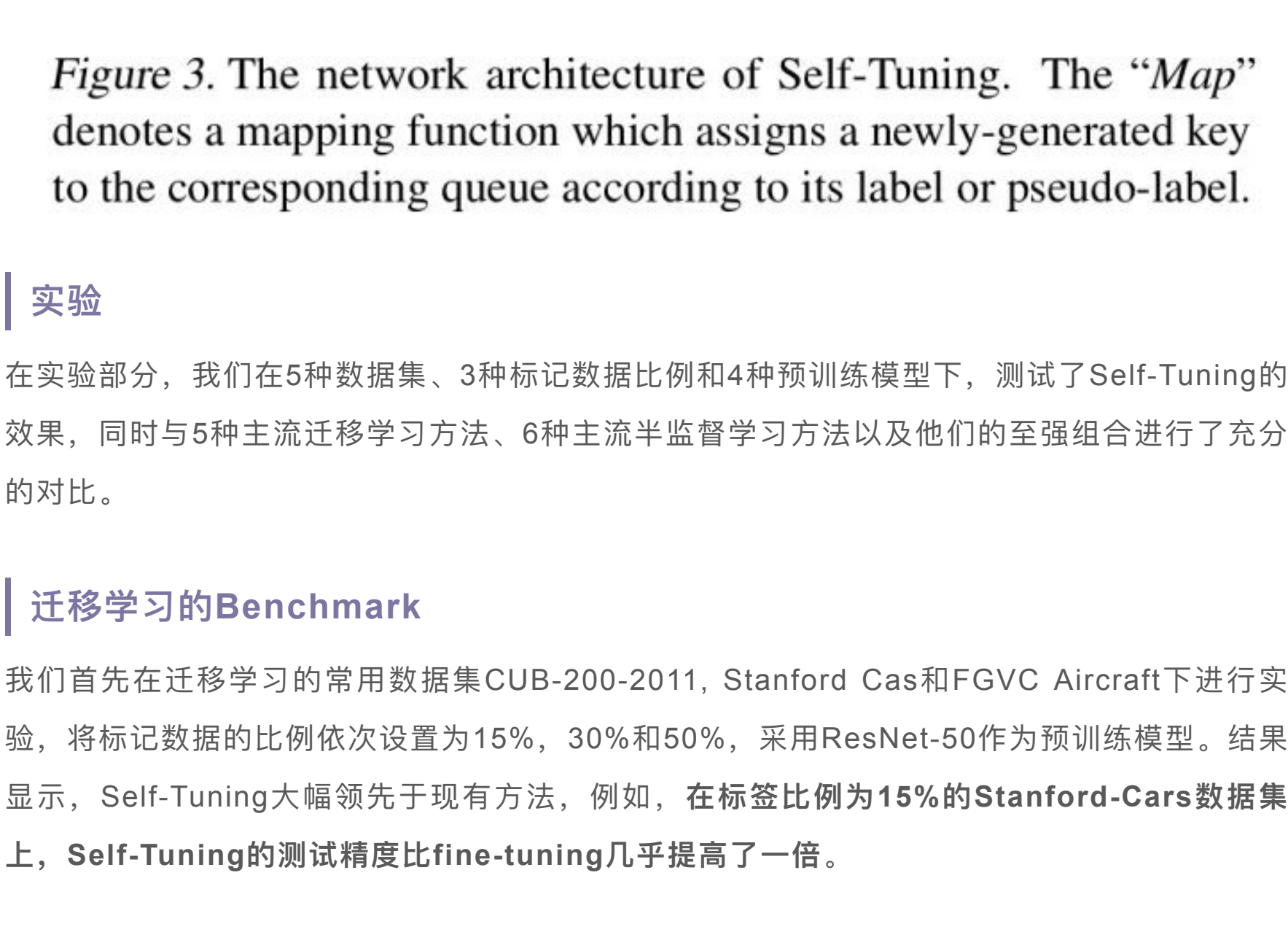
$$L_{CL} = - \log \frac{\exp(\mathbf{q} \cdot \mathbf{k}_0 / \tau)}{\exp(\mathbf{q} \cdot \mathbf{k}_0 / \tau) + \sum_{d=1}^D \exp(\mathbf{q} \cdot \mathbf{k}_d / \tau)}$$

可以看出，对比学习旨在最大化同一样本在两个不同数据源下的表征相似性，而最小化不同样本之间的表征相似性，从而实现样本鉴别。挖掘数据中隐藏的流形结构。这种设计与伪标签无关，天然地不受错误的伪标签的影响。然而，标准的对比学习损失函数未能将标记和伪标签嵌入到模型训练中，从而使有用的鉴别信息束之高阁。



