

具身智能-06

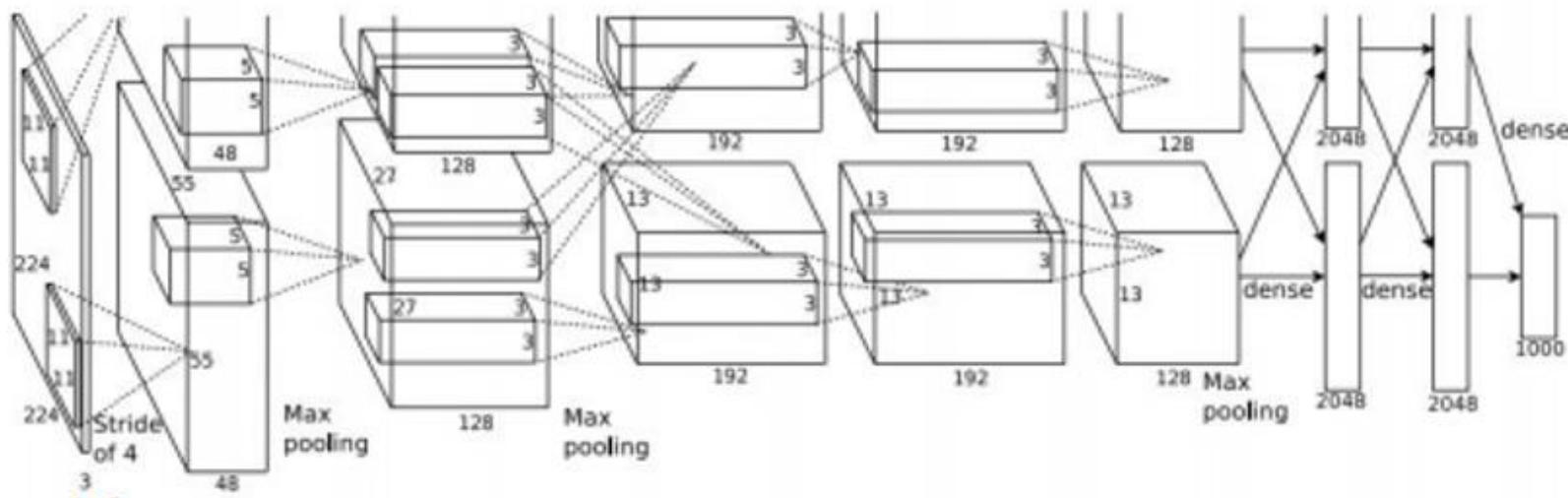
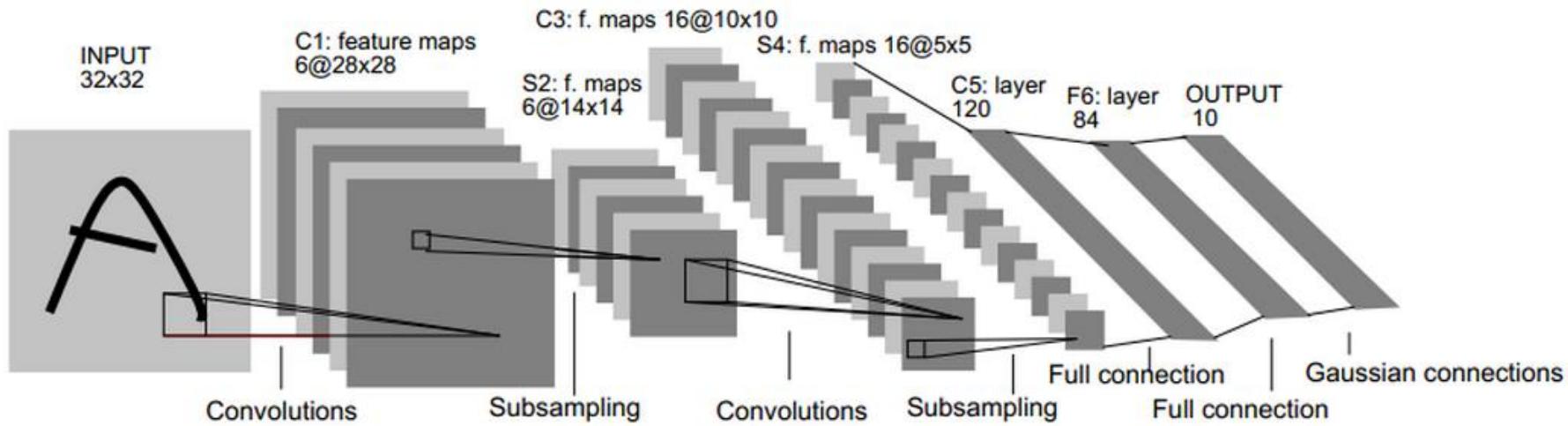
刘华平

2025年3月26日

课次	上课内容
1	绪论
2	深度学习
3	强化学习1
4	强化学习2
5	作业准备
6	自监督与持续学习
7	开题
8	形态智能
9	视觉导航：VLN
10	主动感知：VSN, EQA
11	五一放假
12	具身学习
13	多体智能
14	面向具身智能的AIGC
15-16	成果准备与展示

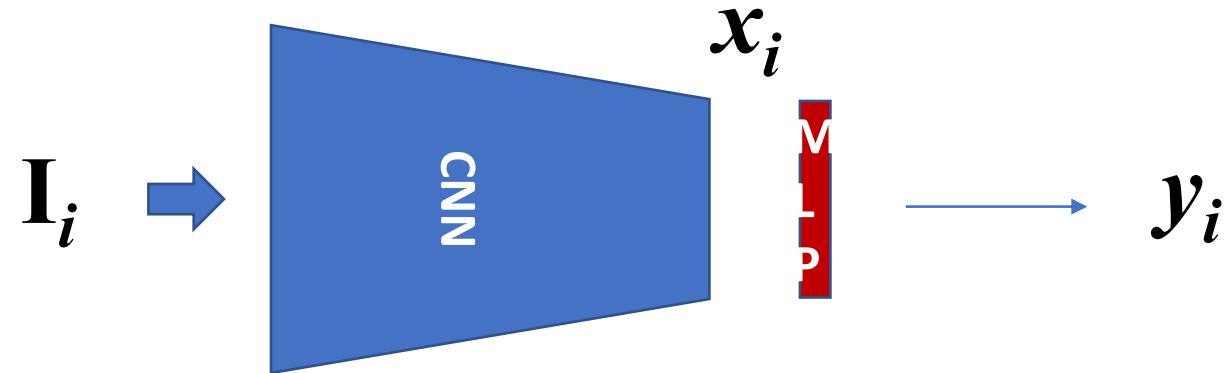
背景

➤ 深度学习



背景

➤ 深度学习

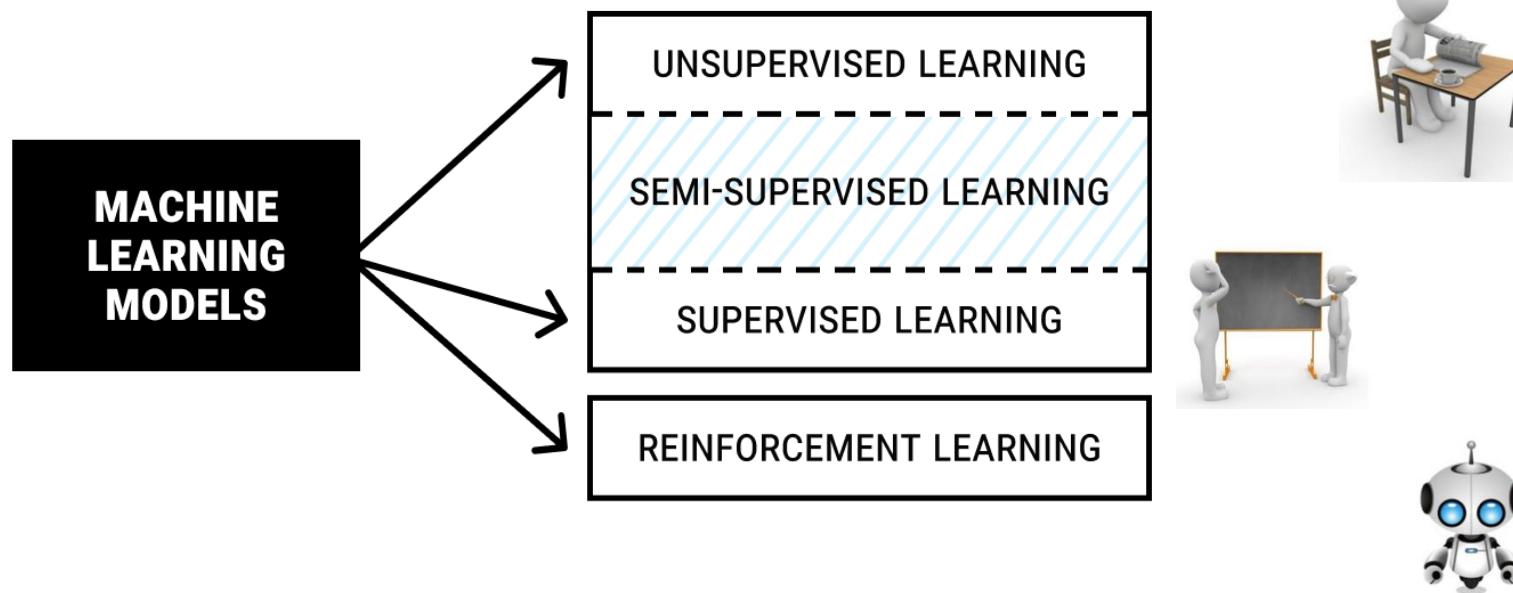


Loss $\{y_i, \textcolor{red}{f}(\textcolor{blue}{x}_i)\}$

Loss $\{y_i, \varphi_{\theta}(I_i)\}$

背景

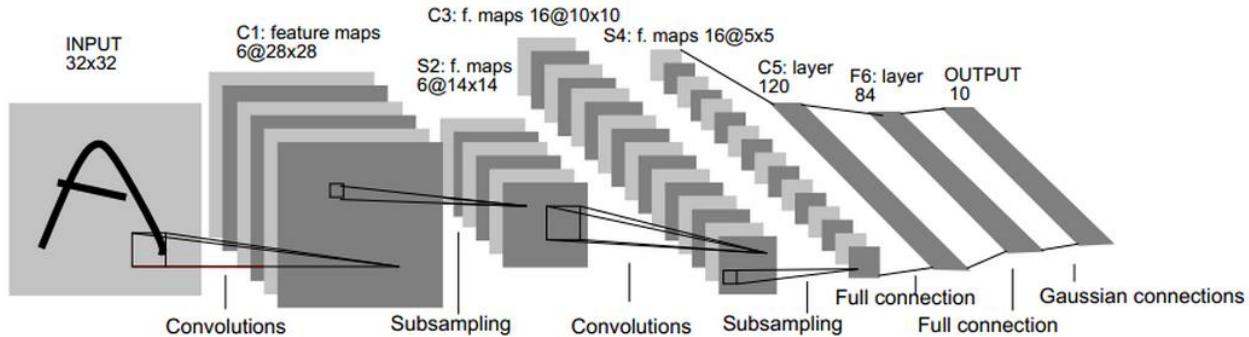
➤ 机器学习



背景

➤ 机器学习

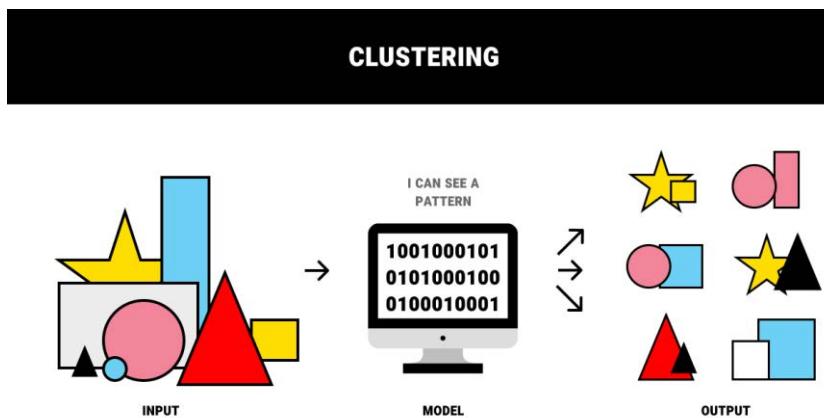
监督式学习



$\{f(x_i), y_i\}$: 特征与分类

无监督学习

自监督学习



$\{f(x_i)\}$: 聚类

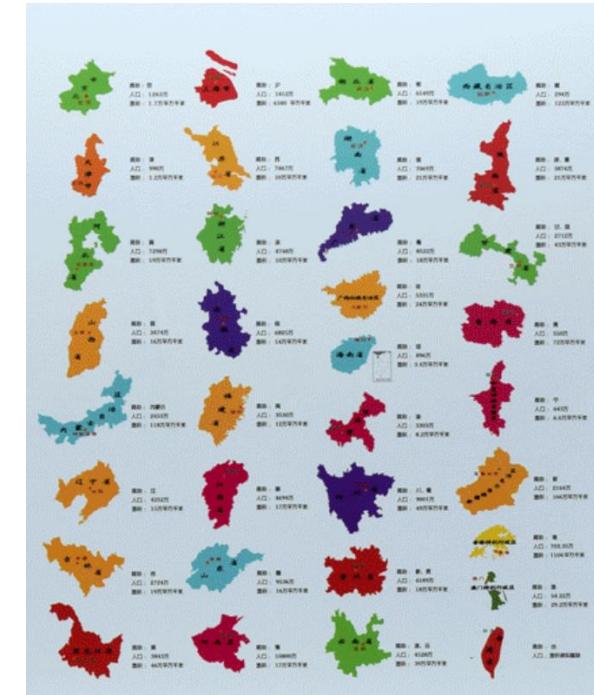


$\{\textcolor{blue}{f}(x_i), \textcolor{red}{y}_i\}$: 特征

背景

- 自监督学习
- 持续学习

自监督学习



自监督学习

➤ 定义

动机：

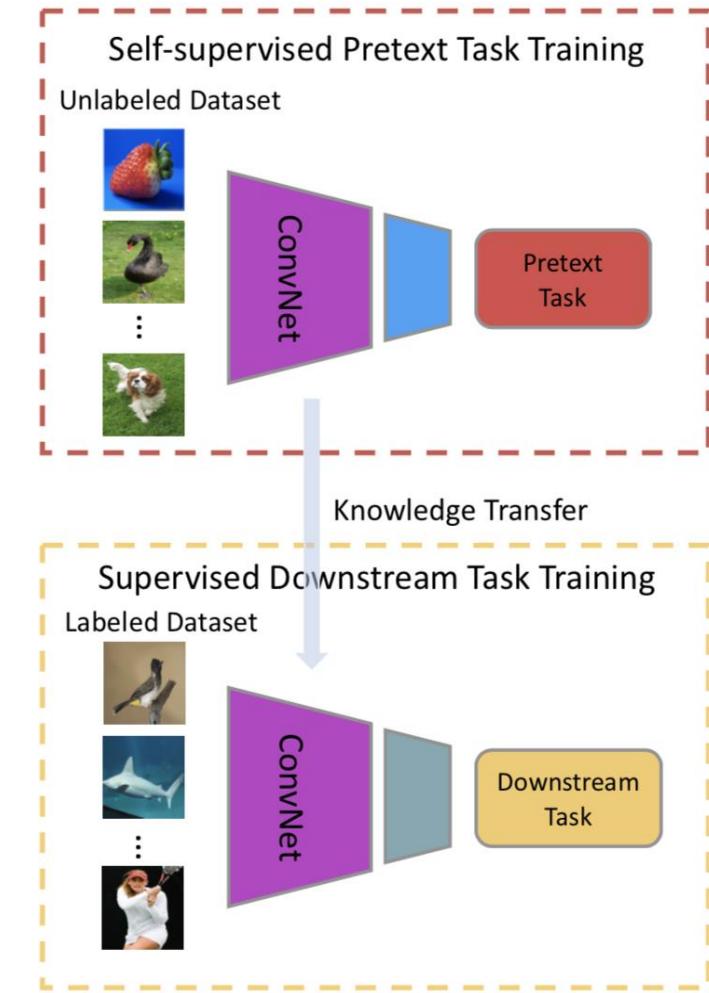
- 监督学习需要带人工标注标签的数据集，标注标签是困难的一步，需要大量的时间和金钱成本。无监督学习/自监督学习的存在便是为了解决这一问题。由于没有人工标注的标签的存在，大部分自监督学习的主要任务在于构建**可自动生成的监督信号**，从而设计出借口任务（pretext task）。

借口任务的特点

- **监督信号可以自动生成**，所以避免了人工标注大量语义标签的过程。借口任务最好是需要一定高级的抽象语义知识才可以解决的，这样学到的网络特征才能更好适应下游任务。

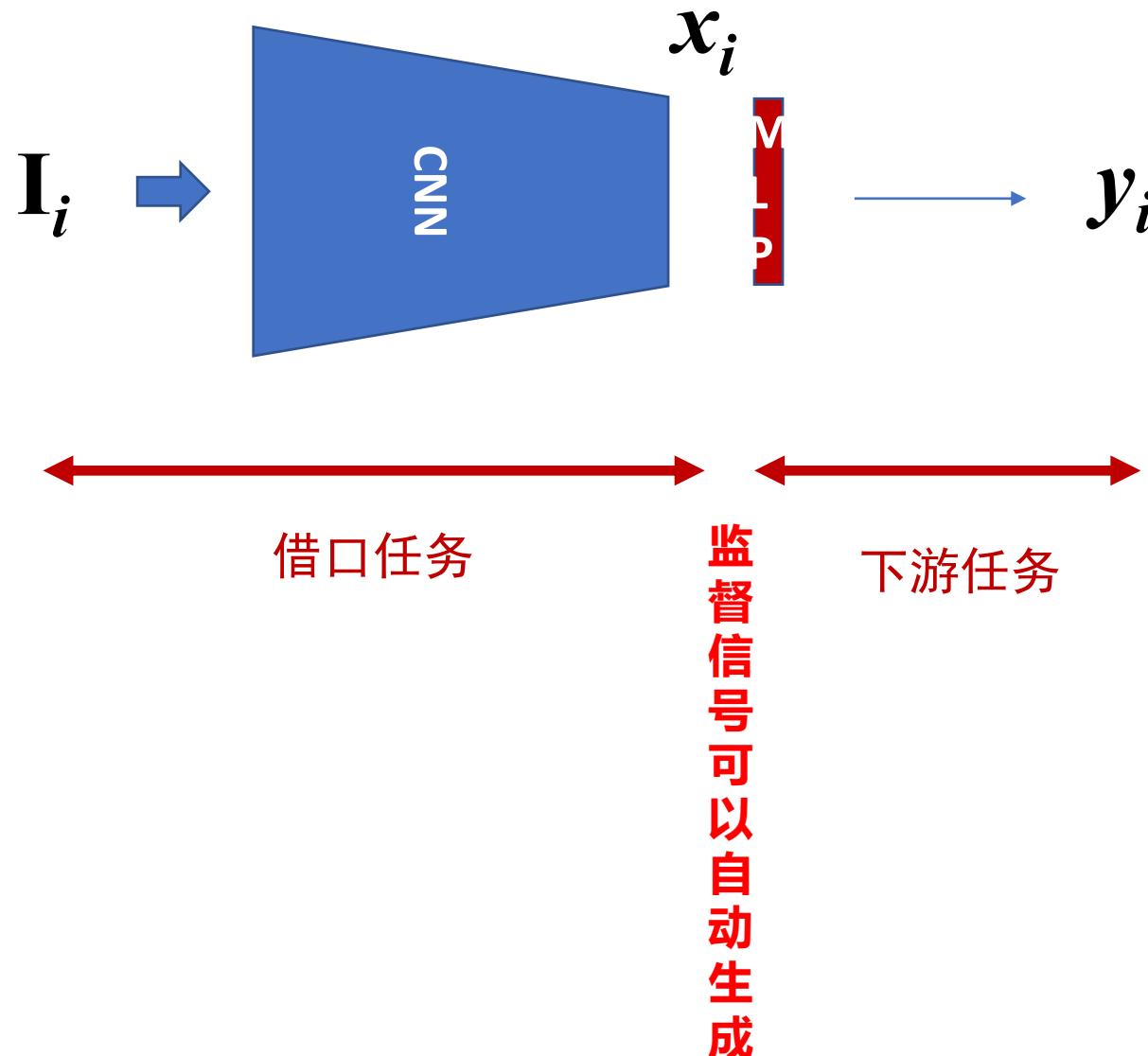
下游任务(Downstream Task)

- 下游任务是指人们真正亟待解决但标注稀缺的问题，例如小数据集的分类问题、检测、分割等等。借口任务学习完成后，通常将学习到的网络作为下游任务的初始化（迁移学习）或者固定参数直接将其作为特征提取器。



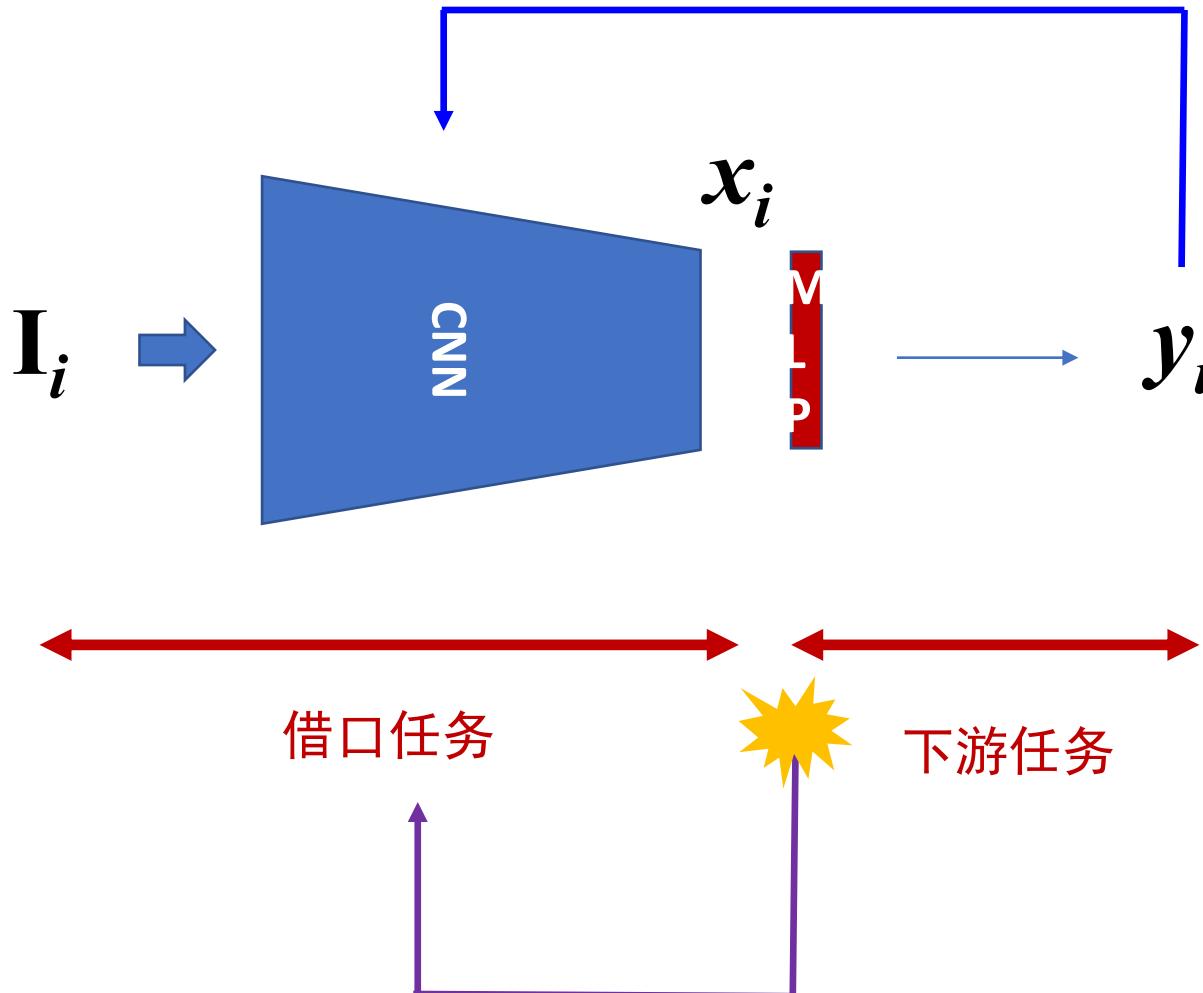
自监督学习

➤ 定义



自监督学习

➤ 定义



自监督学习

➤ 下游任务

下游任务 (benchmark) :

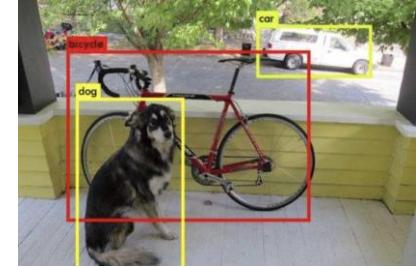
在图像方面的应用包括**图像分类、目标检测、语义分割**。在下游任务上的表现通常是衡量自监督学习方法好坏的依据。常用的数据集包括Pascal VOC2007, ImageNet和Places 205场景数据集。

Pascal VOC2007: 训练集包括5011张图片，测试集4952张图片。共包含20个类别。包括分类、检测、分割三个任务。通常**采用finetune方式来判断自监督模型的好坏**。

ImageNet: 包括约130万张物体图像，共1000类。

Places 205: 包括约250万张场景图像，共205类。

ImageNet和Places 205为大型数据集，通常采用**固定自监督模型参数，并增加一个线性层的方法来衡量自监督模型的性能**。



目标检测



语义分割

自监督学习

➤ 借口任务

借口任务 (pretext tasks) :

- Exemplar CNN
- Context Prediction
- Jigsaw Puzzle
- Colorization
- Rotation
- Learning to Count
- Learning by inpainting
- Split-brain AutoEncoder
- Deep Cluster
- Non parametric Instance Discrimination
- CPC (Contrastive Predictive Coding)
- CMC(Contrastive Multiview Coding)



开脑洞



理论

自监督学习

➤ Exemplar-CNN

取32x32的有较大梯度区域的图像patch（称为Exemplar patch）进行数据增广，包括color jitter, translation, rotation, scaling等等。

数据集中有**N张图像**，则每一张图像的Exemplar patch和它的所有数据增广为同一个类别，进行**N分类**。目的是学习图像patch对数据增广的不变性，从而得到较为鲁棒的特征。

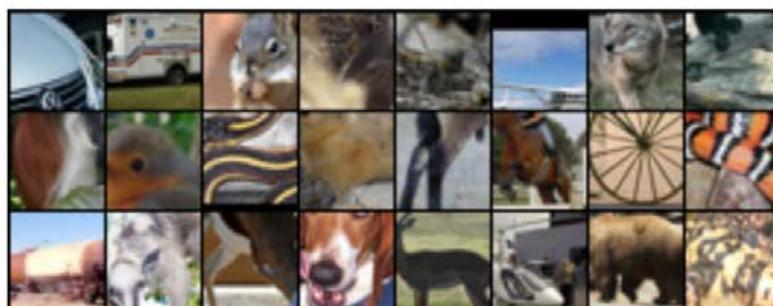


Figure 1: Exemplary patches sampled from the STL unlabeled dataset which are later augmented by various transformations to obtain surrogate data for the CNN training.



Figure 2: Several random transformations applied to one of the patches extracted from the STL unlabeled dataset. The original ('seed') patch is in the top left corner.

自监督学习

➤ Exemplar-CNN

Exemplar-CNN

该方法在STL-10, CIFAR-10, Caltech-101等小型数据集上做了实验，结果如下。可以看出，Exemplar-CNN超过了大部分方法，并且在STL-10上的表现超过了Supervised-learning

Table 1: Classification accuracies on several datasets (in percent). † Average per-class accuracy² $78.0\% \pm 0.4\%$. ‡ Average per-class accuracy $84.4\% \pm 0.6\%$.

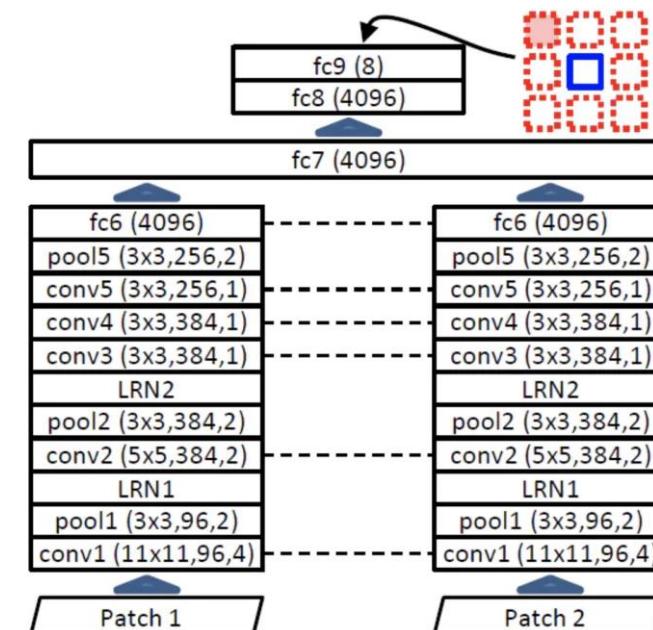
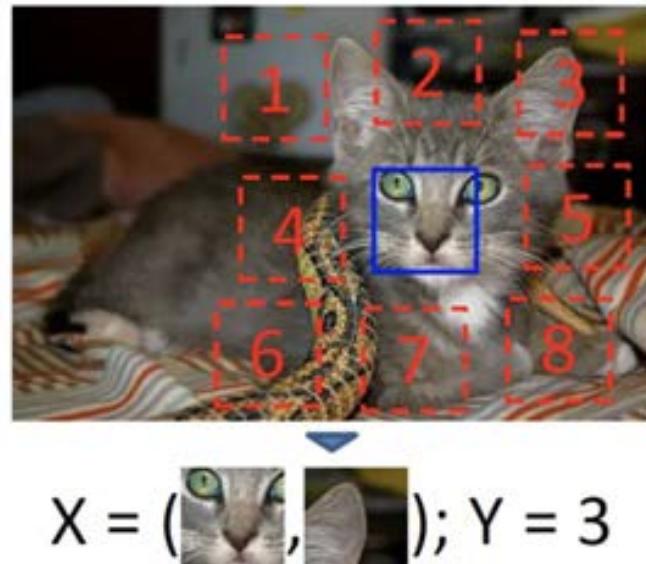
Algorithm	STL-10	CIFAR-10(400)	CIFAR-10	Caltech-101	#features
Convolutional K-means Network [26]	60.1 ± 1	70.7 ± 0.7	82.0	—	8000
Multi-way local pooling [28]	—	—	—	77.3 ± 0.6	1024×64
Slowness on videos [10]	61.0	—	—	74.6	556
Hierarchical Matching Pursuit (HMP) [27]	64.5 ± 1	—	—	—	1000
Multipath HMP [29]	—	—	—	82.5 ± 0.5	5000
View-Invariant K-means [12]	63.7	72.6 ± 0.7	81.9	—	6400
Exemplar-CNN (64c5-64c5-128f)	67.1 ± 0.3	69.7 ± 0.3	75.7	79.8 ± 0.5 †	256
Exemplar-CNN (64c5-128c5-256c5-512f)	72.8 ± 0.4	75.3 ± 0.2	82.0	85.5 ± 0.4‡	960
Supervised state of the art	70.1 [30]	—	91.2 [31]	91.44 [32]	—

自监督学习

➤ Context Prediction

Context Prediction

预测图像patch之间的相对位置，对于一张图像如下左图，划分9个patch，然后对相邻的patch进行方位的预测，例如2在中间蓝色patch的上方，3在右上方。由此变成一个方位的分类问题。具体的网络如下右图，采用一个 siamese network得到两个块的特征并concatenate进行分类



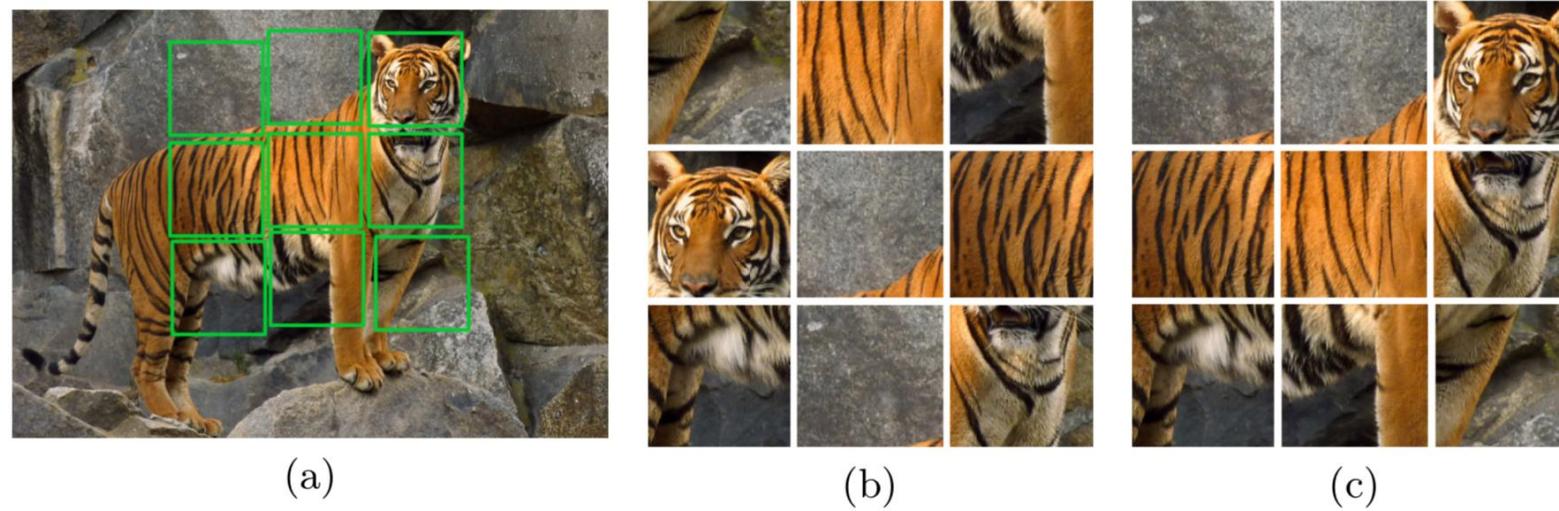
自监督学习

➤ Jigsaw Puzzle

将图像分成 3×3 的图像块，如下图(a)，然后把这9块打乱，变成 (b) 的样子。
借口任务就是能够“还原”到原来的有序的样子 (c)。

具体的做法是，把所有的排列 ($9! = 362880$) 进行采样，取出其中的100种较为混乱的，然后把这100个置換作为监督信号。即将还原问题转化为一个100类的多分类问题。

这个任务希望通过学习jigsaw puzzle来学习到关于物体部分与部分之间的关联。



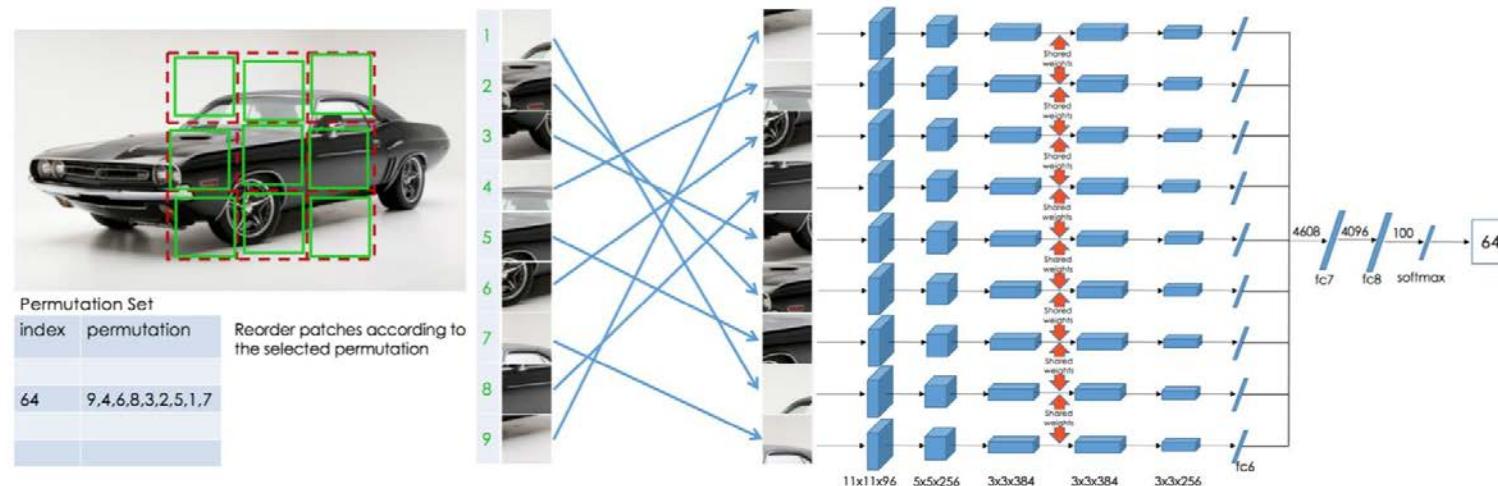
自监督学习

➤ Jigsaw Puzzle

对于具体的网络设计，结构如下图，采用siamese网络来提取各个部分的特征，最后进行融合来判断打乱的置换。

Cheating

关于作弊的问题，假如正常的切分patch，那么网络会根据一些很简单的底层特征，例如边缘信息，色差等来进行作弊，瞬间便可以学会，那么整个网络并没有学到什么语义信息。



Doersch, Carl, Abhinav Gupta, and Alexei A. Efros. "Unsupervised visual representation learning by context prediction." *Proceedings of the IEEE international conference on computer vision*. 2015.

自监督学习

➤ Jigsaw Puzzle

防止作弊：将一切可能的低级特征隐藏起来，让网络无法去通过这些shortcut来判断。

1. Random crop，得到的patch进行随机的crop，例如得到每个patch为80x80，那么从中random crop得到64x64，这样避免了通过边缘信息来作弊。
2. 随机转换为灰度图或者进行color jitter，消除色差的影响。
3. 对每个patch进行normalize（减均值除标准差），避免通过整体的颜色来判断。

防作弊对于downstream task的效果如下图

Table 5: Ablation study on the impact of the shortcuts.

Gap	Normalization	Color jittering	Jigsaw task accuracy	Detection performance
✗	✓	✓	98	47.7
✓	✗	✓	90	43.5
✓	✓	✗	89	51.1
✓	✓	✓	88	52.6

自监督学习

➤ Jigsaw Puzzle

实验结果，在ImageNet2012上进行了固定层训练，下面分别是固定conv1-conv5，结果发现Jigsaw (CFN) 表现最好

Table 2: Comparison of classification results on ImageNet 2012 [9]. The numbers are obtained by averaging 10 random crops predictions.

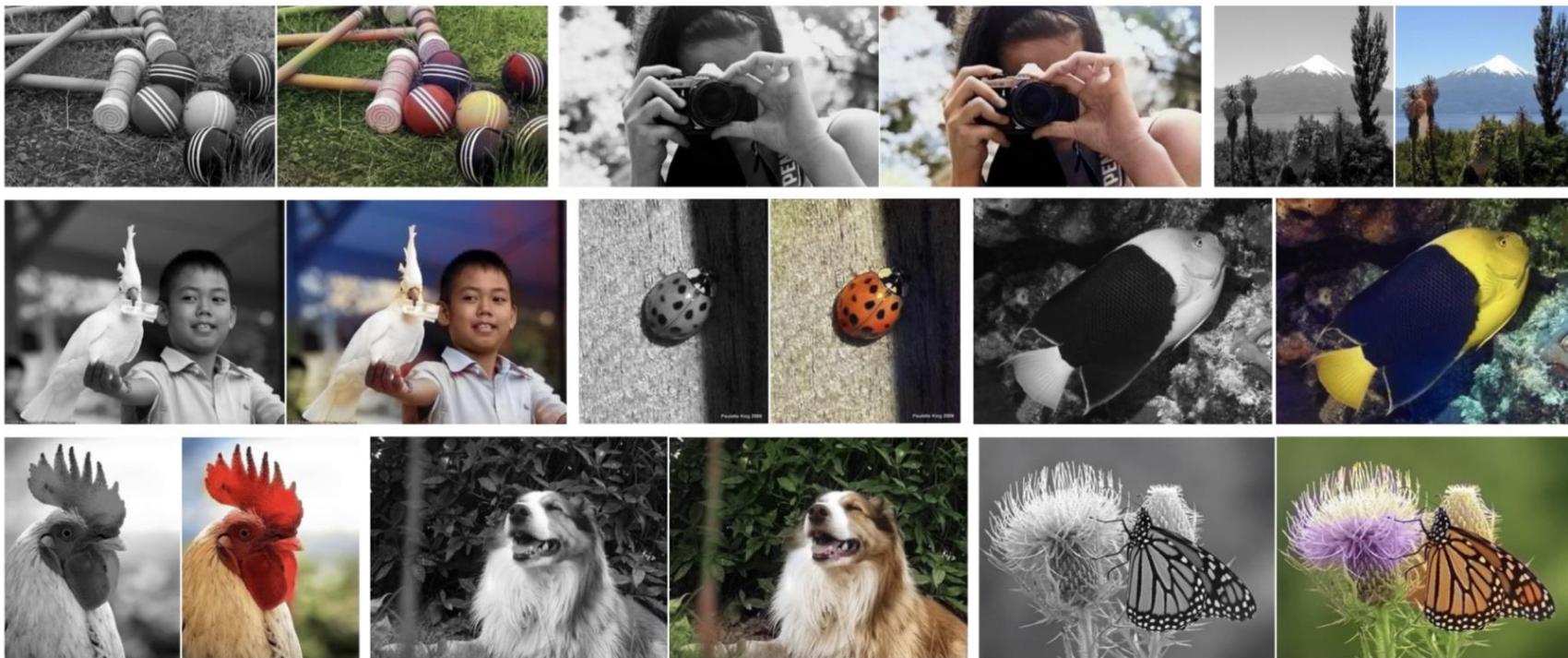
	conv1	conv2	conv3	conv4	conv5
CFN	54.7	52.8	49.7	45.3	34.6
Doersch <i>et al.</i> [10]	53.1	47.6	48.7	45.6	30.4
Wang and Gupta [39]	51.8	46.9	42.8	38.8	29.8
Random	48.5	41.0	34.8	27.1	12.0

Method	Pretraining time	Supervision	Classification	Detection	Segmentation
Krizhevsky <i>et al.</i> [25]	3 days	1000 class labels	78.2%	56.8%	48.0%
Wang and Gupta [39]	1 week	motion	58.4%	44.0%	-
Doersch <i>et al.</i> [10]	4 weeks	context	55.3%	46.6%	-
Pathak <i>et al.</i> [30]	14 hours	context	56.5%	44.5%	29.7%
Ours	2.5 days	context	67.6%	53.2%	37.6%

自监督学习

➤ Colorization

通过给灰度图上色来进行自监督学习。如下图所示，灰度图的颜色信息往往需要网络对纹理有初步的认识。学习到一些关于纹理相关的特征，这些特征对于下游任务例如分类检测等也是有帮助的。

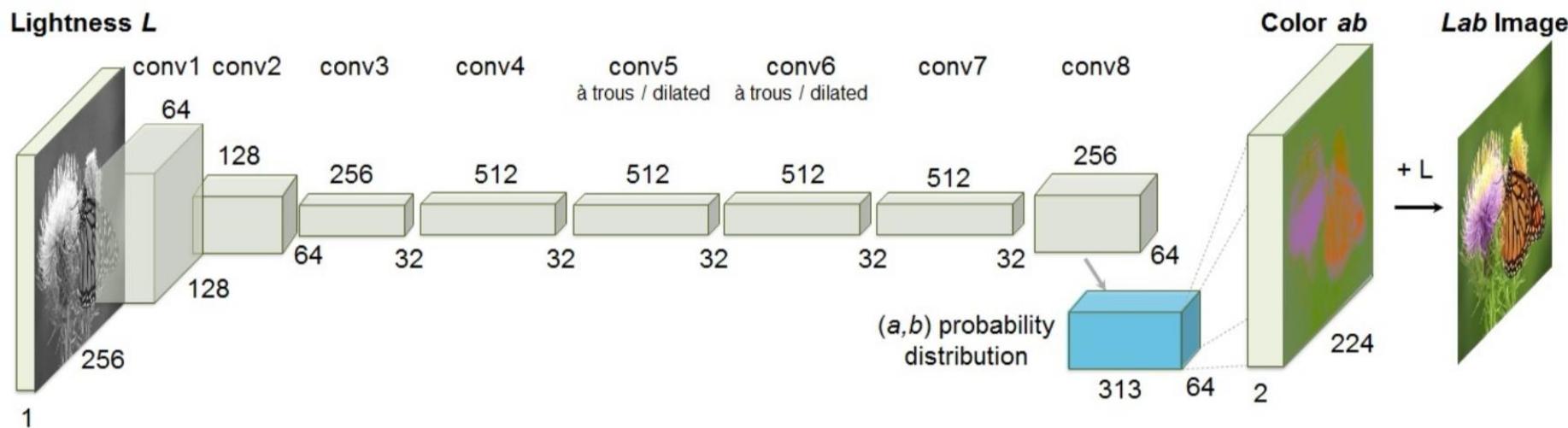


自监督学习

➤ Colorization

网络结构如下图所示，为了更好的表示颜色，采用比较适合人眼对色彩感知的**Lab**颜色系统，输入L channel，预测ab两个channel。

由于一个像素的颜色是**multimodal**的（即一个像素的**可能颜色有多个**），如果采用简单的回归，则会导致收敛到多个可能颜色的均值，总体会导致饱和度偏低。为此作者将此变成一个分类问题。



自监督学习

➤ Colorization

染色效果如下，Colorization不仅仅是一个自监督方法，同时自身的借口任务（染色）也是一个有意思的问题。

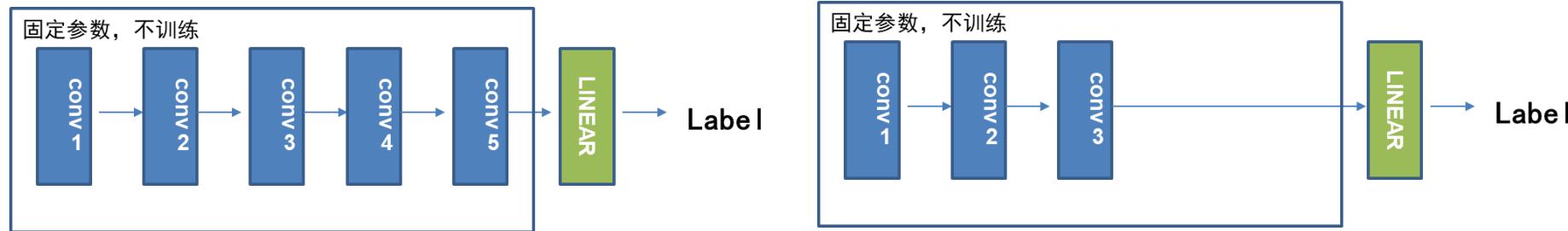


自监督学习

➤ Colorization

此外，此论文还提出了被后面的各种自监督算法广泛采用的验证方式，即在ImageNet上的线性分类任务。具体来说，将卷积层固定参数，并添加一个线性层，下游任务使用ImageNet标签来训练，但是只训练线性层。此任务主要是检验学到的表征在复杂数据标签下的线性可分性。

当然也可以测试前面的卷积层产生的特征的线性可分性，例如前三层，此时需要把conv4, conv5移除。



自监督学习

➤ Colorization

实验结果如下，左边是所有conv层都固定，接一个线性层到conv1、conv2... conv5上，线性分类ImageNet的结果。右图是在PASCAL VOC2007三个任务上的Finetune表现。结果显示colorization超过了其他的自监督方法。

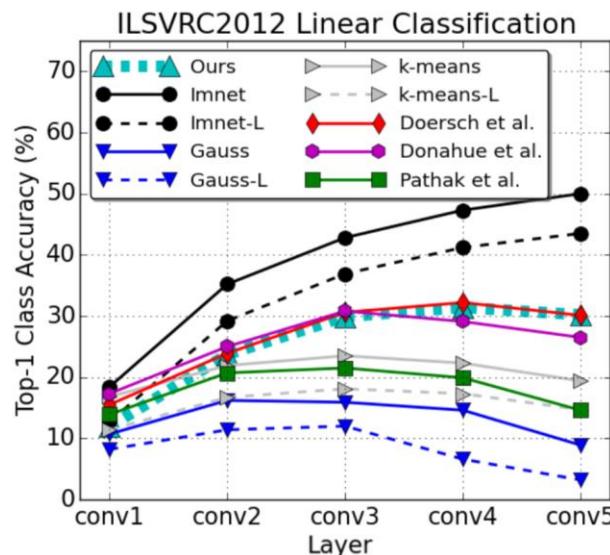


Fig. 7. ImageNet Linear Classification

Dataset and Task Generalization on PASCAL [37]								
fine-tune layers	[Ref]	Class. (%mAP)			Det. (%mAP)		Seg. (%mIU)	
		fc8	fc6-8	all	[Ref]	all	[Ref]	all
ImageNet [38]	-	76.8	78.9	79.9	[36]	56.8	[42]	48.0
Gaussian	[10]	-	-	53.3	[10]	43.4	[10]	19.8
Autoencoder	[16]	24.8	16.0	53.8	[10]	41.9	[10]	25.2
k-means [36]	[16]	32.0	39.2	56.6	[36]	45.6	[16]	32.6
Agrawal et al. [8]	[16]	31.2	31.0	54.2	[36]	43.9	-	-
Wang & Gupta [15]	-	28.1	52.2	58.7	[36]	47.4	-	-
*Doersch et al. [14]	[16]	44.7	55.1	65.3	[36]	51.1	-	-
*Pathak et al. [10]	[10]	-	-	56.5	[10]	44.5	[10]	29.7
*Donahue et al. [16]	-	38.2	50.2	58.6	[16]	46.2	[16]	34.9
Ours (gray)	-	52.4	61.5	65.9	-	46.1	-	35.0
Ours (color)	-	52.4	61.5	65.6	-	46.9	-	35.6

Table 2. PASCAL Tests

自监督学习

➤ Rotation

将无标签数据集中的图像进行旋转0, 90, 180, 270度，旋转后的图像经过网络后变成一个四分类问题。

Rotation的一个优点是不需要防止作弊，因为旋转很难留下容易作弊的线索。

Rotation虽然非常简单，但是却效果非常好。



Figure 1: Images rotated by random multiples of 90 degrees (e.g., 0, 90, 180, or 270 degrees). The core intuition of our self-supervised feature learning approach is that if someone is not aware of the concepts of the objects depicted in the images, he cannot recognize the rotation that was applied to them.

Gidaris, Spyros, Praveer Singh, and Nikos Komodakis. "Unsupervised representation learning by predicting image rotations." *arXiv preprint arXiv:1803.07728* (2018).