为了解决这一挑战，我们提出了一种通用的伪标签组对比机制(Pseudo Group Contrast, PGC)。对于任何一个查询样本 $q_i$ ，它的伪标签用 $\hat{y}$ 表示。PGC将具有相同伪标签( $\hat{y}$ )的样本都视为正样本，而具有不同伪标签( $\{1, 2, \dots, C\} \setminus \hat{y}$ )的样本则组成了负样本，从而最大化查询样本与具有相同伪标签的正样本的表征相似性，实现伪标签组对比。

$$\hat{L}_{PGC} = - \frac{1}{D+1} \sum_{d=0}^D \log \frac{\exp(\mathbf{q} \cdot \mathbf{k}_d^+ / \tau)}{\text{Pos} + \text{Neg}}$$
$$\text{Pos} = \exp(\mathbf{q} \cdot \mathbf{k}_0^+ / \tau) + \sum_{j=1}^D \exp(\mathbf{q} \cdot \mathbf{k}_j^+ / \tau)$$
$$\text{Neg} = \sum_{c=1}^{\{1,2,\dots,C\} \setminus \hat{y}} \sum_{j=1}^D \exp(\mathbf{q} \cdot \mathbf{k}_j^c / \tau),$$



## 实验

在实验部分，我们在6种数据集、3种标记数据比例和4种预训练模型下，测试了Self-Tuning的效果，同时与5种主流迁移学习方法、6种主流半监督学习方法以及他们的强强组合进行了充分的对比。

## 迁移学习的Benchmark

我们首先在迁移学习的常用数据集CUB-200-2011、Stanford Cars和FGVC Aircraft下进行实验，将标记数据的比例依次设置为15%、30%和50%，采用ResNet-50作为预训练模型。结果显示，Self-Tuning大幅领先于现有方法，例如，在标签比例为15%的Stanford-Cars数据集上，Self-Tuning的测试精度比fine-tuning几乎提高了一倍。

Dataset	Type	Method	Label Proportion			
			15%	30%	50%	100%
CUB-200-2011	TL	Fine-Tuning (baseline)	45.25±0.19	59.68±0.11	70.12±0.19	78.01±0.16
		L2-SP (Li et al., 2018)	45.08±0.19	57.78±0.24	69.47±0.29	78.44±0.17
		DELTA (Li et al., 2019)	46.83±0.19	60.37±0.25	71.38±0.29	78.83±0.18
		BSS (Chen et al., 2019)	47.74±0.19	63.38±0.29	72.56±0.17	78.85±0.19
		Co-Tuning (You et al., 2020)	52.58±0.18	66.47±0.17	74.64±0.16	81.24±0.14
	SSL	IL-model (Laine & Aila, 2017)	45.20±0.21	56.20±0.29	64.07±0.32	—
		Pseudo-Labeling (Lee, 2013)	45.33±0.18	62.02±0.31	72.30±0.29	—
		Mean Teacher (Tarvainen & Valpola, 2017)	53.36±0.19	66.06±0.29	74.37±0.18	—
		UDA (Xie et al., 2020)	46.90±0.19	61.16±0.35	71.86±0.47	—
		FixMatch (Sohn et al., 2020)	44.06±0.19	63.54±0.18	75.96±0.29	—
Stanford Cars	TL	IL-model (Laine & Aila, 2017)	45.20±0.21	56.20±0.29	64.07±0.32	—
		Pseudo-Labeling (Lee, 2013)	45.33±0.18	62.02±0.31	72.30±0.29	—
		Mean Teacher (Tarvainen & Valpola, 2017)	53.36±0.19	66.06±0.29	74.37±0.18	—
		UDA (Xie et al., 2020)	46.90±0.19	61.16±0.35	71.86±0.47	—
		FixMatch (Sohn et al., 2020)	44.06±0.19	63.54±0.18	75.96±0.29	—
	SSL	IL-model (Laine & Aila, 2017)	45.20±0.21	56.20±0.29	64.07±0.32	—
		Pseudo-Labeling (Lee, 2013)	45.33±0.18	62.02±0.31	72.30±0.29	—
		Mean Teacher (Tarvainen & Valpola, 2017)	53.36±0.19	66.06±0.29	74.37±0.18	—
		UDA (Xie et al., 2020)	46.90±0.19	61.16±0.35	71.86±0.47	—
		FixMatch (Sohn et al., 2020)	44.06±0.19	63.54±0.18	75.96±0.29	—
FGVC Aircraft	TL	IL-model (Laine & Aila, 2017)	45.20±0.21	56.20±0.29	64.07±0.32	—
		Pseudo-Labeling (Lee, 2013)	45.33±0.18	62.02±0.31	72.30±0.29	—
		Mean Teacher (Tarvainen & Valpola, 2017)	53.36±0.19	66.06±0.29	74.37±0.18	—
		UDA (Xie et al., 2020)	46.90±0.19	61.16±0.35	71.86±0.47	—
		FixMatch (Sohn et al., 2020)	44.06±0.19	63.54±0.18	75.96±0.29	—
	SSL	IL-model (Laine & Aila, 2017)	45.20±0.21	56.20±0.29	64.07±0.32	—
		Pseudo-Labeling (Lee, 2013)	45.33±0.18	62.02±0.31	72.30±0.29	—
		Mean Teacher (Tarvainen & Valpola, 2017)	53.36±0.19	66.06±0.29	74.37±0.18	—
		UDA (Xie et al., 2020)	46.90±0.19	61.16±0.35	71.86±0.47	—
		FixMatch (Sohn et al., 2020)	44.06±0.19	63.54±0.18	75.96±0.29	—