自监督学习

➤ Rotation

下面是Rotation在ImageNet和Places205数据集上的线性分类表现。可以看出来，Rotation确实超过了其他的方法很多。Rotation实现起来非常简单，效果确实最好的。

ImageNet

Method	Conv1	Conv2	Conv3	Conv4	Conv5
ImageNet labels	19.3	36.3	44.2	48.3	50.5
Random	11.6	17.1	16.9	16.3	14.1
Random rescaled [Krähenbühl et al. (2015)]	17.5	23.0	24.5	23.2	20.6
Context (Doersch et al., 2015)	16.2	23.3	30.2	31.7	29.6
Context Encoders (Pathak et al., 2016b)	14.1	20.7	21.0	19.8	15.5
Colorization (Zhang et al., 2016a)	12.5	24.5	30.4	31.5	30.3
Jigsaw Puzzles (Noroozi & Favaro, 2016)	18.2	28.8	34.0	33.9	27.1
BIGAN (Donahue et al., 2016)	17.7	24.5	31.0	29.9	28.0
Split-Brain (Zhang et al., 2016b)	17.7	29.3	35.4	35.2	32.8
Counting (Noroozi et al., 2017)	18.0	30.6	34.3	32.5	25.7
(Ours) RotNet	18.8	31.7	38.7	38.2	36.5

Places205

Method	Conv1	Conv2	Conv3	Conv4	Conv5
Places labels [Zhou et al. (2014)]	22.1	35.1	40.2	43.3	44.6
ImageNet labels	22.7	34.8	38.4	39.4	38.7
Random	15.7	20.3	19.8	19.1	17.5
Random rescaled [Krähenbühl et al. (2015)]	21.4	26.2	27.1	26.1	24.0
Context (Doersch et al., 2015)	19.7	26.7	31.9	32.7	30.9
Context Encoders (Pathak et al., 2016b)	18.2	23.2	23.4	21.9	18.4
Colorization (Zhang et al., 2016a)	16.0	25.7	29.6	30.3	29.7
Jigsaw Puzzles (Noroozi & Favaro, 2016)	23.0	31.9	35.0	34.2	29.3
BIGAN (Donahue et al., 2016)	22.0	28.7	31.8	31.3	29.7
Split-Brain (Zhang et al., 2016b)	21.3	30.7	34.0	34.1	32.5
Counting (Noroozi et al., 2017)	23.3	33.9	36.3	34.7	29.6
(Ours) RotNet	21.5	31.0	<u>35.1</u>	<u>34.6</u>	<u>33.7</u>

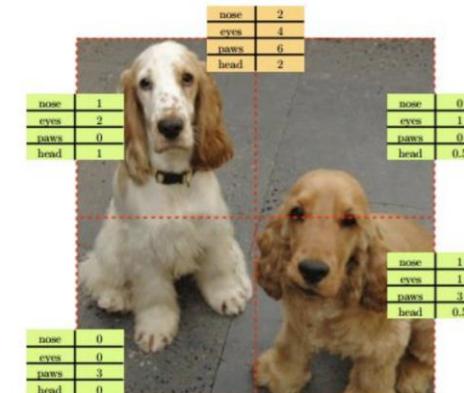
自监督学习

➤ Learning by Learning to Count

通过计数一个图像patch中primitive elements。

将一个图像分为 2×2 个patch，每个patch中的primitive elements的数目之和应该等同于整张图像的primitive elements的数目，即一张图片分成若干部分后得到的特征求和，应当等于整个图片直接的特征求和，如下式所示，指代数元素的函数，用神经网络来代替，

$$\phi(D \circ \mathbf{x}) = \sum_{j=1}^4 \phi(T_j \circ \mathbf{x}).$$



可以容易看出，四个patch的基本元素的个数之和应该等与整张图像的个数之和

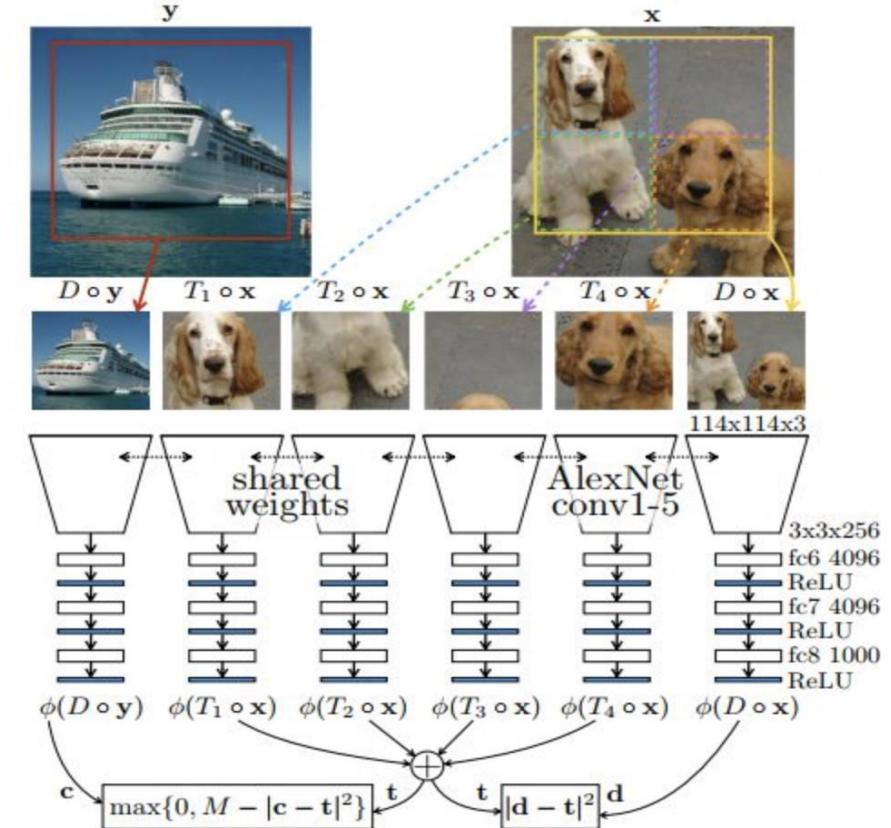
自监督学习

➤ Learning by Learning to Count

$$\phi(D \circ \mathbf{x}) = \sum_{j=1}^4 \phi(T_j \circ \mathbf{x}).$$

然而，若只用神经网络进行上述等式的建模，会造成trivial solution。即全部元素数目均为0即可满足等式。

为避免这一点，引入负样本，即另一张图像的元素总和和这张图像四个patch之和是尽量不等的。所以整体网络也很简单，如右图所示。采用contrastive loss来做到对正样本距离尽量小，对负样本距离尽可能大（实际上超过一个阈值M即可）。



自监督学习

➤ Learning by Learning to Count

下图为ImageNet线性分类的结果，右图为Pascal VOC2007 fine-tune的结果，结果还不错。但是并不是对比中最好的结果。

Method	conv1	conv2	conv3	conv4	conv5
Supervised [20]	19.3	36.3	44.2	48.3	50.5
Random	11.6	17.1	16.9	16.3	14.1
Context [9]	16.2	23.3	30.2	31.7	29.6
Jigsaw [30]	18.2	28.8	34.0	<u>33.9</u>	27.1
ContextEncoder [33]	14.1	20.7	21.0	<u>19.8</u>	15.5
Adversarial [10]	17.7	24.5	31.0	29.9	28.0
Colorization [43]	12.5	24.5	30.4	31.5	<u>30.3</u>
Split-Brain [44]	17.7	<u>29.3</u>	35.4	35.2	32.8
Counting	<u>18.0</u>	30.6	<u>34.3</u>	32.5	25.7

Method	Ref	Class.	Det.	Segm.
Supervised [20]	[43]	79.9	56.8	48.0
Random	[33]	53.3	43.4	19.8
Context [9]	[19]	55.3	46.6	-
Context [9]*	[19]	65.3	51.1	-
Jigsaw [30]	[30]	<u>67.6</u>	53.2	<u>37.6</u>
ego-motion [1]	[1]	52.9	41.8	-
ego-motion [1]*	[1]	54.2	43.9	-
Adversarial [10]*	[10]	58.6	46.2	34.9
ContextEncoder [33]	[33]	56.5	44.5	29.7
Sound [31]	[44]	54.4	44.0	-
Sound [31]*	[44]	61.3	-	-
Video [41]	[19]	62.8	47.4	-
Video [41]*	[19]	63.1	47.2	-
Colorization [43]*	[43]	65.9	46.9	35.6
Split-Brain [44]*	[44]	67.1	46.7	36.0
ColorProxy [22]	[22]	65.9	-	38.0
WatchingObjectsMove [32]	[32]	61.0	<u>52.2</u>	-
Counting		67.7	51.4	36.6

Table 1: Evaluation of transfer learning on PASCAL. Classification and detection are evaluated on PASCAL VOC 2007 in the frameworks introduced in [19] and [11] respectively. Both tasks are evaluated using mean average precision (mAP) as a performance measure. Segmentation is evaluated on PASCAL VOC 2012 in the framework of [26], which reports mean intersection over union (mIoU). (*) denotes the use of the data initialization method [19].

自监督学习

➤ Learning by Inpainting

思路很简单，通过填补图像中空缺的部分来学习，如右图所示。该方法提出context encoder，即将缺失的图像输入，尽量还原出缺失的部分的autoencoder。如下图所示。

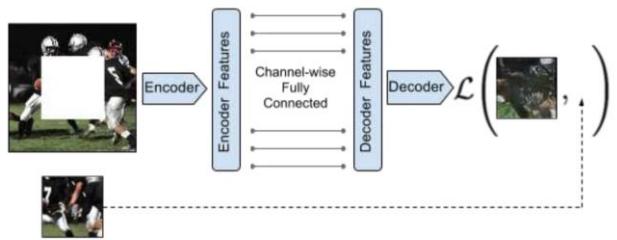
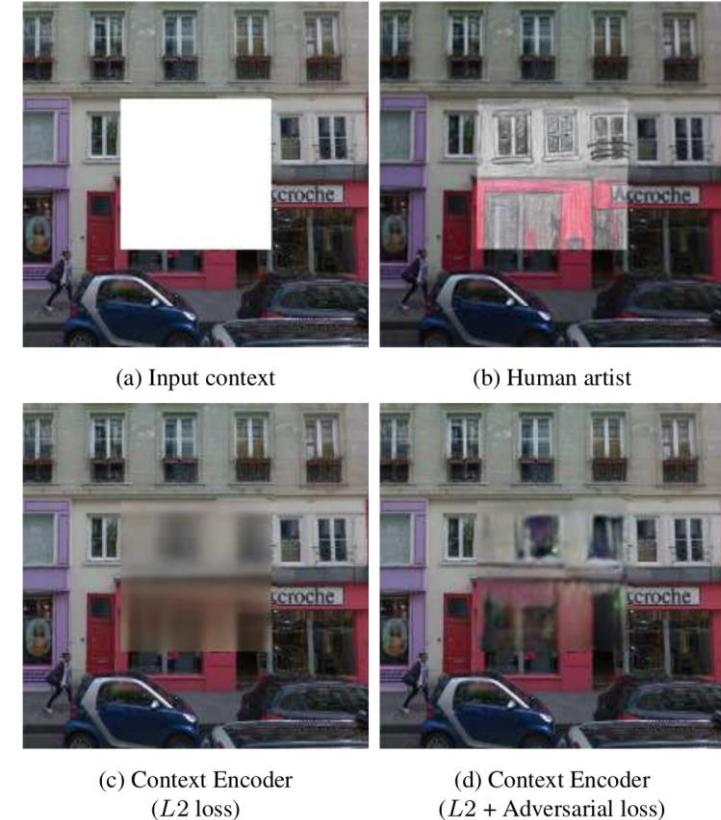


Figure 2: Context Encoder. The context image is passed through the encoder to obtain features which are connected to the decoder using channel-wise fully-connected layer as described in Section 3.1. The decoder then produces the missing regions in the image.



自监督学习

➤ Learning by Inpainting

Inpainting的结果如下图所示，可以看出，采用L2会导致**模糊**的现象，采用对抗损失会导致**不连续**的现象，而两者结合起来结果更好些，尤其是在中间一行的表现。而提出的方法比其他方法要好很多

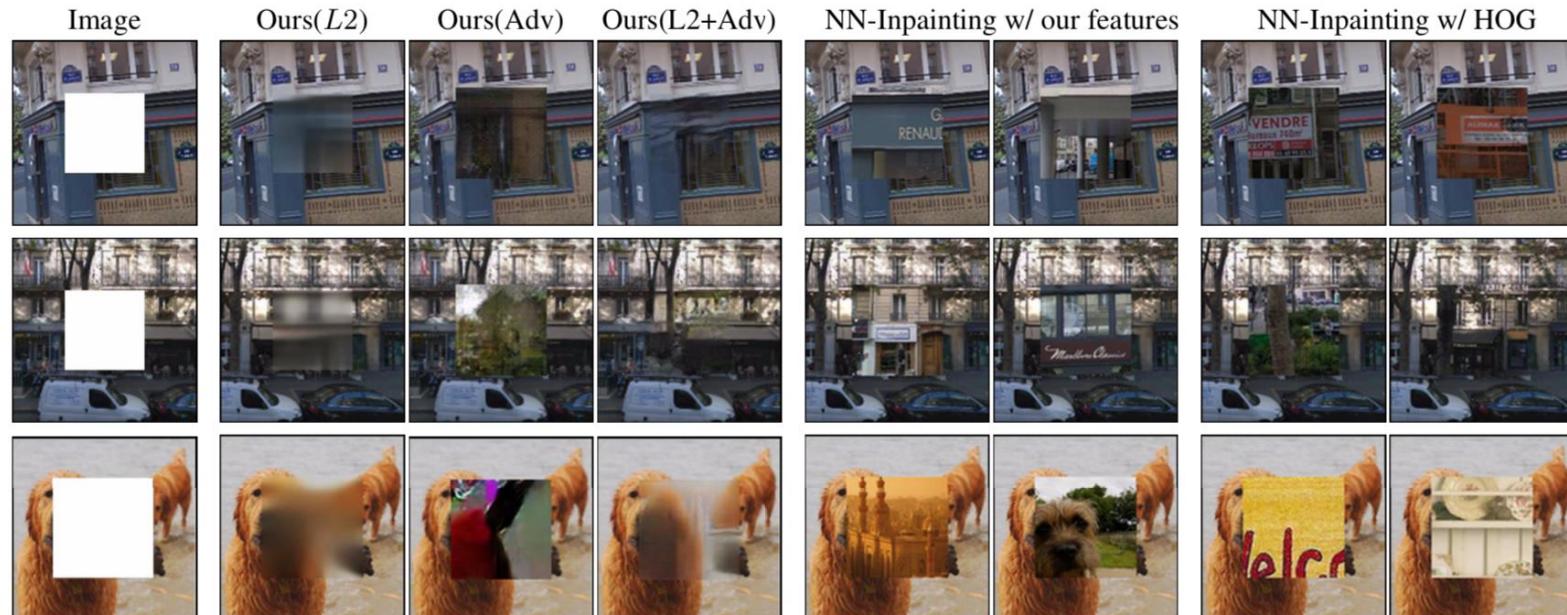


Figure 6: Semantic Inpainting using different methods. Context Encoder with just L2 are well aligned, but not sharp. Using adversarial loss, results are sharp but not coherent. Joint loss alleviate the weaknesses of each of them. The last two columns are the results if we plug-in the best nearest neighbor (NN) patch in the masked region.

自监督学习

➤ Learning by Inpainting

无监督结果如下，在Pascal VOC2007上的实验。Inpainting作为一个无监督方法效果非常一般。生成模型在无监督学习方面的表现普遍较差。

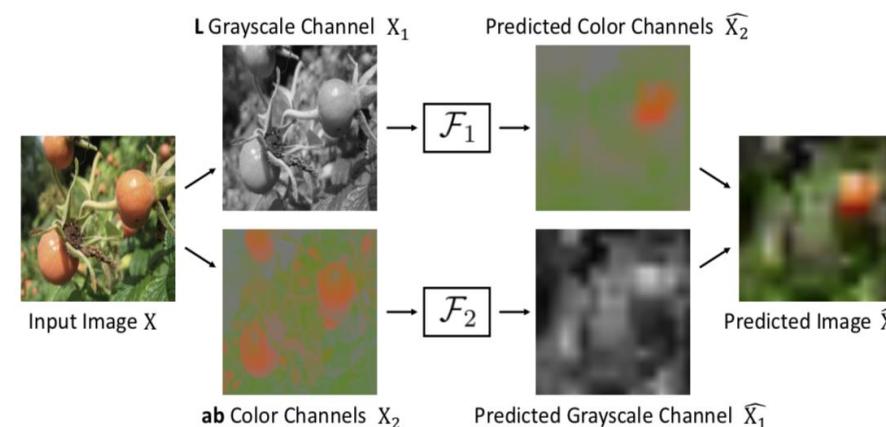
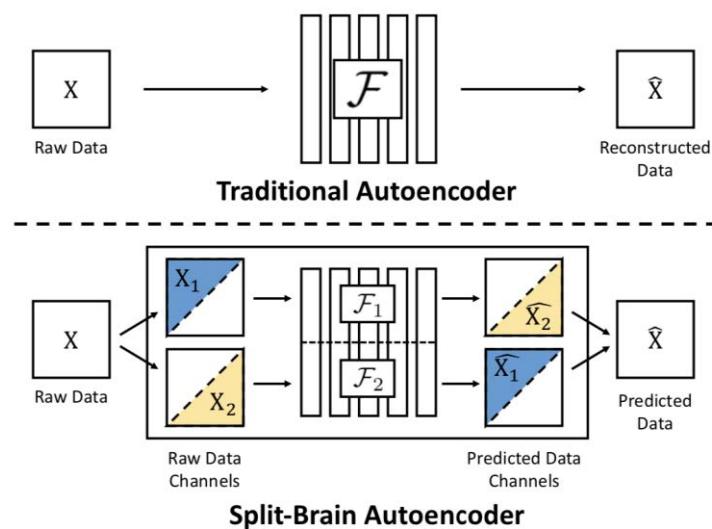
Pretraining Method	Supervision	Pretraining time	Classification	Detection	Segmentation
ImageNet [26]	1000 class labels	3 days	78.2%	56.8%	48.0%
Random Gaussian	initialization	< 1 minute	53.3%	43.4%	19.8%
Autoencoder	-	14 hours	53.8%	41.9%	25.2%
Agrawal <i>et al.</i> [1]	egomotion	10 hours	52.9%	41.8%	-
Doersh <i>et al.</i> [7]	context	4 weeks	55.3%	46.6%	-
Wang <i>et al.</i> [39]	motion	1 week	58.4%	44.0%	-
Ours	context	14 hours	56.5%	44.5%	29.7%

自监督学习

➤ Split-Brain Autoencoders

思路也很简单，使两个channel的数据互相预测彼此。如下面右图所示，将一个图像分割成两部分，一部分是灰度，另一部分是色彩。两者彼此互相预测。

作者使用的网络比较独特，将一个完整网络一分为二，但这并不重要，使用两个独立的网络和这种设计是基本差不多的。



自监督学习

➤ Split-Brain Autoencoders

实验结果如右图所示，作者发现类似于color之间的转换（或者就是colorization）是有利无监督学习的，并且方法的效果也是很不错的。

Task Generalization on ImageNet Classification [37]					
Method	conv1	conv2	conv3	conv4	conv5
ImageNet-labels [26]	19.3	36.3	44.2	48.3	50.5
Gaussian	11.6	17.1	16.9	16.3	14.1
Krähenbühl et al. [25]	17.5	23.0	24.5	23.2	20.6
¹ Noroozi & Favaro [31]	19.2	30.1	34.7	33.9	28.3
Doersch et al. [8]	16.2	23.3	30.2	31.7	29.6
Donahue et al. [9]	17.7	24.5	31.0	29.9	28.0
Pathak et al. [35]	14.1	20.7	21.0	19.8	15.5
Zhang et al. [49]	13.1	24.8	31.0	32.6	31.8
Lab→Lab	12.9	20.1	18.5	15.1	11.5
Lab(drop50)→Lab	12.1	20.4	19.7	16.1	12.3
L→ab(cl)	12.5	25.4	32.4	33.1	32.0
L→ab(reg)	12.3	23.5	29.6	31.1	30.1
ab→L(cl)	11.6	19.2	22.6	21.7	19.2
ab→L(reg)	11.5	19.4	23.5	23.9	21.7
(L,ab)→(ab,L)	15.1	22.6	24.4	23.2	21.1
(L,ab,Lab)→(ab,L,Lab)	15.4	22.9	24.0	22.0	18.9
Ensembled L→ab	11.7	23.7	30.9	32.2	31.3
Split-Brain Auto (reg,reg)	17.4	27.9	33.6	34.2	32.3
Split-Brain Auto (cl,cl)	17.7	29.3	35.4	35.2	32.8

Table 2: Task Generalization on ImageNet Classification

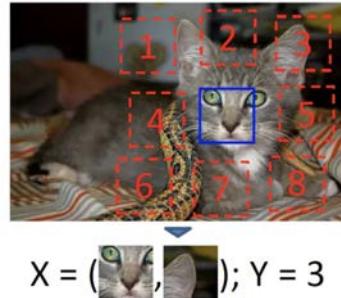
To test unsupervised feature representations, we train linear logistic regression classifiers on top of each layer to perform 1000-way ImageNet classification, as proposed in [49]. All weights are frozen and feature maps spatially resized to be ~9000 dimensions. All methods use AlexNet variants [26], and were pre-trained on ImageNet without labels, except for **ImageNet-labels**. Note that the proposed split-brain autoencoder achieves the best performance on all layers across unsupervised methods.

自监督学习

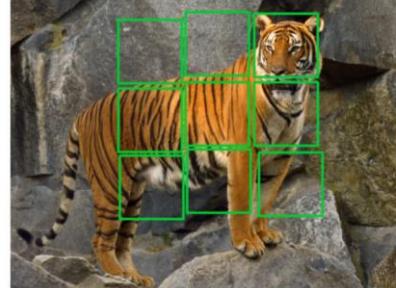
➤ 总结



Exemplar-CNN



Context Prediction



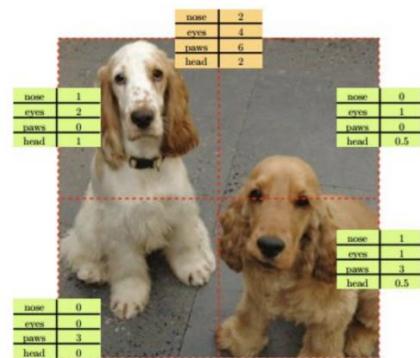
Jigsaw Puzzle



Colorization



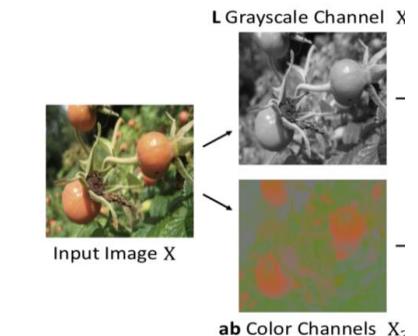
Rotation



Learning by Learning to Count

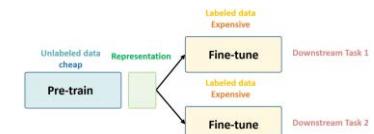


Learning by Inpainting



Split-Brain Autoencoders

Unsupervised Pre-train, Supervised Fine-tune.



自监督学习

➤ 总结

借口任务 (pretext tasks) :

- Exemplar CNN
- Context Prediction
- Jigsaw Puzzle
- Colorization
- Rotation
- Learning to Count
- Learning by inpainting
- Split-brain AutoEncoder
- **Deep Cluster**
- Non parametric Instance Discrimination
- CPC (Contrastive Predictive Coding)
- CMC(Contrastive Multiview Coding)



开脑洞

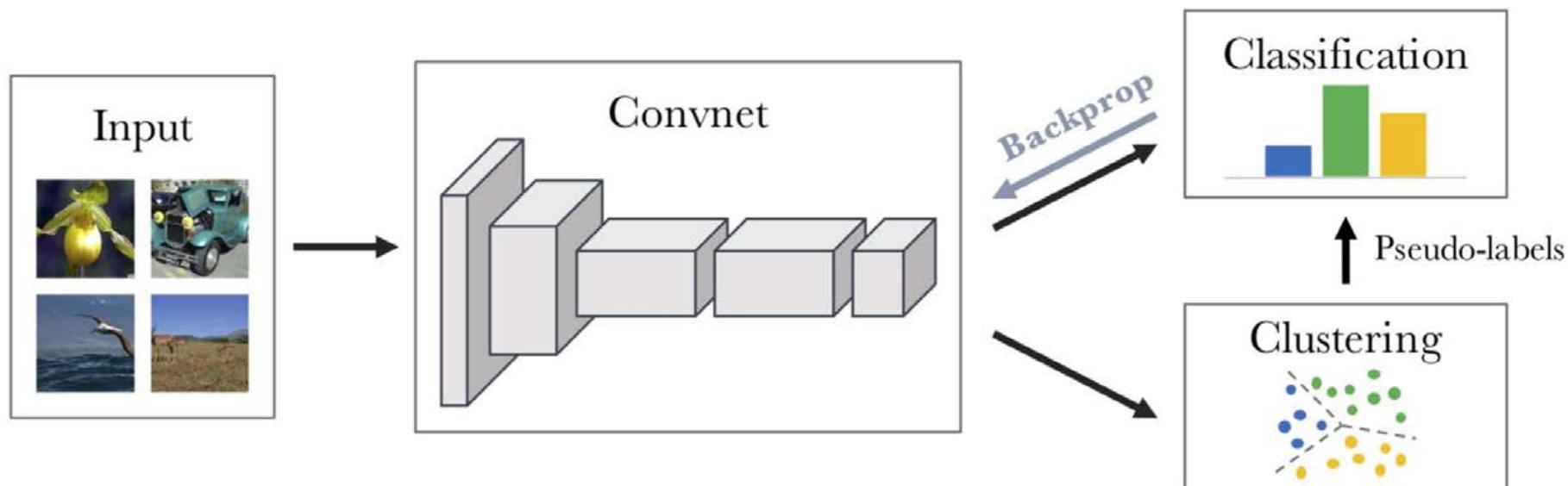


理论

自监督学习

➤ Deep Cluster

对图像经过CNN后的特征进行聚类，并把得到的**聚类编号作为监督信号**进行分类，如此迭代。如下图所示。

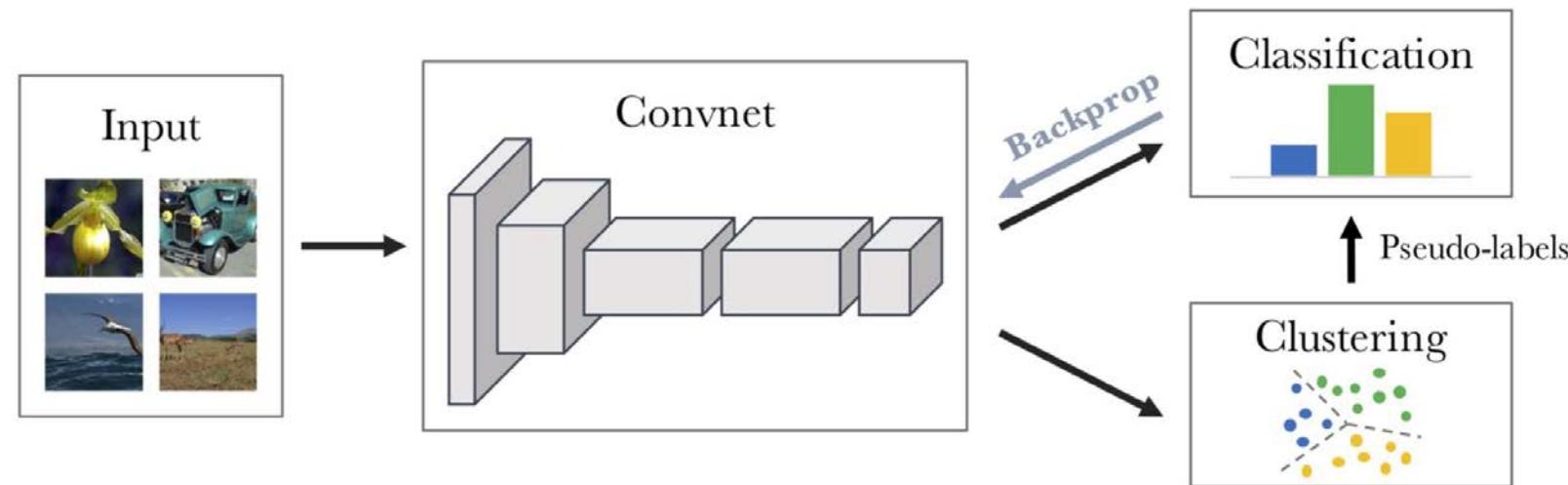


自监督学习

➤ Deep Cluster

Deep Cluster相对于其他的自监督学习，更像是一种无监督学习方法，即并没有刻意地创造与模态（例如图像）相关的先验知识，而是利用通用的聚类方法来找到聚类中心，并使得每个样本离聚类中心更近一些。至于为什么这样的方法会有效（自己为自己提供监督信息）。

聚类使得每个聚类中心产生，而让每个样本离聚类中心更近一些，就让整体的分类边界不会穿过数据的高密度区域。而一个好的分类器的分类边界不穿过高密度区域是一个十分有用的假设（在半监督学习中经常涉及）。



自监督学习

➤ Deep Cluster

DeepCluster超过了绝大多数方法。这是一个很有趣的结果，只用聚类来提供伪标签却超过了大部分“精心设计”的算法。往往简单的是更有效的。

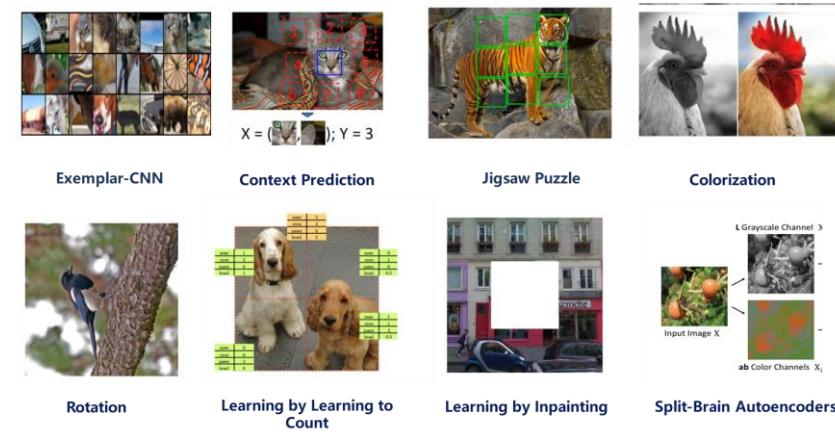
Method	ImageNet					Places				
	conv1	conv2	conv3	conv4	conv5	conv1	conv2	conv3	conv4	conv5
Places labels	—	—	—	—	—	22.1	35.1	40.2	43.3	44.6
ImageNet labels	19.3	36.3	44.2	48.3	50.5	22.7	34.8	38.4	39.4	38.7
Random	11.6	17.1	16.9	16.3	14.1	15.7	20.3	19.8	19.1	17.5
Pathak <i>et al.</i> [38]	14.1	20.7	21.0	19.8	15.5	18.2	23.2	23.4	21.9	18.4
Doersch <i>et al.</i> [25]	16.2	23.3	30.2	31.7	29.6	19.7	26.7	31.9	32.7	30.9
Zhang <i>et al.</i> [28]	12.5	24.5	30.4	31.5	30.3	16.0	25.7	29.6	30.3	29.7
Donahue <i>et al.</i> [20]	17.7	24.5	31.0	29.9	28.0	21.4	26.2	27.1	26.1	24.0
Noroozi and Favaro [26]	18.2	28.8	34.0	33.9	27.1	23.0	32.1	35.5	34.8	31.3
Noroozi <i>et al.</i> [45]	18.0	30.6	34.3	32.5	25.7	23.3	33.9	36.3	34.7	29.6
Zhang <i>et al.</i> [43]	17.7	29.3	35.4	35.2	32.8	21.3	30.7	34.0	34.1	32.5
DeepCluster	12.9	29.2	38.2	39.8	36.1	18.6	30.8	37.0	37.5	33.1

Table 1: Linear classification on ImageNet and Places using activations from the convolutional layers of an AlexNet as features. We report classification accuracy on the central crop. Numbers for other methods are from Zhang *et al.* [43].

自监督学习

借口任务 (pretext tasks) :

- Exemplar CNN
- Context Prediction
- Jigsaw Puzzle
- Colorization
- Rotation
- Learning to Count
- Learning by inpainting
- Split-brain AutoEncoder
- Deep Cluster
- Non parametric Instance Discrimination
- CPC (Contrastive Predictive Coding)
- CMC(Contrastive Multiview Coding)



对比学习

自监督学习

➤ 重温Exemplar-CNN

取32x32的有较大梯度区域的图像patch (称为Exemplar patch) 进行数据增广，包括color jitter, translation, rotation, scaling等等。

数据集中有N张图像，则每一张图像的Exemplar patch和它的所有数据增广为同一个类别，进行N分类。目的是学习图像patch对数据增广的不变性，从而得到较为鲁棒的特征。



Figure 1: Exemplary patches sampled from the STL unlabeled dataset which are later augmented by various transformations to obtain surrogate data for the CNN training.



Figure 2: Several random transformations applied to one of the patches extracted from the STL unlabeled dataset. The original ('seed') patch is in the top left corner.

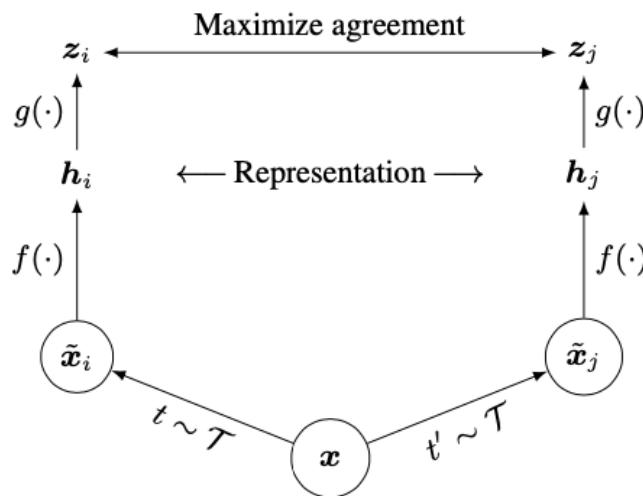
自监督学习

➤ 对比学习

对比学习学到的表征效果在大部分任务上已经超过了监督学习

Contrastive Learning思想非常简单：将同一个样本进行两次随机的数据增广，祈求两个数据增广之间的距离尽可能小。

核心思想就是学习对数据增广的不变性。



自监督学习

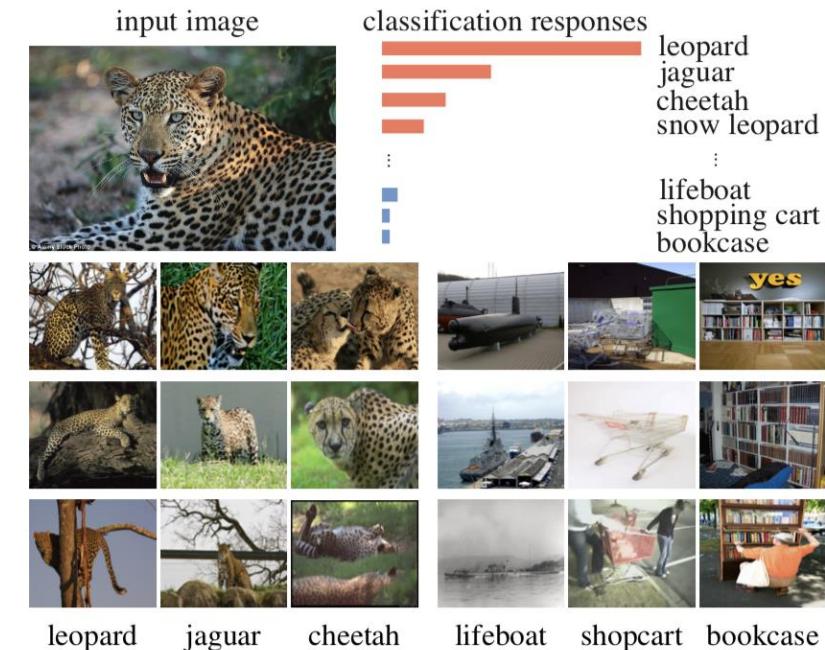
➤ 对比学习

Non parametric Instance Discrimination

实际上是Exemplar CNN的思想，对augmentation产生不变表示。但作者给了另一种解释，并且使用了Non parametric方法，性能表现非常好。

作者的解释：如右图所示，对于普通的分类问题，一张图像的top5分类应该是语义相近的，例如右图的美洲豹(leopard)，top5分类为：美洲豹(leopard)，美洲虎(jaguar)，猎豹(cheetah)...

作者认为如果将每一个instance作为一个类别，那么相似的instance也会出现这种现象。也就是相似的instance被映射到相似的特征。



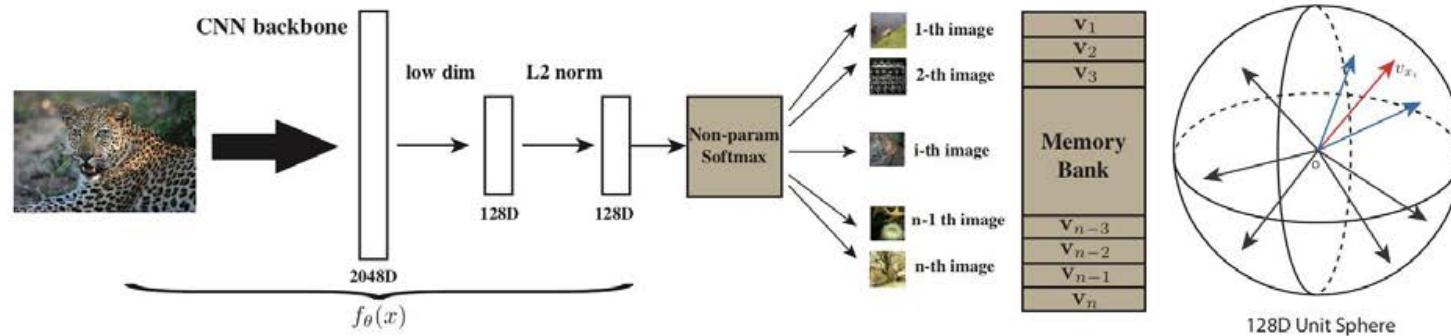
自监督学习

➤ 对比学习

Non parametric Instance Discrimination

作者将每一个图像作为一个类别，那么对于ImageNet来说，便有大约130万个类别。如果采用常用的参数方法，那么假设倒数第二层的神经元个数为4096，采用32位浮点数，那么需要 $4096 \times 1.3 \times 10^6 \times 4$ 约为 19 GB。现有GPU很难容纳如此多的参数。

为了解决这个困难，作者采用了非参数方法，如下图所示



自监督学习

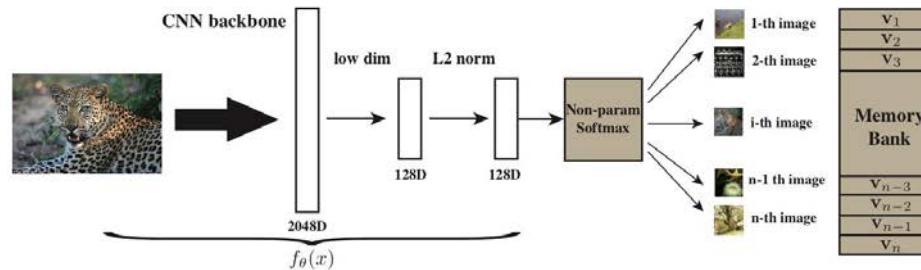
➤ 对比学习

Non parametric Instance Discrimination

参数分类器与非参数分类器：

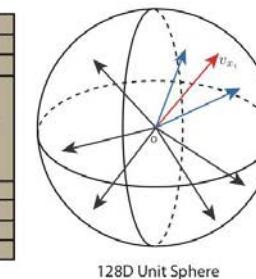
对于N分类问题，参数方法是最后一层的分类器具有参数w，w往往是(FxC)的矩阵(F为特征维度，C为类别数)。

非参数方法不同特征进行相似度计算对于v，计算其对所有其他类别的特征相似度，最后softmax分类。可以看出，非参数方法并没有参数w，而是仅仅用相似度来计算概率。



$$P(i|\mathbf{v}) = \frac{\exp(\mathbf{w}_i^T \mathbf{v})}{\sum_{j=1}^n \exp(\mathbf{w}_j^T \mathbf{v})}.$$

$$P(i|\mathbf{v}) = \frac{\exp(\mathbf{v}_i^T \mathbf{v}/\tau)}{\sum_{j=1}^n \exp(\mathbf{v}_j^T \mathbf{v}/\tau)}$$



Memory Bank

为了避免计算所有instance的特征，利用一个memory bank将所有的特征记录下来，每次更新特征便将其保存到memory bank中

自监督学习

➤ 对比学习

Non parametric Instance Discrimination

结果如右图，结果其实非常好，这类方法不仅效果好，而且对于大网络的适应性也很好，即用更好的网络得到更好的结果。这一结论并不是平凡的，事实上很多self-supervised learning的方法（如rotation）在更好的网络上并不能很好的提升效果。

method	Image Classification Accuracy on ImageNet						
	conv1	conv2	conv3	conv4	conv5	kNN	#dim
Random	11.6	17.1	16.9	16.3	14.1	3.5	10K
Data-Init [16]	17.5	23.0	24.5	23.2	20.6	-	10K
Context [2]	16.2	23.3	30.2	31.7	29.6	-	10K
Adversarial [4]	17.7	24.5	31.0	29.9	28.0	-	10K
Color [47]	13.1	24.8	31.0	32.6	31.8	-	10K
Jigsaw [27]	19.2	30.1	34.7	33.9	28.3	-	10K
Count [28]	18.0	30.6	34.3	32.5	25.7	-	10K
SplitBrain [48]	17.7	29.3	35.4	35.2	32.8	11.8	10K
Exemplar[3]			31.5			-	4.5K
Ours Alexnet	16.8	26.5	31.8	34.1	35.6	31.3	128
Ours VGG16	16.5	21.4	27.6	35.1	39.2	33.9	128
Ours Resnet18	16.0	19.9	29.8	39.0	44.5	41.0	128
Ours Resnet50	15.3	18.8	24.9	40.6	54.0	46.5	128

Table 2: Top-1 classification accuracy on ImageNet.

method	Image Classification Accuracy on Places						
	conv1	conv2	conv3	conv4	conv5	kNN	#dim
Random	15.7	20.3	19.8	19.1	17.5	3.9	10K
Data-Init [16]	21.4	26.2	27.1	26.1	24.0	-	10K
Context [2]	19.7	26.7	31.9	32.7	30.9	-	10K
Adversarial [4]	17.7	24.5	31.0	29.9	28.0	-	10K
Video [44]	20.1	28.5	29.9	29.7	27.9	-	10K
Color [47]	22.0	28.7	31.8	31.3	29.7	-	10K
Jigsaw [27]	23.0	32.1	35.5	34.8	31.3	-	10K
SplitBrain [48]	21.3	30.7	34.0	34.1	32.5	10.8	10K
Ours Alexnet	18.8	24.3	31.9	34.5	33.6	30.1	128
Ours VGG16	17.6	23.1	29.5	33.8	36.3	32.8	128
Ours Resnet18	17.8	23.0	30.1	37.0	38.1	38.6	128
Ours Resnet50	18.1	22.3	29.7	42.1	45.5	41.6	128

Table 3: Top-1 classification accuracy on Places, based directly on features learned on ImageNet, without any fine-tuning.

自监督学习

➤ 对比学习

CPC(Contrastive Predictive Coding)

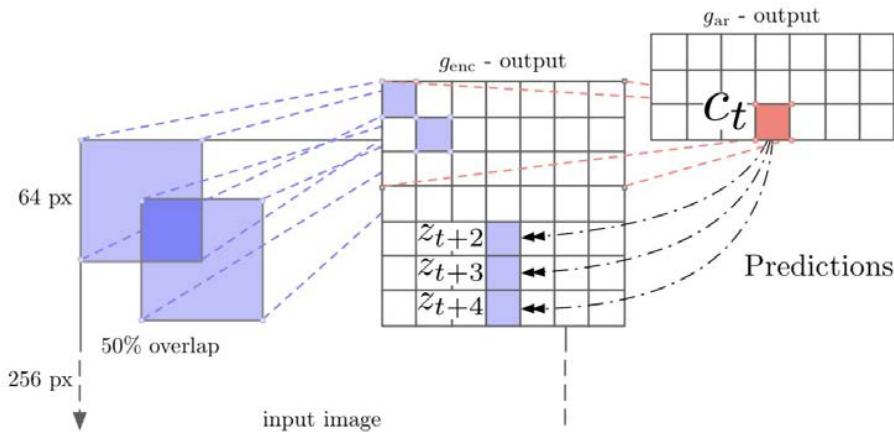


Figure 4: Visualization of Contrastive Predictive Coding for images (2D adaptation of Figure 1).

如上图所示，把一张图像进行分割，由不同的patch进行预测。利用RNN将当前的状态作为条件，预测未来状态。在图像中则用来预测后面的patch。

CPC利用Contrastive Learning来进行预测。并给出了与Mutual Information之间的联系

自监督学习

➤ 对比学习

CPC(Contrastive Predictive Coding)

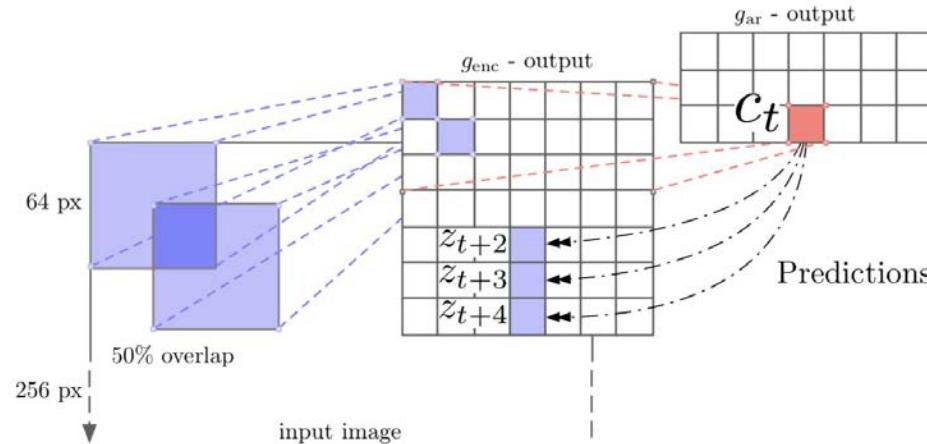


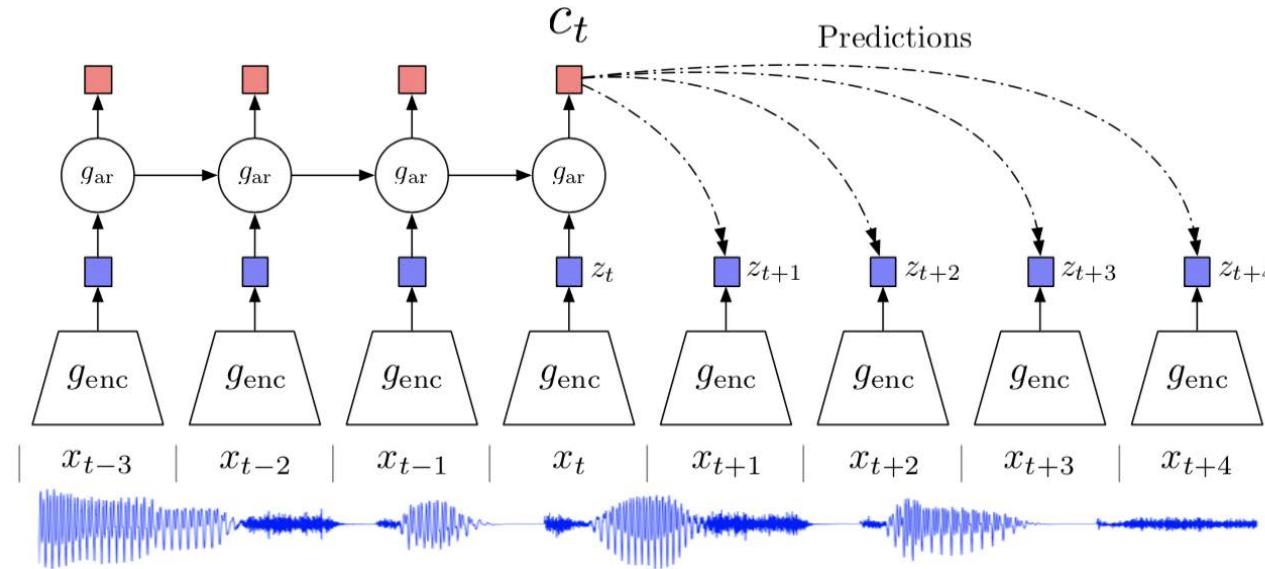
Figure 4: Visualization of Contrastive Predictive Coding for images (2D adaptation of Figure 1).

实际上，CPC的道理很简单，就是用一个图像中的一部分去预测另一部分。这个本没有什么特殊的贡献，但是我们之前提过了，生成模型的效果普遍很差，为了可以得到不错的效果，本文采用了contrastive learning，即并不是去根据一部分来生成另一部分，而是给一堆图像，让网络可以成功挑选出匹配的那一个。

自监督学习

➤ 对比学习

CPC(Contrastive Predictive Coding)



整体的网络也很简单，用一个RNN来生成context，并根据context预测下面的图像部分（或者声音，任何的模态都可以）。前面讲了，并不是真的去生成下面的图像，而是辨别哪一个是真实的对应样本。这就需要负样本（噪声样本）的存在

自监督学习

➤ 对比学习

CPC(Contrastive Predictive Coding)

效果如下，看似好像很好，但是用自己的ResNet实现去和其他论文的AlexNet实现对比难免有失公平。

但是，CPC还是很有用的，尤其是提供了非常具有启发意义的contrastive loss，被很多其他的论文所采纳

Method	Top-1 ACC
Using AlexNet conv5	
Video [28]	29.8
Relative Position [11]	30.4
BiGan [35]	34.8
Colorization [10]	35.2
Jigsaw [29] *	38.1
Using ResNet-V2	
Motion Segmentation [36]	27.6
Exemplar [36]	31.5
Relative Position [36]	36.2
Colorization [36]	39.6
CPC	48.7

Table 3: ImageNet top-1 unsupervised classification results. *Jigsaw is not directly comparable to the other AlexNet results because of architectural differences.