Method	Network	2.5k	10k
IL-Model	WRN-28-8	57.25	37.88
Pseudo-Labeling		57.38	36.21
Mean Teacher		53.91	35.83
MixMatch		39.94	28.31
UDA		33.13	24.50
ReMixMatch	#Para: 11.76M	27.43	23.03
FixMatch		28.64	23.18

FixMatch	EfficientNet-B2	29.99	21.69
Fine-Tuning		31.69	21.74
Co-Tuning		30.94	22.22
Self-Tuning		24.16	17.57

Method	Network	2.5k	10k
IL-Model	WRN-28-8	57.25	37.88
Pseudo-Labeling		57.38	36.21
Mean Teacher		53.91	35.83
MixMatch		39.94	28.31
UDA		33.13	24.50
ReMixMatch	#Para: 11.76M	27.43	23.03
FixMatch		28.64	23.18

Method	Network	2.5k	10k
IL-Model	WRN-28-8	57.25	37.88
Pseudo-Labeling		57.38	36.21
Mean Teacher		53.91	35.83
MixMatch		39.94	28.31
UDA		33.13	24.50
ReMixMatch	#Para: 11.76M	27.43	23.03
FixMatch		28.64	23.18

Method	Network	2.5k	10k
IL-Model	WRN-28-8	57.25	37.88
Pseudo-Labeling		57.38	36.21
Mean Teacher		53.91	35.83
MixMatch		39.94	28.31
UDA		33.13	24.50
ReMixMatch	#Para: 11.76M	27.43	23.03
FixMatch		28.64	23.18

Method	Network	2.5k	10k
IL-Model	WRN-28-8	57.25	37.88
Pseudo-Labeling		57.38	36.21
Mean Teacher		53.91	35.83
MixMatch		39.94	28.31
UDA		33.13	24.50
ReMixMatch	#Para: 11.76M	27.43	23.03
FixMatch		28.64	23.18

Method	Network	2.5k	10k
IL-Model	WRN-28-8	57.25	37.88
Pseudo-Labeling		57.38	36.21
Mean Teacher		53.91	35.83
MixMatch		39.94	28.31
UDA		33.13	24.50
ReMixMatch	#Para: 11.76M	27.43	23.03
FixMatch		28.64	23.18

Method	Network	2.5k	10k
IL-Model	WRN-28-8	57.25	37.88
Pseudo-Labeling		57.38	36.21
Mean Teacher		53.91	35.83
MixMatch		39.94	28.31
UDA		33.13	24.50
ReMixMatch	#Para: 11.76M	27.43	23.03
FixMatch		28.64	23.18

Method	Network	2.5k	10k
IL-Model	WRN-28-8	57.25	37.88
Pseudo-Labeling		57.38	36.21
Mean Teacher		53.91	35.83
MixMatch		39.94	28.31
UDA		33.13	24.50
ReMixMatch	#Para: 11.76M	27.43	23.03
FixMatch		28.64	23.18

Method	Network	2.5k	10k
IL-Model	WRN-28-8	57.25	37.88
Pseudo-Labeling		57.38	36.21
Mean Teacher		53.91	35.83
MixMatch		39.94	28.31
UDA		33.13	24.50
ReMixMatch	#Para: 11.76M	27.43	23.03
FixMatch		28.64	23.18