Method	Top-5 ACC
Motion Segmentation (MS)	48.3
Exemplar (Ex)	53.1
Relative Position (RP)	59.2
Colorization (Col)	62.5
Combination of MS + Ex + RP + Col	69.3
CPC	73.6

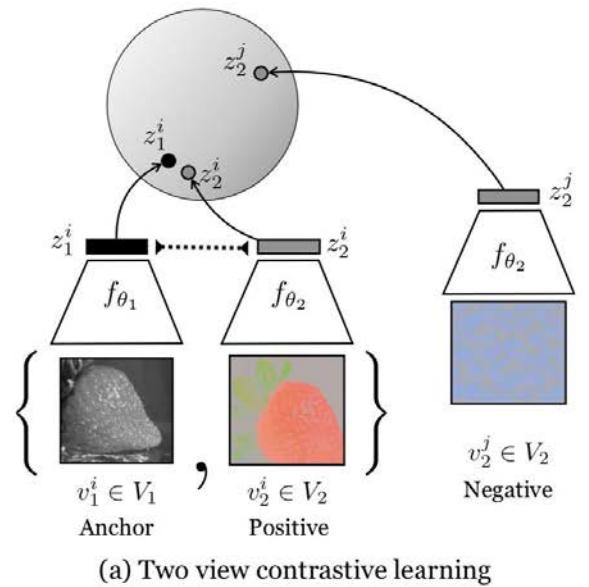
Table 4: ImageNet top-5 unsupervised classification results. Previous results with MS, Ex, RP and Col were taken from [36] and are the best reported results on this task.

$$P(i|\mathbf{v}) = \frac{\exp(\mathbf{v}_i^T \mathbf{v} / \tau)}{\sum_{j=1}^n \exp(\mathbf{v}_j^T \mathbf{v} / \tau)}, \quad (2)$$

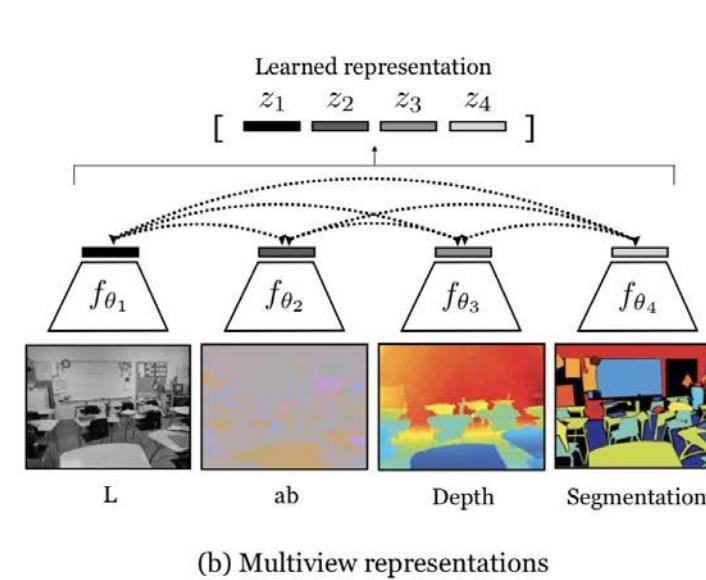
自监督学习

➤ 对比学习

CMC(Contrastive Multiview Coding)



(a) Two view contrastive learning



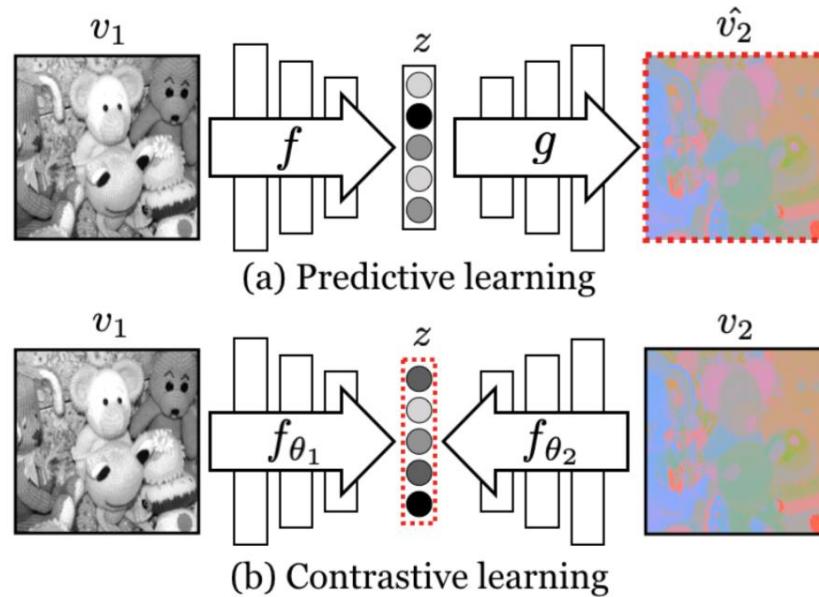
(b) Multiview representations

采用了CPC的损失函数和Split-brain的任务。由于CPC损失函数的意义在于最大化当前状态与预测状态之间的互信息，CMC的方法相当于优化不同模态之间的互信息

自监督学习

➤ 对比学习

CMC(Contrastive Multiview Coding)



具体地，将图像转化为Lab表示，L channel表示灰度，ab channel表示颜色信息。L channel的正样本为对应的ab channel 负样本为其他图像的ab 样本

自监督学习

➤ 对比学习

CMC(Contrastive Multiview Coding)

在Imagenet上的实验，结果相当好。

Method	ImageNet Classification Accuracy				
	conv1	conv2	conv3	conv4	conv5
ImageNet-Labels	19.3	36.3	44.2	48.3	50.5
Random	11.6	17.1	16.9	16.3	14.1
Data-Init (Krähenbühl et al., 2015)	17.5	23.0	24.5	23.2	20.6
Context (Doersch et al., 2015)	16.2	23.3	30.2	31.7	29.6
Colorization (Zhang et al., 2016)	13.1	24.8	31.0	32.6	31.8
Jigsaw (Noroozi & Favaro, 2016)	19.2	30.1	34.7	33.9	28.3
BiGAN (Donahue et al., 2017)	17.7	24.5	31.0	29.9	28.0
SplitBrain [†] (Zhang et al., 2017)	17.7	29.3	35.4	35.2	32.8
Counting (Noroozi et al., 2017)	18.0	30.6	34.3	32.5	25.7
Inst-Dis (Wu et al., 2018)	16.8	26.5	31.8	34.1	35.6
RotNet (Gidaris et al., 2018)	18.8	31.7	38.7	38.2	36.5
DeepCluster (Caron et al., 2018)	12.9	29.2	38.2	39.8	36.1
AET (Zhang et al., 2019)	19.3	32.8	40.6	39.7	37.7
CMC	18.4	33.5	38.1	40.4	42.6

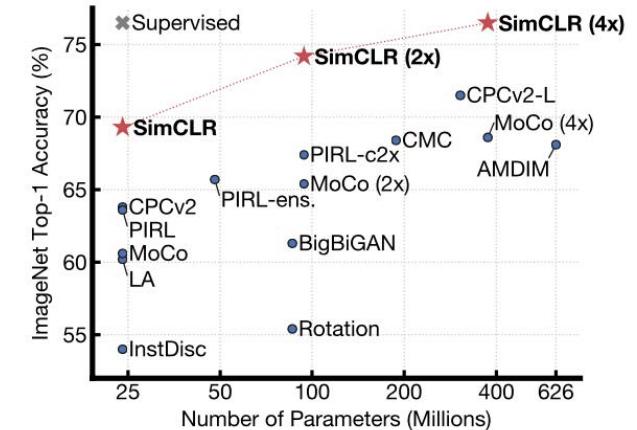
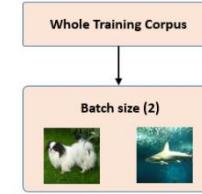
自监督学习

➤ 对比学习

SimCLR

SimCLR的思想非常简单，**抛弃掉Memory Bank**，直接使用超大的Batch，来实时的选取负样本。

前面介绍的Instance Discrimination方法提出了将正负样本存储在Memory Bank中，这样的好处是可以选择更多的负样本，来学到更uniform的特征。然而缺点是Memory Bank中存储的特征都是过时的。SimCLR的解决办法简单粗暴，采用超大的BatchSize (4096)，便可以为当前的样本选择4095个当前batch中的其他样本作为负样本来训练（其实大batch的想法CVPR2019有篇论文已经阐述了，只不过是在小数据集上做的）。

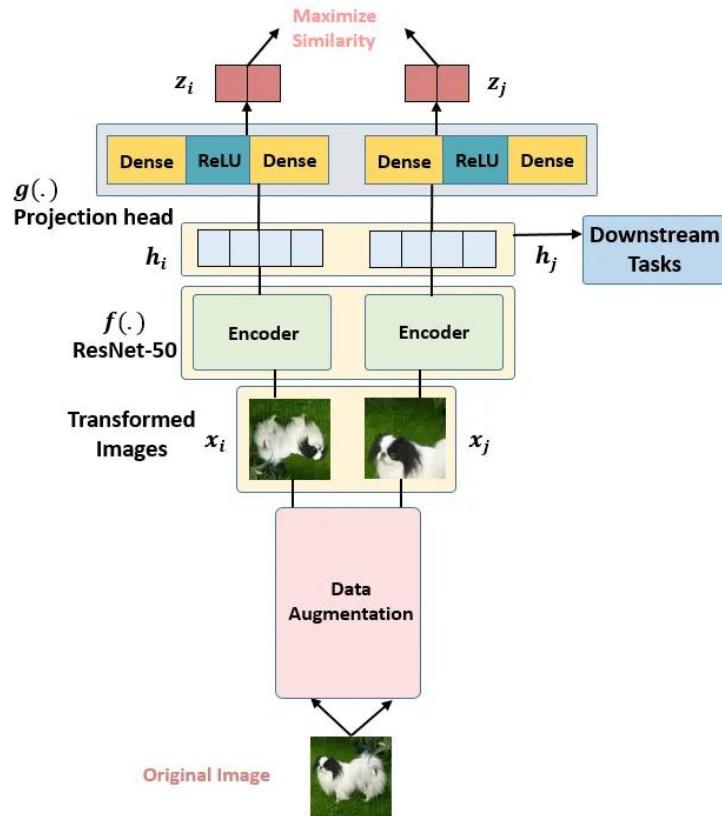


除此之外，SimCLR提出了一些可以提高性能的技术，例如把原来的一层神经网络变成了两层，效果提高很多。另外加强了Data Augmentation，训练更多的epochs。再加上大量TPU资源的加持，SimCLR跑出了当时最好的效果。如此技术性的工作，发表在了ICML2020上。

自监督学习

➤ 对比学习

SimCLR



Algorithm 1 SimCLR's main learning algorithm.

```
input: batch size  $N$ , constant  $\tau$ , structure of  $f, g, \mathcal{T}$ .
for sampled minibatch  $\{x_k\}_{k=1}^N$  do
    for all  $k \in \{1, \dots, N\}$  do
        draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$ 
        # the first augmentation
         $\tilde{x}_{2k-1} = t(x_k)$ 
         $h_{2k-1} = f(\tilde{x}_{2k-1})$                                 # representation
         $z_{2k-1} = g(h_{2k-1})$                                 # projection
        # the second augmentation
         $\tilde{x}_{2k} = t'(x_k)$ 
         $h_{2k} = f(\tilde{x}_{2k})$                                 # representation
         $z_{2k} = g(h_{2k})$                                 # projection
    end for
    for all  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  do
         $s_{i,j} = z_i^\top z_j / (\|z_i\| \|z_j\|)$       # pairwise similarity
    end for
    define  $\ell(i, j)$  as  $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$ 
     $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$ 
    update networks  $f$  and  $g$  to minimize  $\mathcal{L}$ 
end for
return encoder network  $f(\cdot)$ , and throw away  $g(\cdot)$ 
```

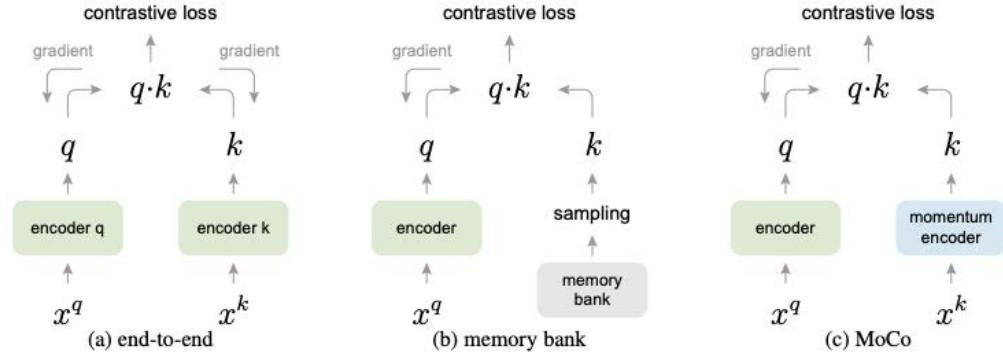
知乎 @nachifur

自监督学习

➤ 对比学习

MoCo

$$\theta_k = m \cdot \theta_{k-1} + (1 - m) \cdot \theta_q$$



Memory Bank中**保存的过时数据**确实是一个亟待解决的问题。MoCo也是为了解决这一问题。相比于SimCLR不讲武德的解决方式，MoCo更有技巧性，也更加亲民，4卡机器也可以来训练，并可以达到和SimCLR旗鼓相当的准确率。

具体地，MoCo**没有放弃**Memory Bank的思想。但是相比于Memory Bank中一个epoch更新一次的更新频率，MoCo可以让存储的样本**更加“新鲜”**。

MoCo采用了两个网络，一个网络是训练网络，另一个网络的参数是第一个网络参数的EMA（Exponential Moving Average）。用该EMA网络来产生样本存储下来作为负样本，并把太久远的样本抛弃，这样，更新样本的频率从epoch级别变到了iteration级别。

自监督学习

➤ 对比学习

MoCo

method	architecture	#params (M)	accuracy (%)
Exemplar [15]	R50w3×	211	46.0 [36]
RelativePosition [11]	R50w2×	94	51.4 [36]
Jigsaw [43]	R50w2×	94	44.6 [36]
Rotation [17]	Rv50w4×	86	55.4 [36]
Colorization [62]	R101*	28	39.6 [12]
DeepCluster [3]	VGG [51]	15	48.4 [4]
BigBiGAN [14]	R50	24	56.6
	Rv50w4×	86	61.3
<i>methods based on contrastive learning follow:</i>			
InstDisc [59]	R50	24	54.0
LocalAgg [64]	R50	24	58.8
CPC v1 [44]	R101*	28	48.7
CPC v2 [33]	R170* _{wider}	303	65.9
CMC [54]	R50 _{L+ab}	47	64.1 [†]
	R50w2× _{L+ab}	188	68.4 [†]
AMDIM [2]	AMDIM _{small}	194	63.5 [†]
	AMDIM _{large}	626	68.1 [†]
MoCo	R50	24	60.6
	RX50	46	63.9
	R50w2×	94	65.4
	R50w4×	375	68.6

Algorithm 1 Pseudocode of MoCo in a PyTorch-like style.

```
# f_q, f_k: encoder networks for query and key
# queue: dictionary as a queue of K keys (CxK)
# m: momentum
# t: temperature

f_k.params = f_q.params # initialize
for x in loader: # load a minibatch x with N samples
    x_q = aug(x) # a randomly augmented version
    x_k = aug(x) # another randomly augmented version

    q = f_q.forward(x_q) # queries: Nx1
    k = f_k.forward(x_k) # keys: Nx1
    k = k.detach() # no gradient to keys

    # positive logits: Nx1
    l_pos = bmm(q.view(N,1,C), k.view(N,C,1))

    # negative logits: NxK
    l_neg = mm(q.view(N,C), queue.view(C,K))

    # logits: Nx(1+K)
    logits = cat([l_pos, l_neg], dim=1)

    # contrastive loss, Eqn.(1)
    labels = zeros(N) # positives are the 0-th
    loss = CrossEntropyLoss(logits/t, labels)

    # SGD update: query network
    loss.backward()
    update(f_q.params)

    # momentum update: key network
    f_k.params = m*f_k.params+(1-m)*f_q.params

    # update dictionary
    enqueue(queue, k) # enqueue the current minibatch
    dequeue(queue) # dequeue the earliest minibatch
```

bmm: batch matrix multiplication; mm: matrix multiplication; cat: concatenation.

可以看出，相比于Memory Bank的解决方案，MoCo已经有了很大的改善。

伪代码解释了一切

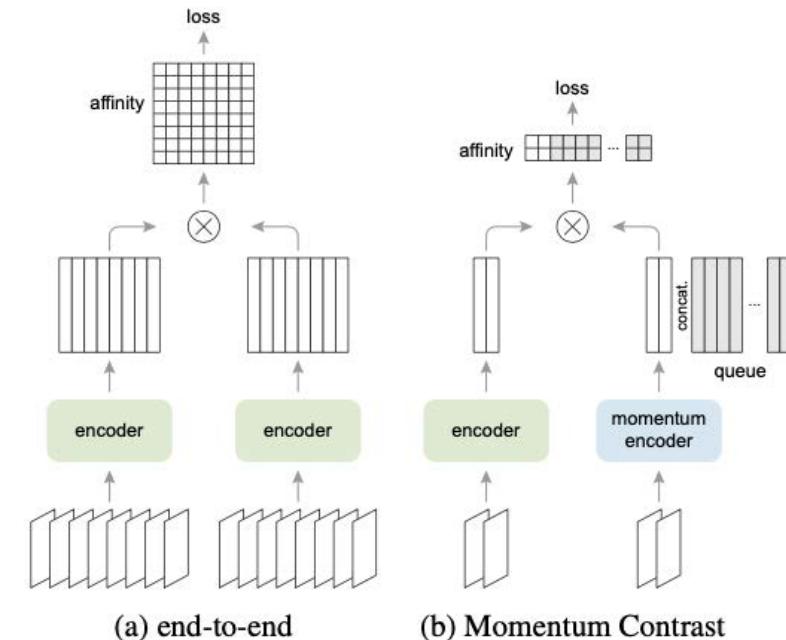
自监督学习

➤ 对比学习

MoCo v2

与MoCo的算法无异，但是把SimCLR提出的技巧性工作都加上了。包括，最后的单层神经网络头部变成了**两层MLP**，增加了Augmentation的强度，改变了温度系数超参数的选取（0.07到0.2），训练更多的epochs。结果好了很多（从60.6到71.1）。

case	unsup. pre-train				ImageNet acc.	VOC detection		
	MLP	aug+	cos	epochs		AP ₅₀	AP	AP ₇₅
supervised					76.5	81.3	53.5	58.8
MoCo v1				200	60.6	81.5	55.9	62.6
(a)	✓			200	66.2	82.0	56.4	62.6
(b)		✓		200	63.4	82.2	56.8	63.2
(c)	✓	✓		200	67.3	82.5	57.2	63.9
(d)	✓	✓	✓	200	67.5	82.4	57.0	63.6
(e)	✓	✓	✓	800	71.1	82.5	57.4	64.0



自监督学习

➤ 对比学习

理解性工作

标准的contrastive loss

$$\begin{aligned} \mathcal{L}_{\text{contrastive}}(f; \tau, M) &\triangleq \\ &\mathbb{E}_{\substack{(x,y) \sim p_{\text{pos}} \\ \{x_i^-\}_{i=1}^M \stackrel{\text{i.i.d.}}{\sim} p_{\text{data}}}} \left[-\log \frac{e^{f(x)^\top f(y)/\tau}}{e^{f(x)^\top f(y)/\tau} + \sum_i e^{f(x_i^-)^\top f(y)/\tau}} \right], \end{aligned} \quad (1)$$

简单的恒等变换

$$\begin{aligned} \mathcal{L}_{\text{contrastive}}(f; \tau, M) &= \mathbb{E}_{(x,y) \sim p_{\text{pos}}} [-f(x)^\top f(y)/\tau] \\ &+ \mathbb{E}_{\substack{(x,y) \sim p_{\text{pos}} \\ \{x_i^-\}_{i=1}^M \stackrel{\text{i.i.d.}}{\sim} p_{\text{data}}}} \left[\log \left(e^{f(x)^\top f(y)/\tau} + \sum_i e^{f(x_i^-)^\top f(x)/\tau} \right) \right]. \end{aligned}$$

自监督学习

➤ 对比学习

理解性工作

$$\begin{aligned}\mathcal{L}_{\text{contrastive}}(f; \tau, M) = & \mathbb{E}_{(x,y) \sim p_{\text{pos}}} \left[-f(x)^T f(y)/\tau \right] \\ & + \mathbb{E}_{\substack{(x,y) \sim p_{\text{pos}} \\ \{x_i^-\}_{i=1}^M \overset{\text{i.i.d.}}{\sim} p_{\text{data}}}} \left[\log \left(e^{f(x)^T f(y)/\tau} + \sum_i e^{f(x_i^-)^T f(x)/\tau} \right) \right].\end{aligned}$$

Contrastive Learning重要的两个性质：

- **对齐性 (Alignment)**
- **均匀性 (Uniformity)**。

即表征需要有数据增广的鲁棒性，又要分布均匀。

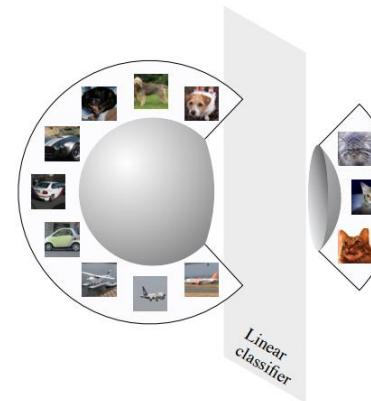
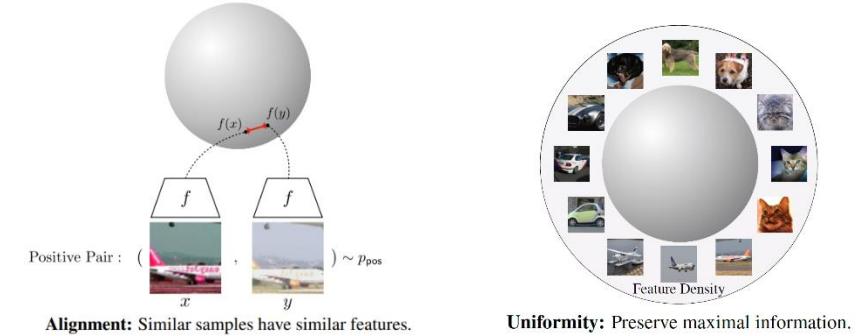


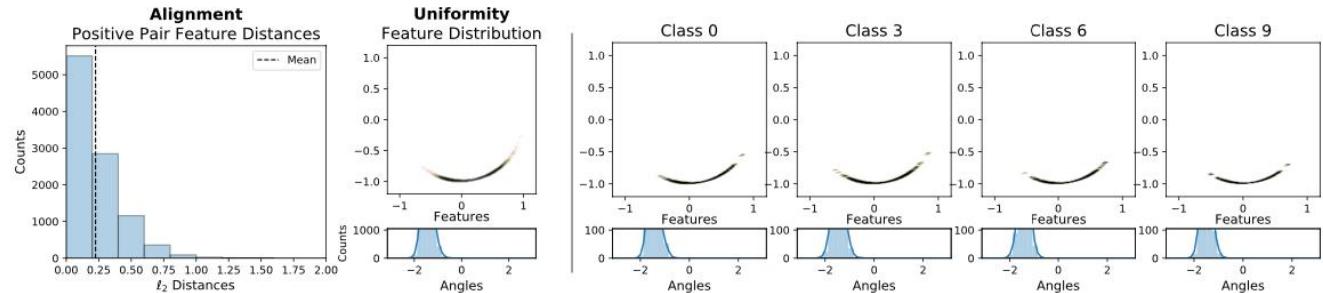
Figure 2: **Hypersphere**: When classes are well-clustered (forming spherical caps), they are linearly separable. The same does not hold for Euclidean spaces.

自监督学习

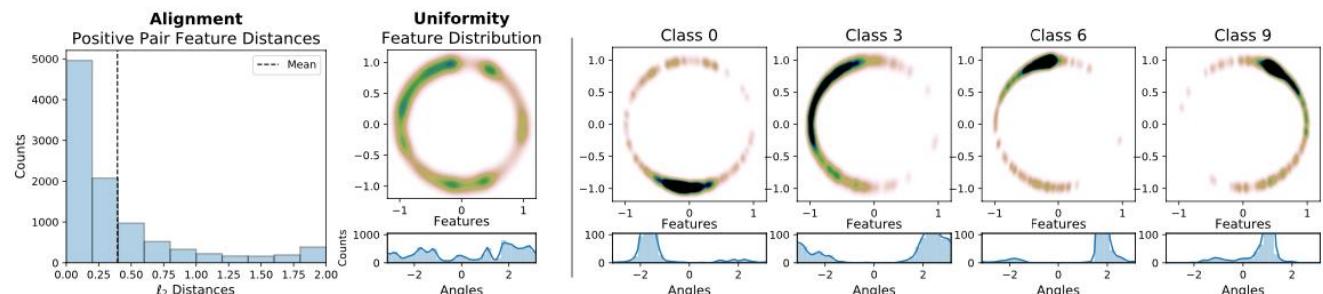
➤ 对比学习

理解性工作

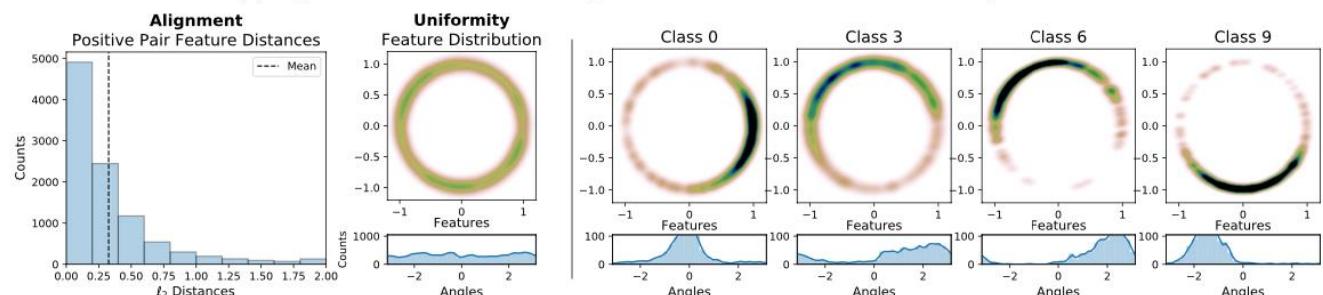
实验验证：总的来说，
uniformity十分重要，这也是为
什么需要很多的负样本。



(a) Random Initialization. Linear classification validation accuracy: 12.71%.



(b) Supervised Predictive Learning. Linear classification validation accuracy: 57.19%.



(c) Unsupervised Contrastive Learning. Linear classification validation accuracy: 28.60%.

自监督学习

➤ 对比学习

理解性工作

Contrastive Learning中的温度系数是比较神奇的参数，大部分论文将之设置为0.07-0.2这个范围，本论文揭示了温度系数的作用，以及Contrastive Loss为什么有效。

$$\begin{aligned} & \lim_{\tau \rightarrow 0^+} -\log \left[\frac{\exp(s_{i,i}/\tau)}{\sum_{k \neq i} \exp(s_{i,k}/\tau) + \exp(s_{i,i}/\tau)} \right] \\ &= \lim_{\tau \rightarrow 0^+} +\log \left[1 + \sum_{k \neq i} \exp((s_{i,k} - s_{i,i})/\tau) \right] \\ &= \lim_{\tau \rightarrow 0^+} +\log \left[1 + \sum_{\substack{k \\ s_{i,k} \geq s_{i,i}}}^k \exp((s_{i,k} - s_{i,i})/\tau) \right] \\ &= \lim_{\tau \rightarrow 0^+} \frac{1}{\tau} \max[s_{\max} - s_{i,i}, 0] \end{aligned} \quad (6)$$

$$\begin{aligned} \mathcal{L}_{\text{contrastive}}(f; \tau, M) &= \mathbb{E}_{(x,y) \sim p_{\text{pos}}} [-f(x)^T f(y)/\tau] \\ &+ \mathbb{E}_{\substack{(x,y) \sim p_{\text{pos}} \\ \{x_i^-\}_{i=1}^M \sim p_{\text{data}}}} \left[\log \left(e^{f(x)^T f(y)/\tau} + \sum_i e^{f(x_i^-)^T f(x)/\tau} \right) \right]. \end{aligned}$$

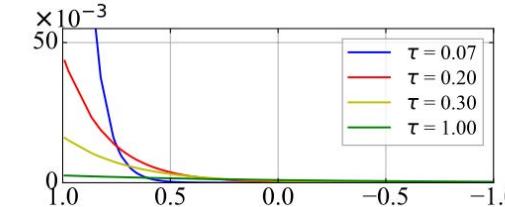


Figure 3. The gradient ratio $r_{i,j}$ with respect to different $s_{i,j}$. We sample the $s_{i,j}$ from a uniform distribution in $[-1, 1]$. As we can see, with lower temperature, the contrastive loss tends to punish more on the hard negative samples.

$$\begin{aligned} & \lim_{\tau \rightarrow +\infty} -\log \left[\frac{\exp(s_{i,i}/\tau)}{\sum_{k \neq i} \exp(s_{i,k}/\tau) + \exp(s_{i,i}/\tau)} \right] \\ &= \lim_{\tau \rightarrow +\infty} -\frac{1}{\tau} s_{i,i} + \log \sum_k \exp(s_{i,k}/\tau) \\ &= \lim_{\tau \rightarrow +\infty} -\frac{1}{\tau} s_{i,i} + \frac{1}{N} \sum_k \exp(s_{i,k}/\tau) - 1 + \log N \\ &= \lim_{\tau \rightarrow +\infty} -\frac{N-1}{N\tau} s_{i,i} + \frac{1}{N\tau} \sum_{k \neq i} s_{i,k} + \log N \end{aligned} \quad (7)$$

自监督学习

➤ 对比学习

理解性工作

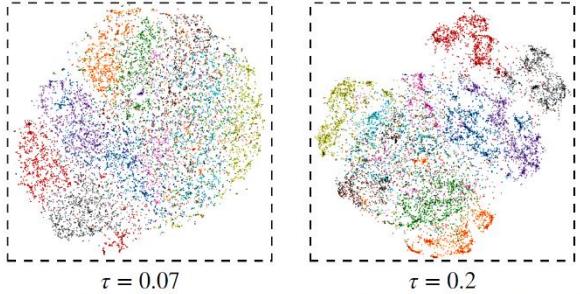


Figure 2. T-SNE [29] visualization of the embedding distribution. The two models are trained on CIFAR10. The temperature is set to 0.07 and 0.2 respectively. Small temperature tends to generate more uniform distribution and be less tolerant to similar samples.

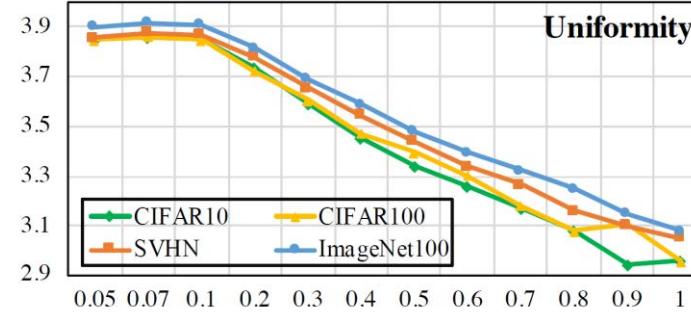


Figure 4. Uniformity of embedding distribution trained with different temperature on CIFAR10, CIFAR100 and SVHN. The x axis represents different temperature, and y axis represents $-\mathcal{L}_{\text{uniformity}}$. Large value means the distribution is more uniform.

Contrastive Loss 是一种困难样本自发现的损失，即Contrastive Loss总能将负样本中比较接近的给予更多的关注（更大的相对梯度）。

- 温度系数决定了对困难样本关注的激烈程度。越小的温度对越困难的样本给予更多的权重。

$$\begin{aligned}\mathcal{L}_{\text{contrastive}}(f; \tau, M) &= \mathbb{E}_{(x, y) \sim p_{\text{pos}}} [-f(x)^T f(y)/\tau] \\ &+ \mathbb{E}_{\substack{(x, y) \sim p_{\text{pos}} \\ \{x_i^-\}_{i=1}^M \stackrel{\text{i.i.d.}}{\sim} p_{\text{data}}}} \left[\log \left(e^{f(x)^T f(y)/\tau} + \sum_i e^{f(x_i^-)^T f(x)/\tau} \right) \right].\end{aligned}$$

$$\mathcal{L}_{\text{hard}}(x_i) = -\log \frac{\exp(s_{i,i}/\tau)}{\sum_{s_{i,k} \geq s_{\alpha}^{(i)}} \exp(s_{i,k}/\tau) + \exp(s_{i,i}/\tau)}$$

自监督学习

➤ 对比学习

理解性工作

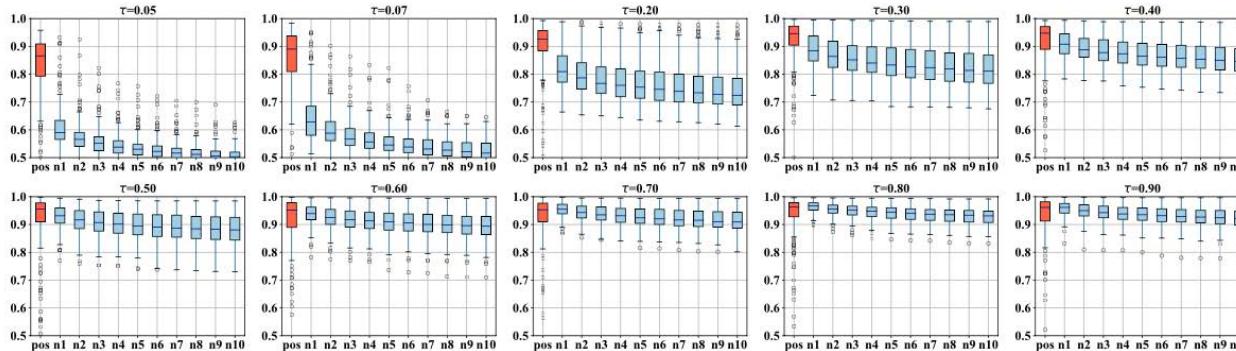


Figure 8. We display the similarity distribution of positive samples and the top-10 nearest negative samples that are marked as 'pos' and 'ni' for the i -th nearest neighbour. All models are trained on CIFAR100. For models trained on other datasets, they present the same pattern with the above figure, and we display them in the supplementary material.

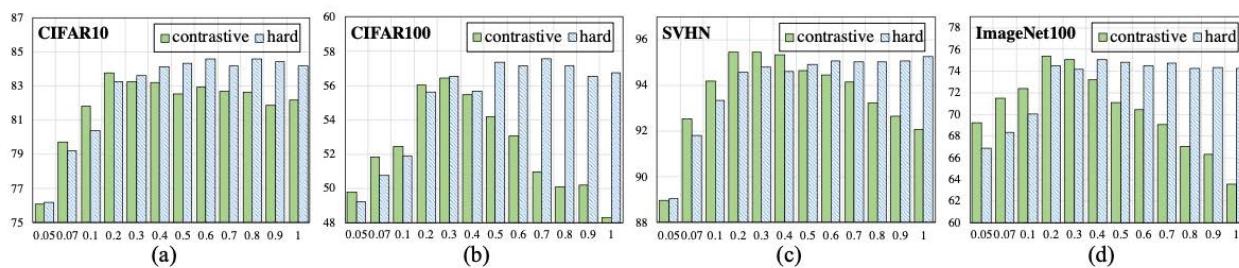


Figure 9. Performance comparison of models trained with different temperatures. For CIFAR10, CIFAR100 and SVHN, the backbone network is ResNet-18, and for ImageNet, the backbone network is ResNet-50. After the pretraining stage, we freeze all convolutional layers and add a linear layer. We report 1-crop top-1 accuracy for all models.

不同温度系数下训练的网络表现：红色盒子为正样本相似度，右侧的蓝色盒子依次是最相似负样本，第二相似负样本。可以看出，温度系数越小，gap越大

不同温度系数下网络在下游任务的表现，四个图分别在四个数据集上，绿色代表采用 contrastive learning，蓝色代表采取显式困难样本挖掘的损失。

$$\mathcal{L}_{\text{hard}}(x_i) = -\log \frac{\exp(s_{i,i}/\tau)}{\sum_{s_{i,k} \geq s_\alpha^{(i)}} \exp(s_{i,k}/\tau) + \exp(s_{i,i}/\tau)}$$

自监督学习

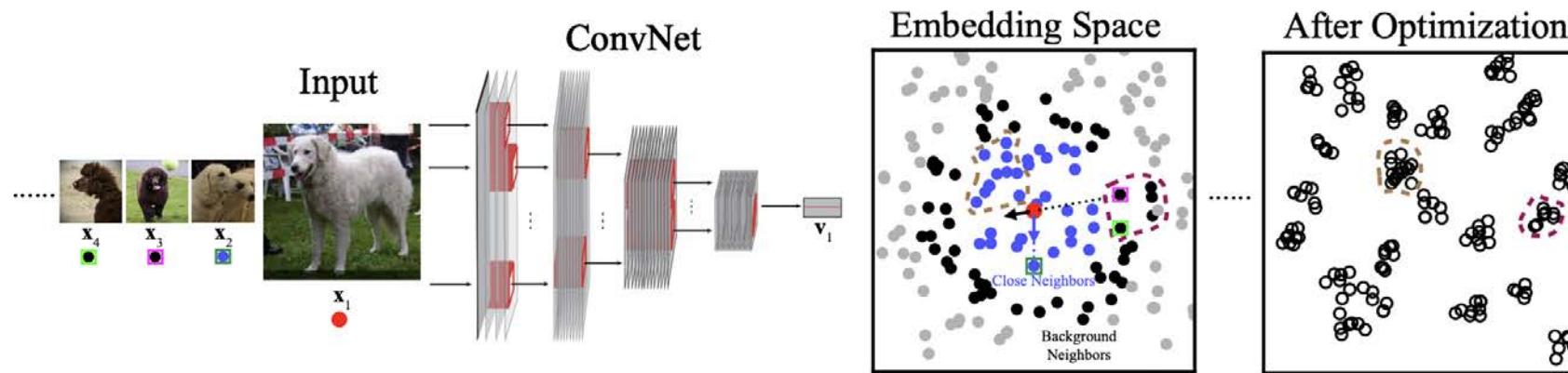


自监督学习

➤ Multiple Positive Samples

前面的方法主要聚焦于用augmentation做多视角，许多工作也聚焦于使用Nearst Neighbour来挖掘正样本。

Local aggregation for unsupervised learning of visual embeddings



使用聚类方法，寻找同一个聚类类别下的samples作为正样本

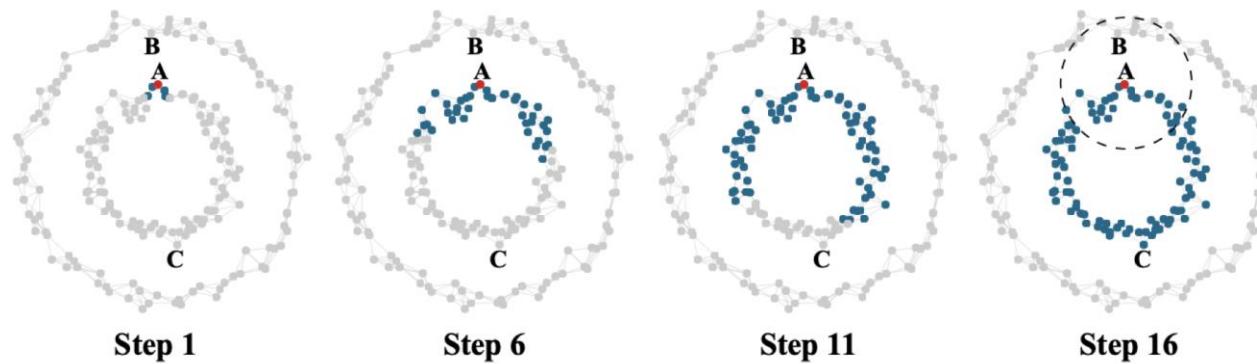
$$P(i|\mathbf{v}) = \frac{\exp(\mathbf{v}_i^T \mathbf{v} / \tau)}{\sum_{j=1}^N \exp(\mathbf{v}_j^T \mathbf{v} / \tau)} \quad (1) \quad P(\mathbf{A}|\mathbf{v}) = \sum_{i \in \mathbf{A}} P(i|\mathbf{v}) \quad (2)$$

自监督学习

➤ Multiple Positive Samples

本工作将正样本进行了扩展。上面的Contrastive Learning方法都是将所有不同的样本都作为负样本。但是将同一个类别的物体也分开是会损害对语义信息的建立的。

采用如下图的方式来构建正样本，选择出每个样本的k近邻，并构造knn graph，并将graph距离小于K的选为正样本。



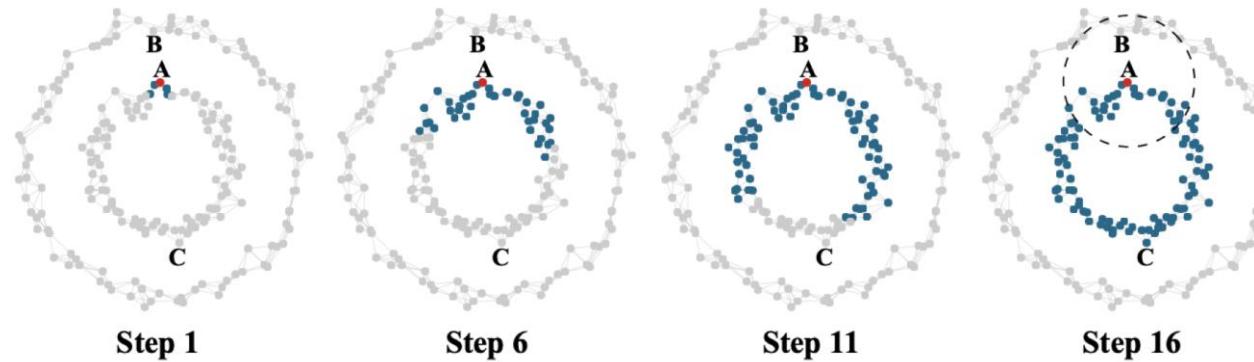
$$\mathcal{N}(i) = \mathcal{N}_k(i) \cup \mathcal{N}_k(\mathcal{N}_k(i)) \cup \dots \cup \underbrace{\mathcal{N}_k(\mathcal{N}_k(\mathcal{N}_k(\dots \mathcal{N}_k(i))))}_{l}$$

- F. Wang, H. Liu, D. Guo, et al. Unsupervised representation learning by invariance propagation[J]. Advances in Neural Information Processing Systems, 2020, 33: 3510-3520.

自监督学习

➤ Multiple Positive Samples

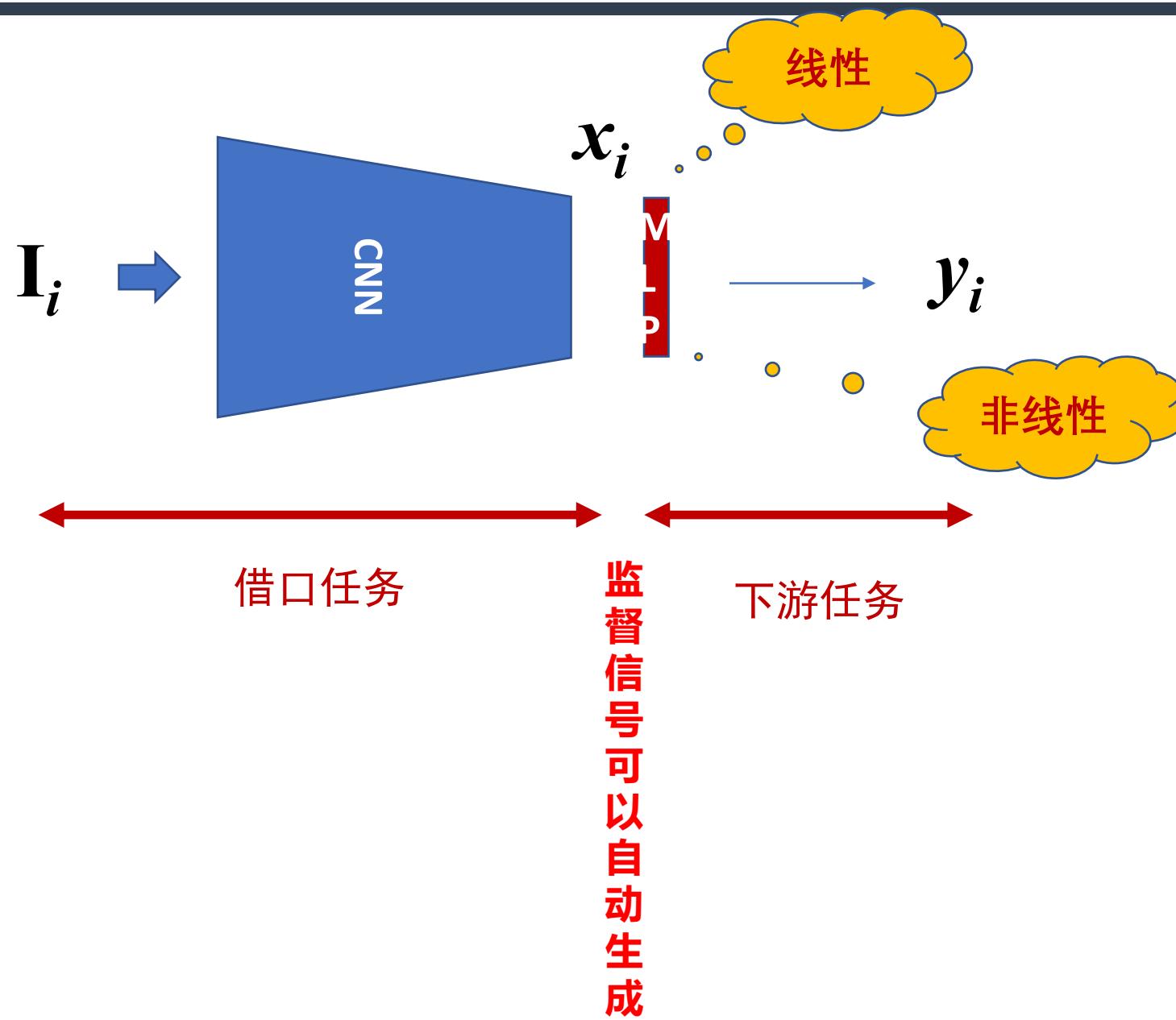
这样选的原因是基于这样一个假设：每个k近邻是语义相似的，理想的embedding空间应该是光滑的。基于传递性，可以将所有k近邻的k近邻也纳入其中，下图展示了这样选取的优点。另外论文也提出了一些技巧，例如选取困难负样本与困难正样本等。总体效果也不错。



CPC v2 [19]	ResNet-170	303	-	65.9	-	-
AMDIM [1]	AMDIM	626	150	68.1	55.0	-
SimCLR [5]	ResNet-50-MLP	28	1000	69.3	-	80.5
MoCo v2 [6]	ResNet-50-MLP	28	800	71.1	-	-
PCL [26]	ResNet-50	24	200	62.2	49.2	82.2
PCL [26]	ResNet-50-MLP	28	200	65.9	49.8	84.0
InvP (Ours)	ResNet-50	24	800	67.7	52.6	84.2
InvP (Ours)	ResNet-50-MLP	28	800	71.3	53.5	84.7

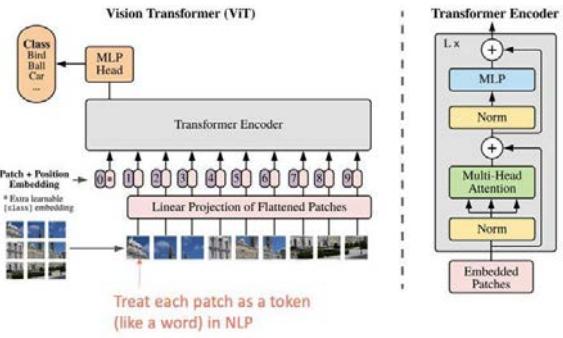
- F. Wang, H. Liu, D. Guo, et al. Unsupervised representation learning by invariance propagation[J]. Advances in Neural Information Processing Systems, 2020, 33: 3510-3520.

自监督学习



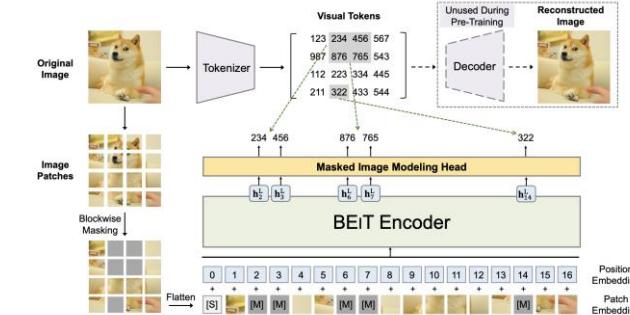
自监督学习

Vision Transformer

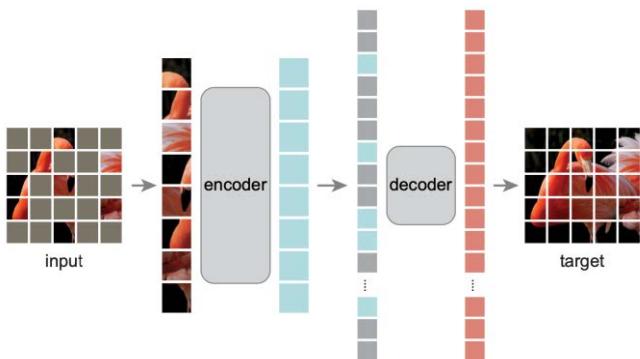


与CNN不同，Vision Transformer使用Attention模块和MLP模块构建网络，将图像划分patch后作为输入。ViT没有池化操作，每层都是相同的特征，与NLP所使用的Transformer结构几乎相同。

Beit: BERT Pre-Training of Image Transformers



Beit的思想是，将图像划分后的patch，mask掉一部分，然后还原这部分，但还原的并不是被mask掉的像素，而是还原被mask的patch的一个离散序号。即首先使用dVAE将每个patch找到一个对应的id，并用此id作为类别来学习。dVAE的作用可以看作给patch聚类得到的一个类别序号。



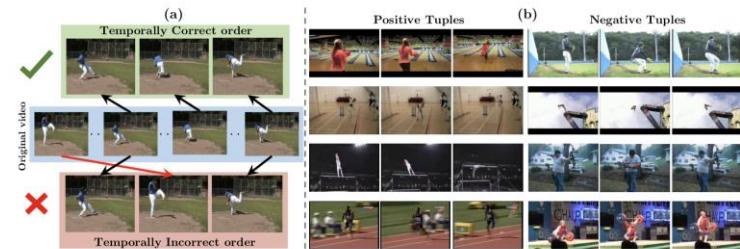
Masked Autoencoders Are Scalable Vision Learners

即便mask掉95%的数据，依然可以还原的不错，这说明图像信息本就是冗余的，通过少量的数据来补全剩下的数据可以学到相关的有用信息

method	pre-train data	ViT-B	ViT-L	ViT-H	ViT-H ₄₄₈
scratch, our impl.	-	82.3	82.6	83.1	-
DINO [5]	IN1K	82.8	-	-	-
MoCo v3 [9]	IN1K	83.2	84.1	-	-
BEiT [2]	IN1K+DALLE	83.2	85.2	-	-
MAE	IN1K	83.6	85.9	86.9	87.8

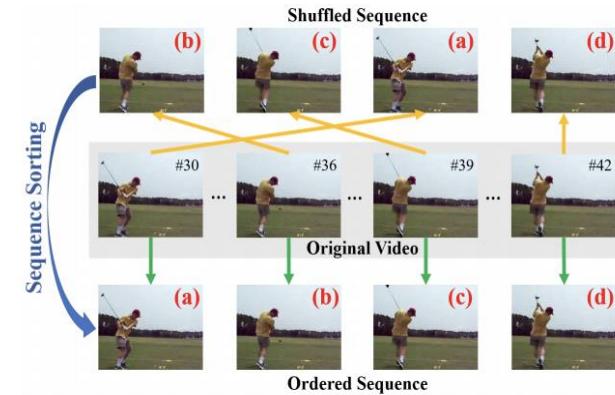
自监督学习

➤ 视频自监督学习 Shuffle and Learn

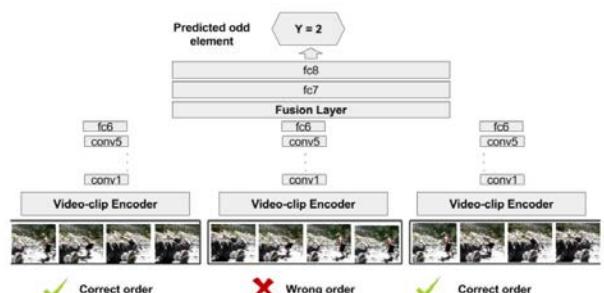


正负样本的实例，正样本为顺序的视频，负样本为乱序的视频

Sequence Order Prediction

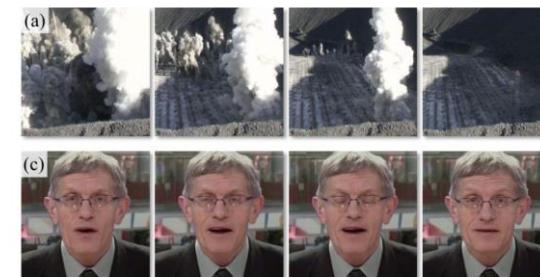


Odd-One-Out Networks



与shuffle and learn类似，都是判断一个视频片段是否是正确的顺序。但是网络结构的设计有所不同，本文采用的是多个样本进行判断，判断哪一个样本是错误打乱的。

Arrow of Time

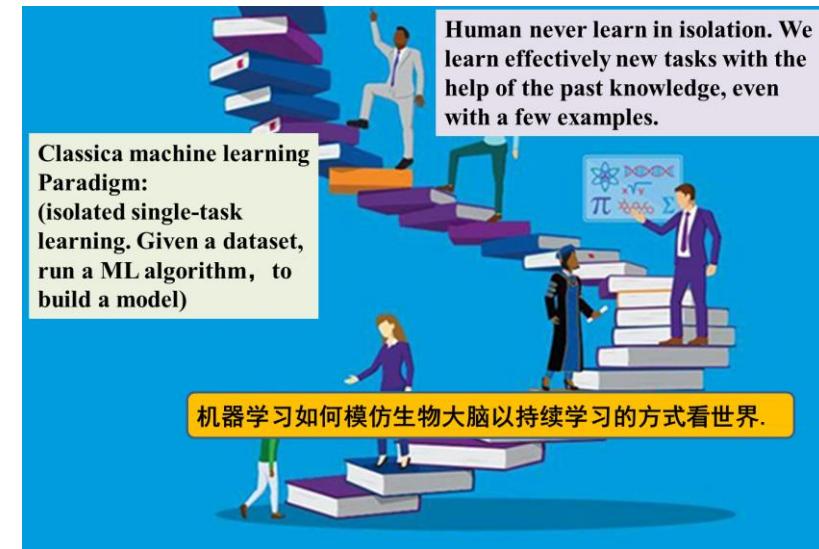
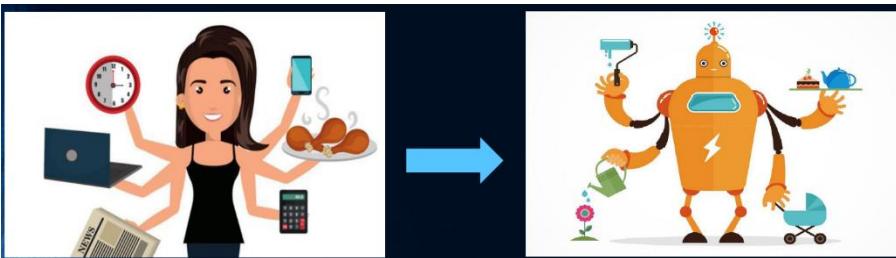


通过选择合理时间的顺序来得到一些较为抽象的信息（热力学第二定律）

AI Today: Impressive... but (Still) "Narrow"



Human-Level AI: "Broad" - Versatile, Multi-Task



持续学习

持续学习

➤ 背景

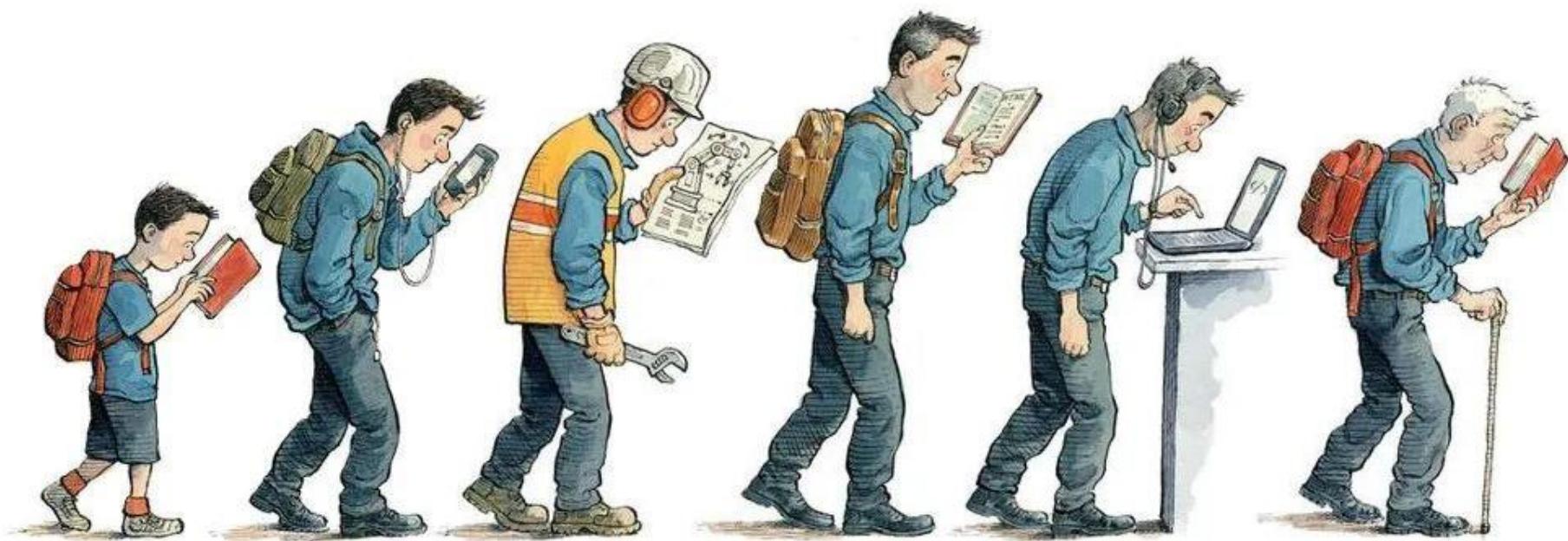


聊天机器人上线24小时被教坏

持续学习

➤ 背景

- Lifelong, Continual Learning



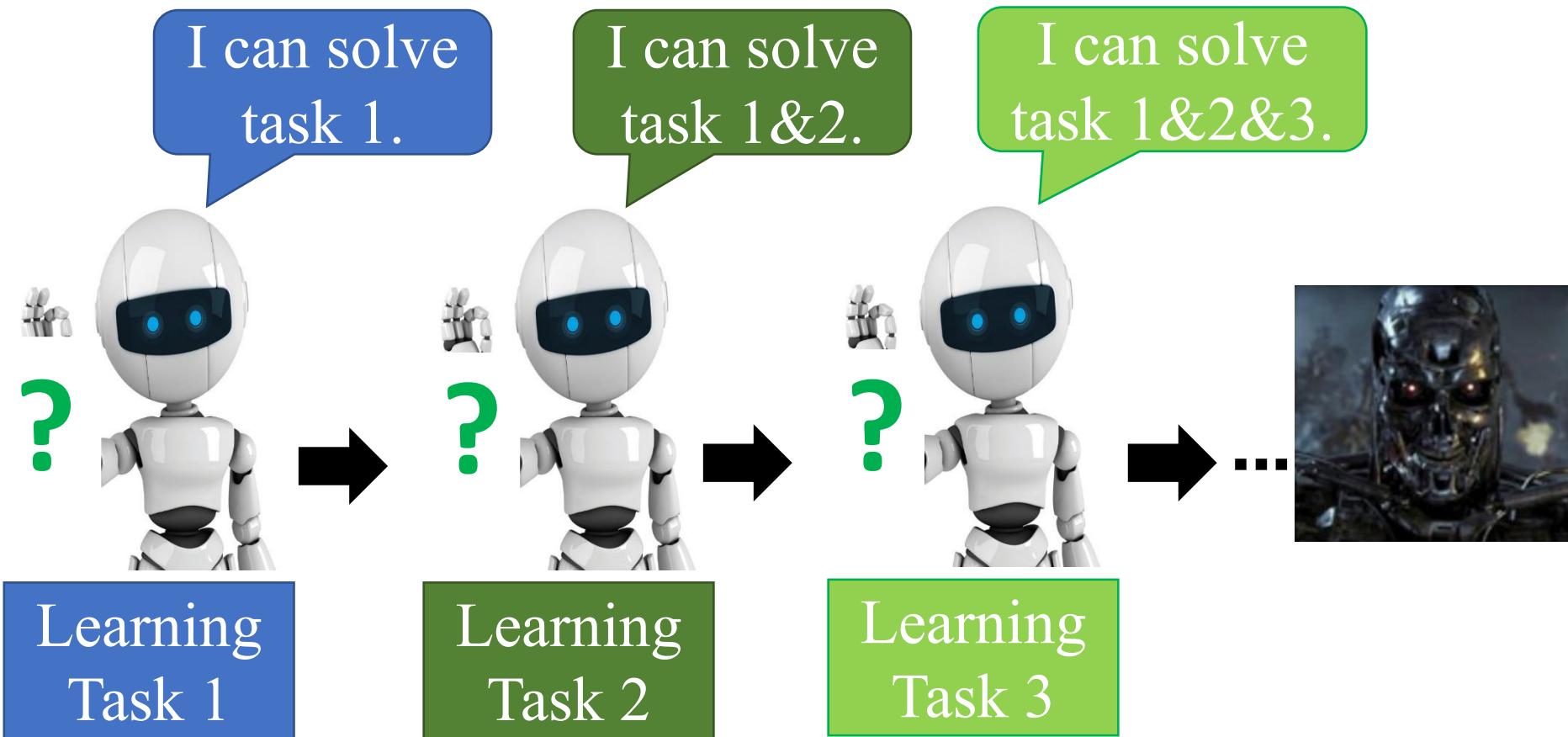
“Continual learning is the constant development of increasingly complex behaviors; the process of building more complicated skills on top of those already developed.”

持续学习

➤ 背景

Life-long Learning (LLL)

(Continuous learning, Never Ending Learning, Incremental Learning)



持续学习

➤ 背景

Lifelong Robot Learning¹

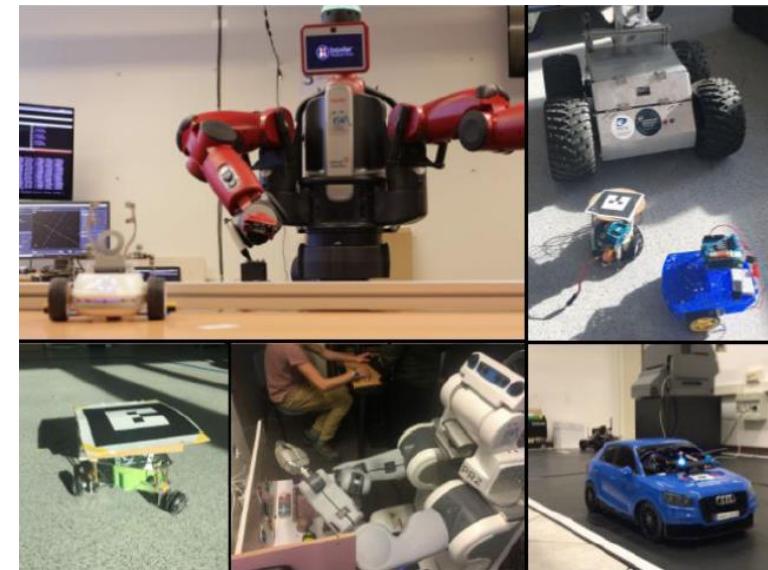
Sebastian Thrun² and Tom M. Mitchell³

² University of Bonn, Institut für Informatik III, Römerstr. 164, 53117 Bonn, Germany
³ School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA

Abstract. Learning provides a useful tool for the automatic design of autonomous robots. Recent research on learning robot control has predominantly focussed on learning single tasks that were studied in isolation. If robots encounter a multitude of control learning tasks over their entire lifetime, however, there is an opportunity to transfer knowledge between them. In order to do so, robots may learn the invariants of the individual tasks and environments. This task-independent knowledge can be employed to bias generalization when learning control, which reduces the need for real-world experimentation. We argue that knowledge transfer is essential if robots are to learn control with moderate learning times in complex scenarios. Two approaches to lifelong robot learning which both capture invariant knowledge about the robot and its environments are presented. Both approaches have been evaluated using a HERO-2000 mobile robot. Learning tasks included navigation in unknown indoor environments and a simple find-and-fetch task.



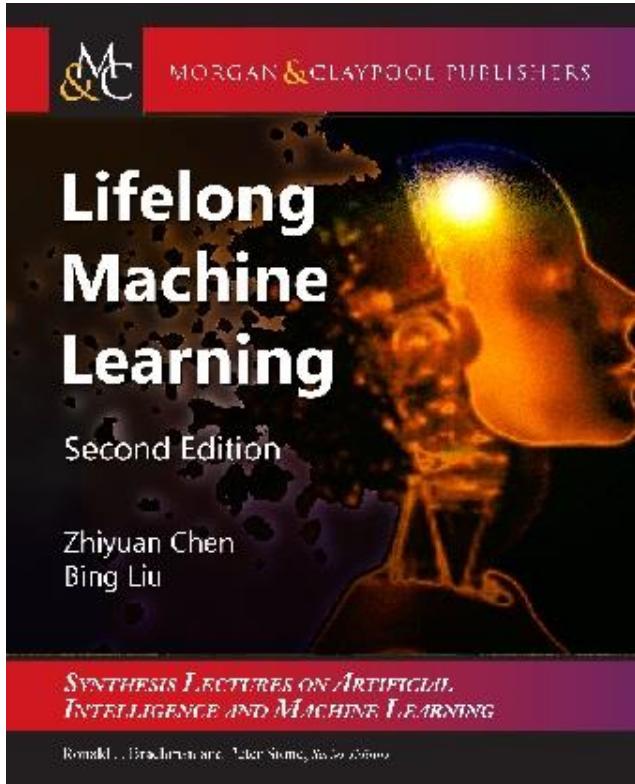
- S. Thrun and T. Mitchell. Lifelong robot learning. *Robotics and Autonomous Systems*, 15:25-46, 1995.



- Lesort, T., Lomonaco, V., Stoian, A., Maltoni, D., Filliat, D. and Díaz-Rodríguez, N., Continual learning for robotics: Definition, framework, learning strategies, opportunities and challenges. *Information fusion*, 2020.

持续学习

➤ 背景



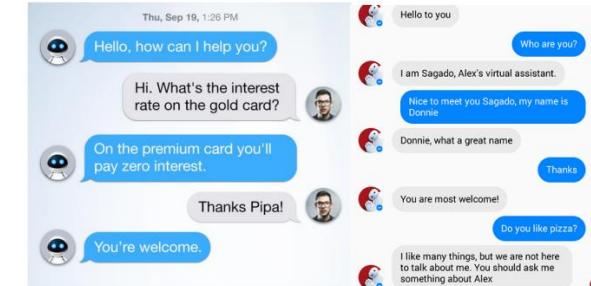
Social chats	
Turn	Dialogue
1	usr: Good morning!
2	agr: Good morning. I hope you had a great weekend.
3	usr: not really.
4	agr: You sound depressed. What can I do for you?
5	usr: I'd like to know where sales are lagging behind our forecast?
6	agr: The worst region is [country], where sales are 15% below projections.
7	usr: Do you know why?
8	agr: The forecast for [product] growth was overly optimistic.
9	usr: How can we turn this around?
10	agr: Here are the 10 customers in [country] with the most growth potential, per our CRM model.
11	usr: Can you set up a meeting with the CTO of [company]?
12	agr: Yes, I've set up a meeting with [person name] for next month when you are in [location].
13	usr: Thanks!

Question-answering	
1	usr: Good morning!
2	agr: Good morning. I hope you had a great weekend.
3	usr: not really.
4	agr: You sound depressed. What can I do for you?
5	usr: I'd like to know where sales are lagging behind our forecast?
6	agr: The worst region is [country], where sales are 15% below projections.
7	usr: Do you know why?
8	agr: The forecast for [product] growth was overly optimistic.
9	usr: How can we turn this around?
10	agr: Here are the 10 customers in [country] with the most growth potential, per our CRM model.
11	usr: Can you set up a meeting with the CTO of [company]?
12	agr: Yes, I've set up a meeting with [person name] for next month when you are in [location].
13	usr: Thanks!

Task completion	
1	usr: Good morning!
2	agr: Good morning. I hope you had a great weekend.
3	usr: not really.
4	agr: You sound depressed. What can I do for you?
5	usr: I'd like to know where sales are lagging behind our forecast?
6	agr: The worst region is [country], where sales are 15% below projections.
7	usr: Do you know why?
8	agr: The forecast for [product] growth was overly optimistic.
9	usr: How can we turn this around?
10	agr: Here are the 10 customers in [country] with the most growth potential, per our CRM model.
11	usr: Can you set up a meeting with the CTO of [company]?
12	agr: Yes, I've set up a meeting with [person name] for next month when you are in [location].
13	usr: Thanks!

Table of Contents

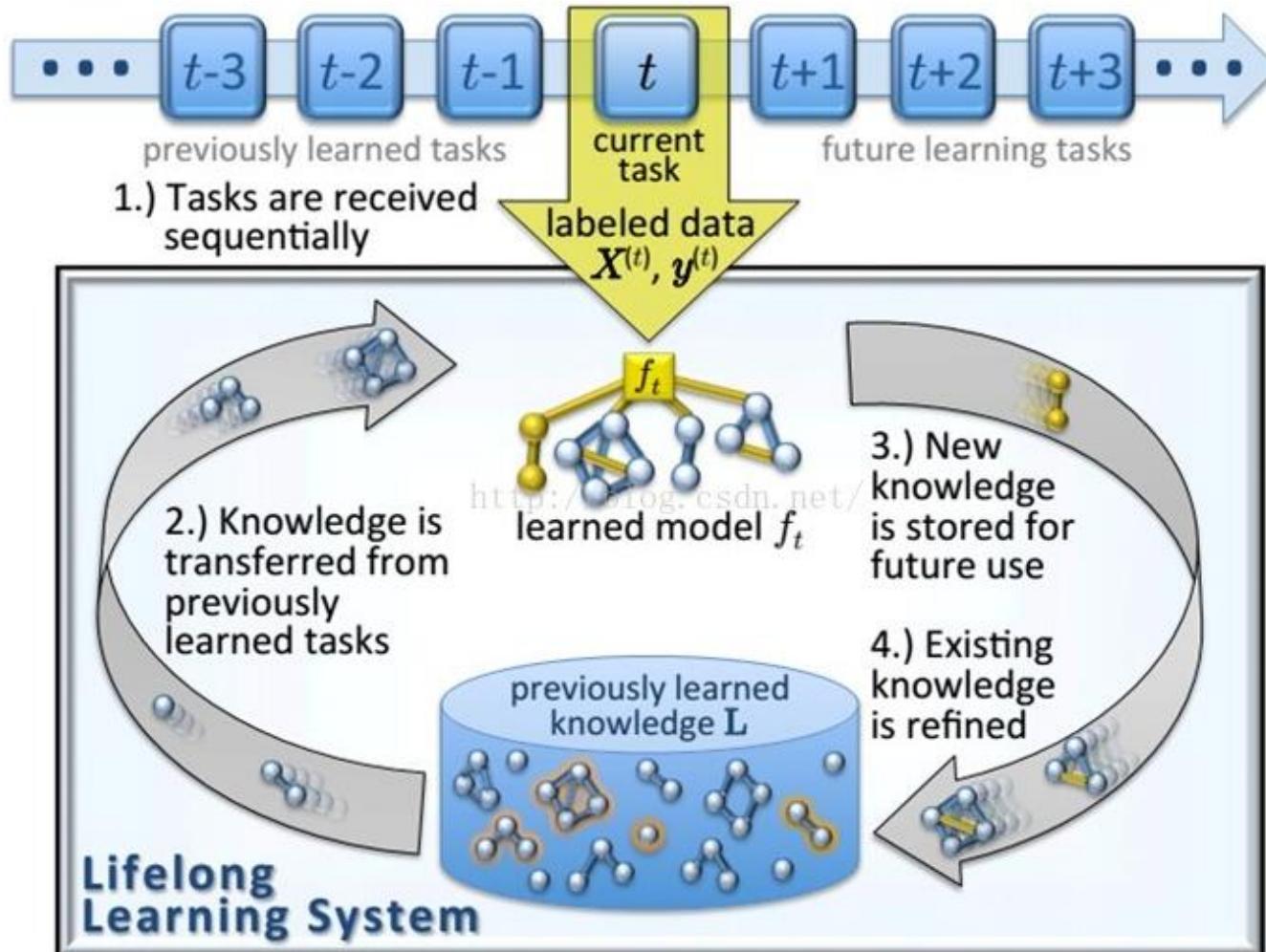
1. [Introduction](#)
2. [Related Learning Paradigms](#)
3. [Lifelong Supervised Learning](#)
4. [Continual Learning and Catastrophic Forgetting](#)
5. [Open-world Learning](#)
6. [Lifelong Topic Modeling](#)
7. [Lifelong Information Extraction](#)
8. [Continuous Knowledge Learning in Chatbots](#)
9. [Lifelong Reinforcement Learning](#)
10. [Conclusion and Future Directions](#)



持续学习

➤ 基本框架

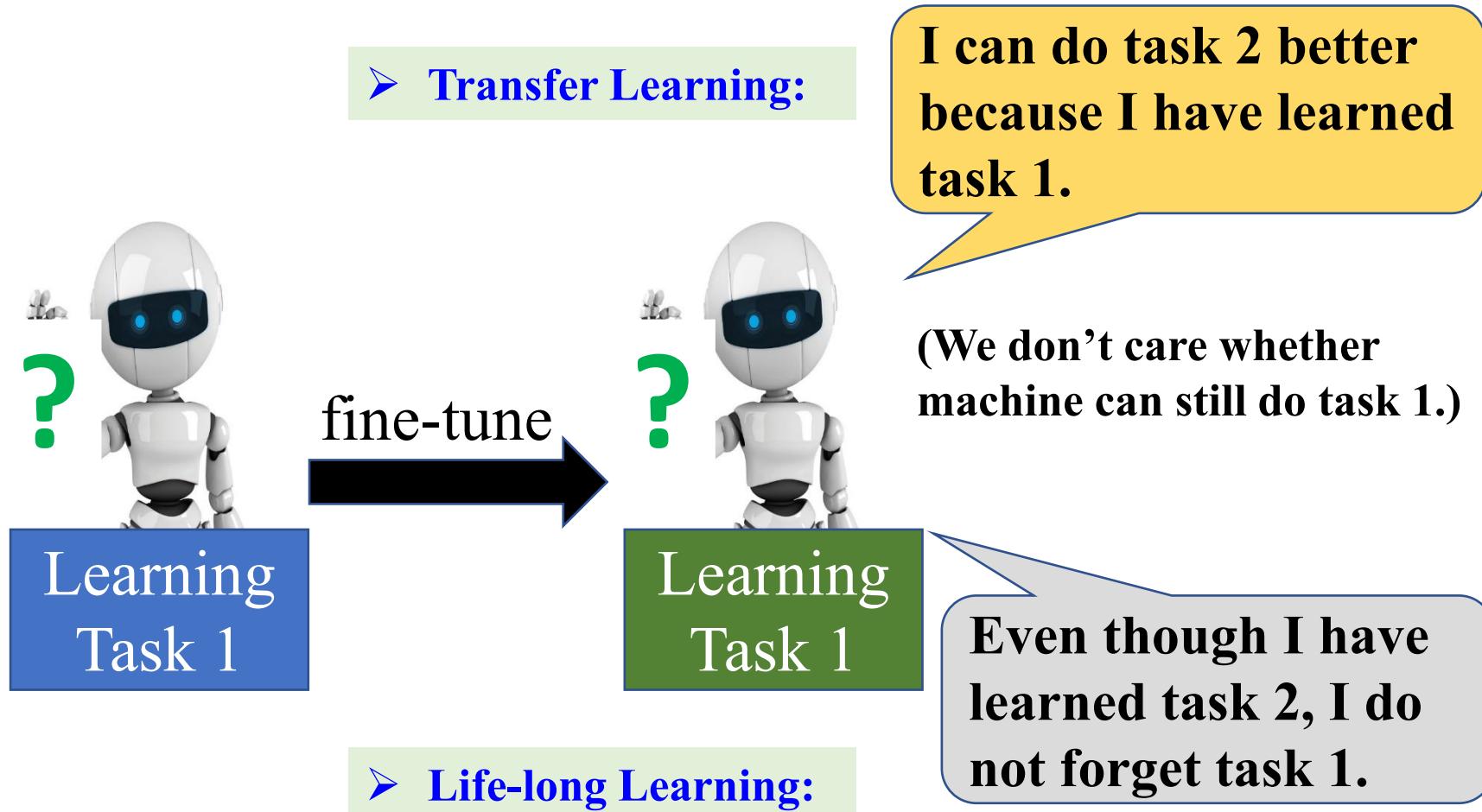
Life-long learning的基本框架



持续学习

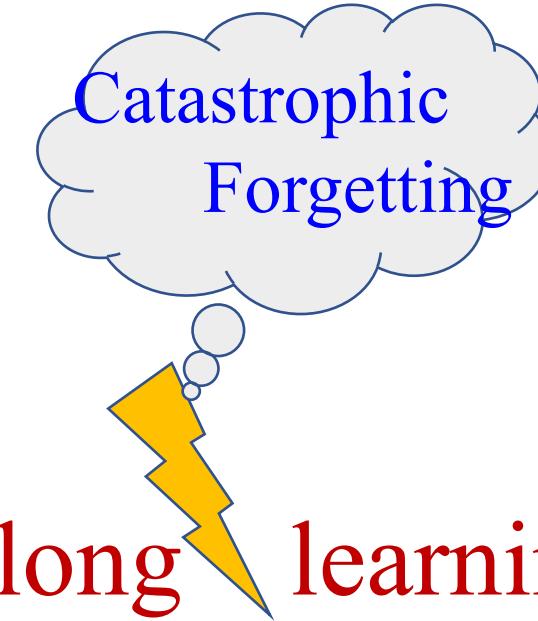
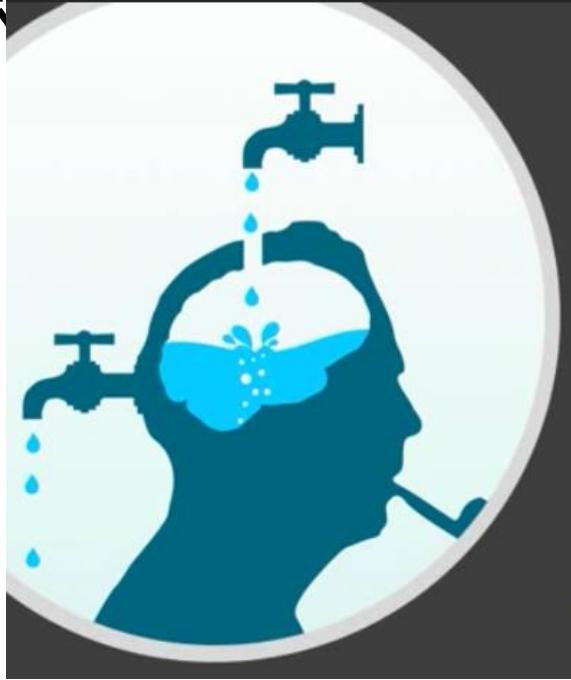
➤ 基本框架

Life-long learning V.S. Transfer Learning

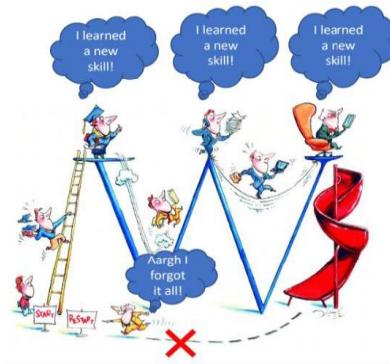


持续学习

➤ 灾难性遗忘

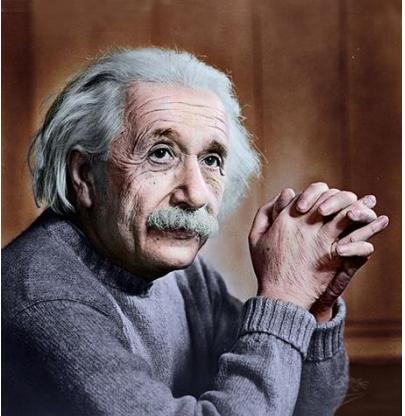


Lifelong learning



持续学习

➤ 灾难性遗忘



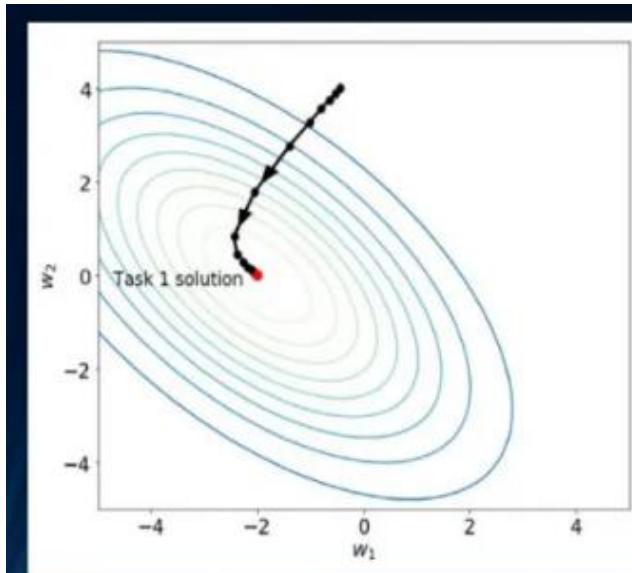
- 教育，就是忘记了在学校所学的一切之后剩下的东西

perform perform perform perform perform perform perform

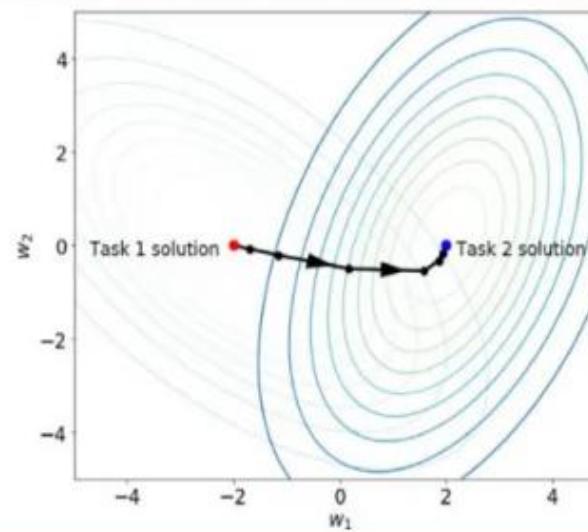


持续学习

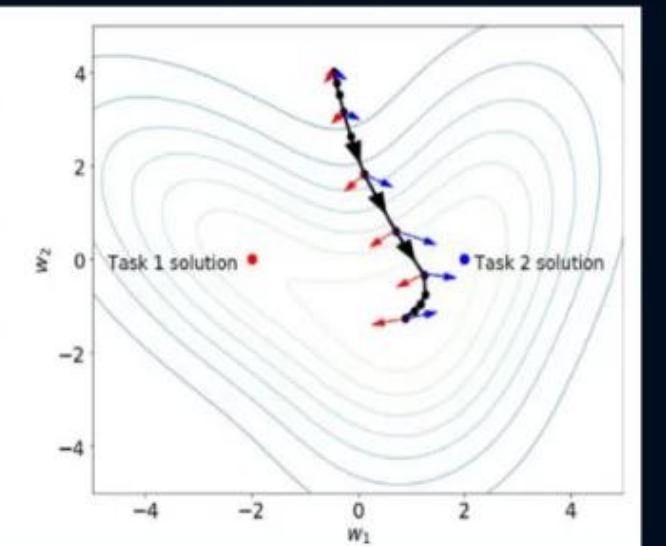
➤ 灾难性遗忘



Loss(Task1)



Loss(Task2)



Loss(Task1) + Loss(Task2)

持续学习

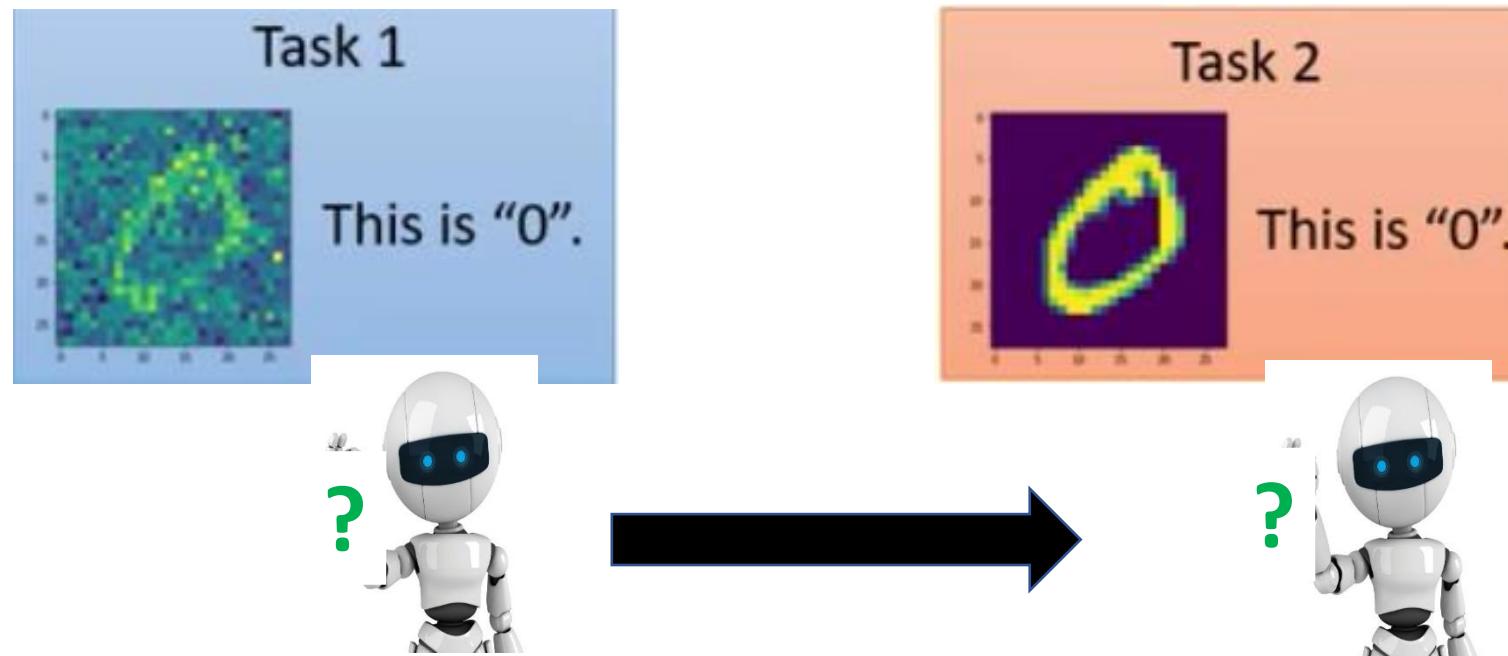
➤ 灾难性遗忘



持续学习

➤ 灾难性遗忘

现有两个数字识别任务，第一个任务的样本加入了某些随机噪声，而第二个任务的样本不做处理，如下图所示：



持续学习

➤ 灾难性遗忘

 = 3 layers, 50 neurons each

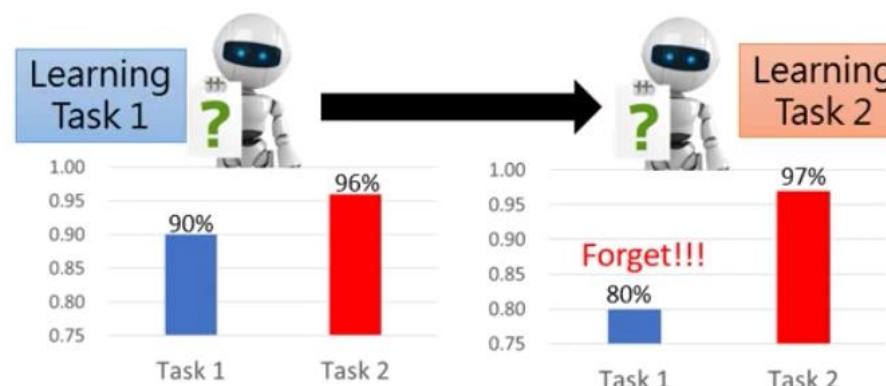


但是经过实验证明，真实情况和预想具有一定的差距。如上所示，当我们在Task 1的训练集上训练完模型后，可以取得90%的准确率，而且将其直接应用到Task 2上可以有96%的准确率。但是当我们继续使用Task 2的训练集训练模型后，Task 2上的准确率改变不大，但是Task 1上的准确率反而大幅的下降了，这显然不是我们想看到的事情。

持续学习

➤ 灾难性遗忘

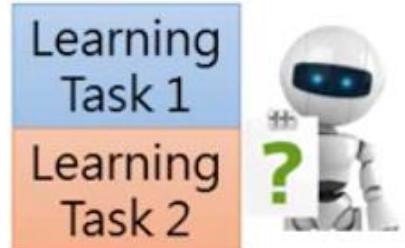
解决上述问题的一个简单的办法就是将Task 1和Task 2的样本混在一起来训练一个模型，希望它可以在两个任务上都表现良好，结果也正是这样。



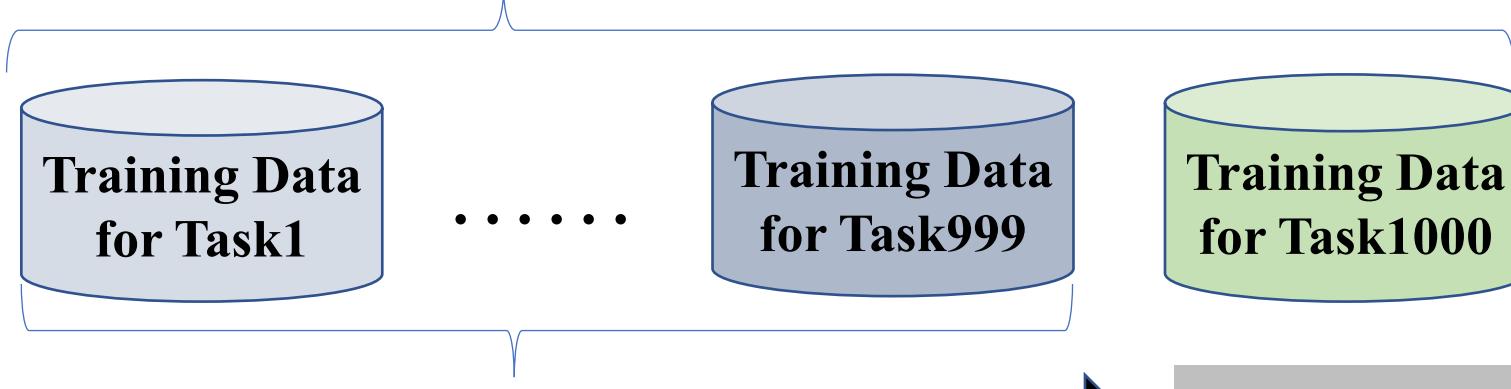
持续学习

➤ 灾难性遗忘

➤ Multi-task training can solve the problem!



Using all the data for training → Computation issue

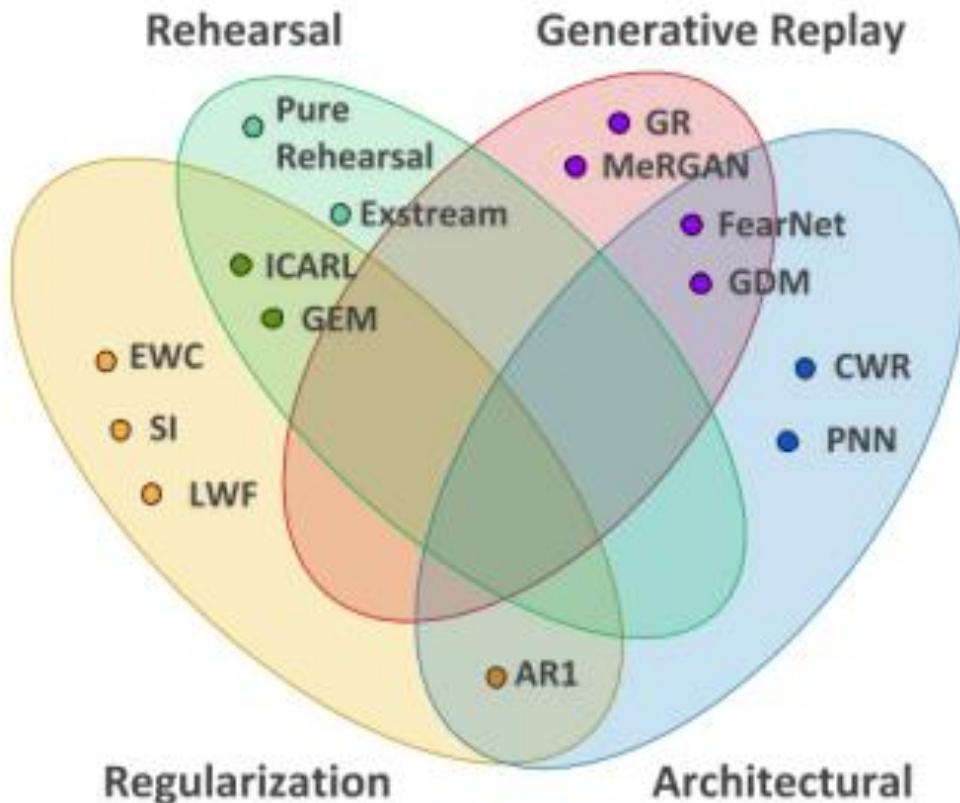


Always keep the data → Storage / data issue

✓ Multi-task training can be considered as the upper bound of LLL.

持续学习

➤ 典型方法

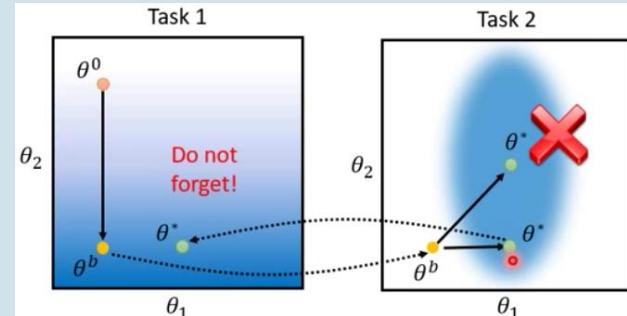


References	Regularization	Rehearsal	Architectural	Generative-Replay
Zhou et al. [178]				✓
Goodfellow et al. [53]	✓			
Lyubova et al. [94]		✓		
Rusu et al. [136]	✓			
Camoriano et al. [17]	✓	✓	✓	
Furlanello et al. [48]	✓			✓
Li et al. [87] (LwF)	✓			✓
Rusu et al. [137] (PN)				✓
Jung et al. [65]	✓			
Aljundi et al. [3]				✓
Rebuffi et al. [125] (Icarl)	✓		✓	
Kirkpatrick et al. [73] (EWC)	✓			
Fernando et al. [43]				✓
Lee et al. [80]	✓			
Lee et al. [174]	✓			
Triki et al. [161]		✓		
Seff et al. [145]		✓		
Shin [150] (DGR)				✓
Veles et al. [165]	✓			
Lopez-Paz et al. [92] (GEM)	✓		✓	
Zenke et al. [176] (SI)	✓			
Nguyen et al. [111] (VCL)	✓	✓	✓	✓
Ramapuram et al. [124]	✓			✓
Mallya et al. [96]				✓
Kamra et al. [69]				✓
Draelos et al. [37]				✓
Serra et al. [146]		✓		
Mallya et al. [95]			✓	
Parisi et al. [115] (GDM)	✓		✓	✓
He et al. [57]	✓		✓	
Hayes et al. [55]			✓	
Wu et al. [172]		✓		✓
Ritter et al. [131]	✓			
Schwarz et al. [144]			✓	
Maltoni et al. [97]	✓			✓
Achille et al. [1]			✓	
Wu et al. [171] (MeRGAN)	✓			✓
Dhar et al.	✓			
Lesort et al. [81]				✓
Caselles-Dupré et al. [20]				✓
Riemer et al. [127] (MER)	✓		✓	
Rios et al. [130] (CloGAN)	✓		✓	✓
Lesort et al. [83]				✓
Sprechmann et al. [154]		✓		✓
Kemker et al. [70] (FearNet)			✓	✓
Chaudhry et al. [23]	✓		✓	
Kalifou1 et al. [68]	✓		✓	

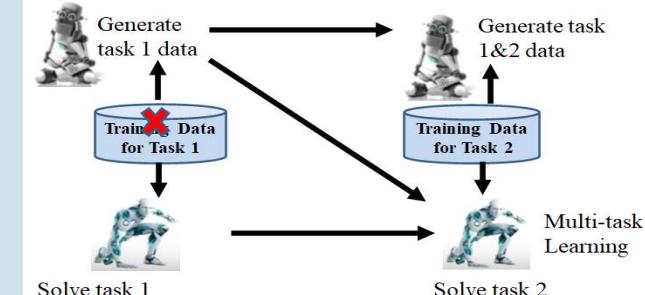
持续学习

➤ 典型方法

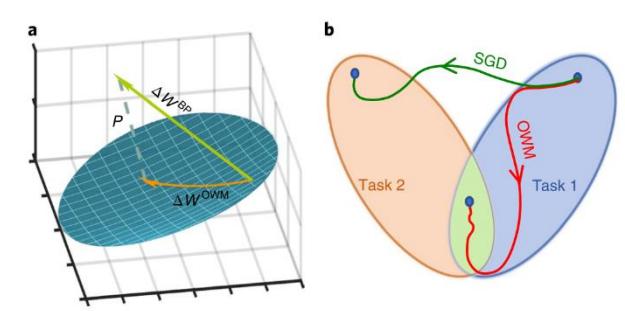
参数加权



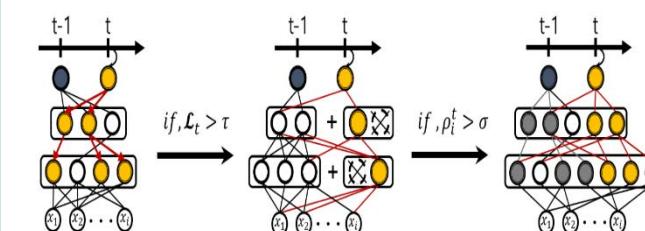
样本回放



约束梯度

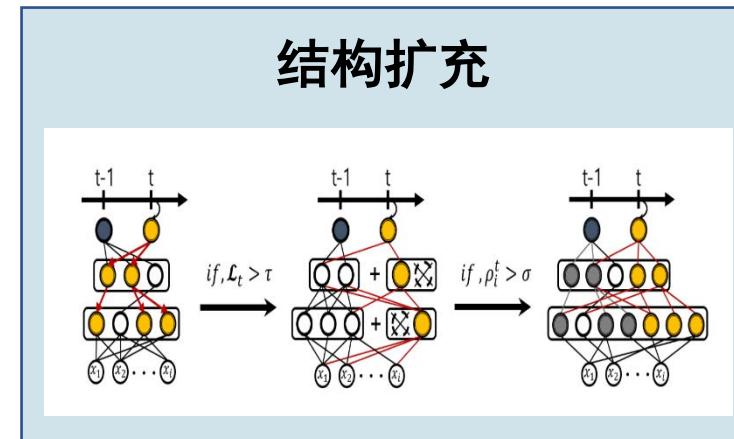
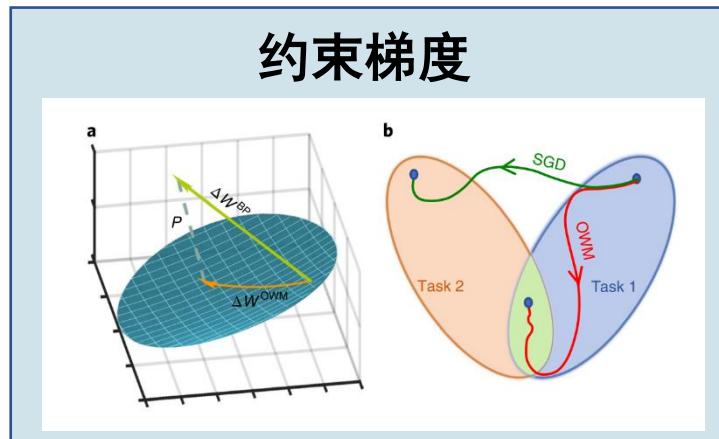
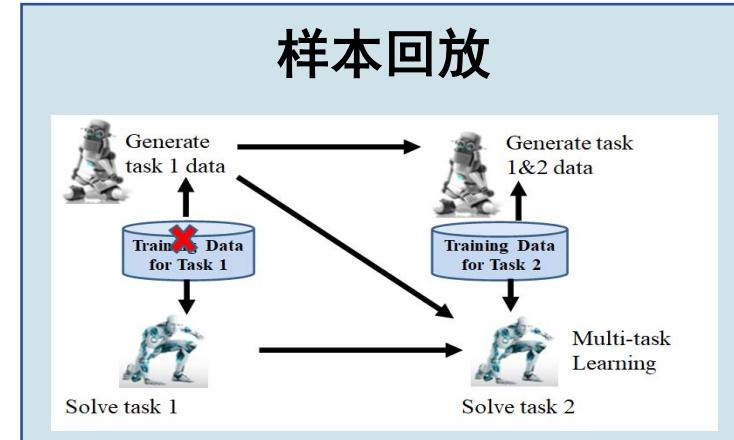
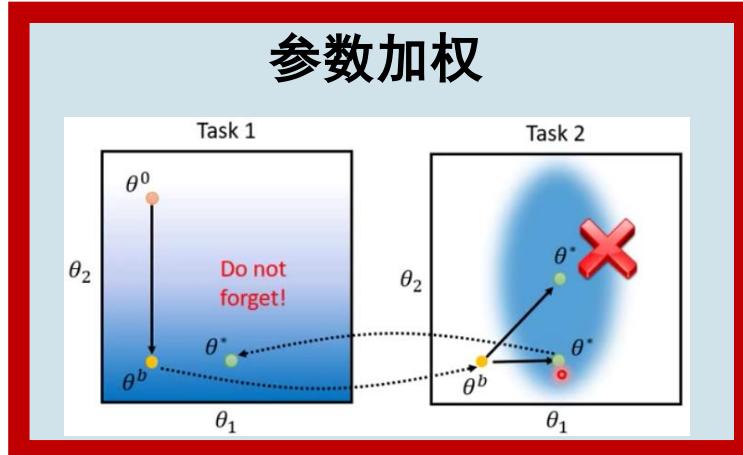


结构扩充



持续学习

➤ 典型方法



持续学习

➤ Elastic Weight Consolidation (EWC)

Overcoming catastrophic forgetting in neural networks

(PNAS-2017)

Basic idea: EWC是通过给权重添加正则，从而控制权重优化方向，从而达到持续学习效果的方法。其方法简单来讲分为以下三个步骤：

- 选择出对于旧任务（old task）比较重要的权重
- 对权重的重要程度进行排序
- 在优化的时候，越重要的权重改变越小，保证其在小范围内改变，不会对旧任务产生较大的影响

$$L'(\theta) = L(\theta) + \lambda \sum_i b_i (\theta_i - \theta_i^b)^2$$

Diagram illustrating the EWC loss function components:

- Loss for current task ($L(\theta)$)
- How important this parameter is (b_i)
- Parameters to be learning (θ_i)
- Parameters learned from previous task (θ_i^b)
- Loss to be optimized ($L'(\theta)$)

https://hand-craft.net/2016/01/11/elastic-weight-consolidation/

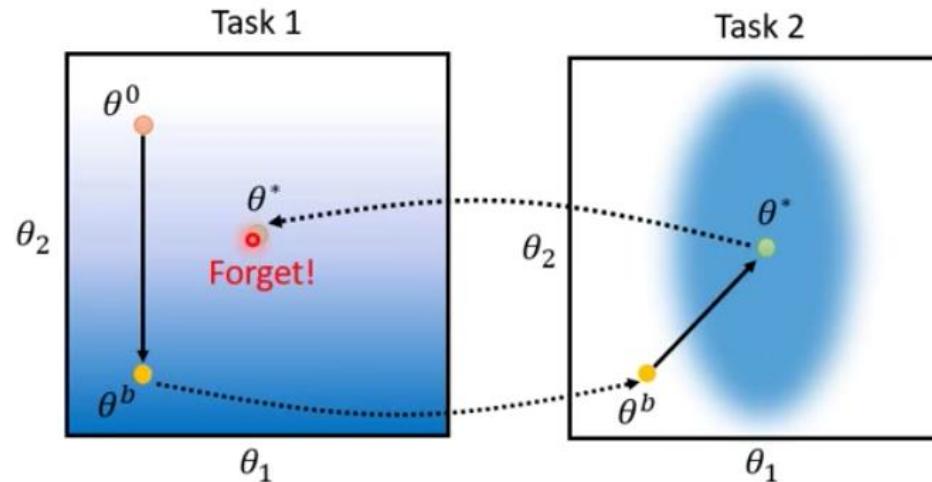
当 $b_i = 0$ 时，表示对于 θ_i 不加约束，它的改变对于模型的效果没什么影响；

当 $b_i = \infty$ 时，表示新模型的参数 θ_i 应该等于原先模型的参数 θ_i^b 。

持续学习

➤ Elastic Weight Consolidation (EWC)

假设使用的模型只有两个参数 θ_1 和 θ_2 ，它们的error surface如下图所示。



The error surfaces of the tasks 1&2
(darker = smaller loss)

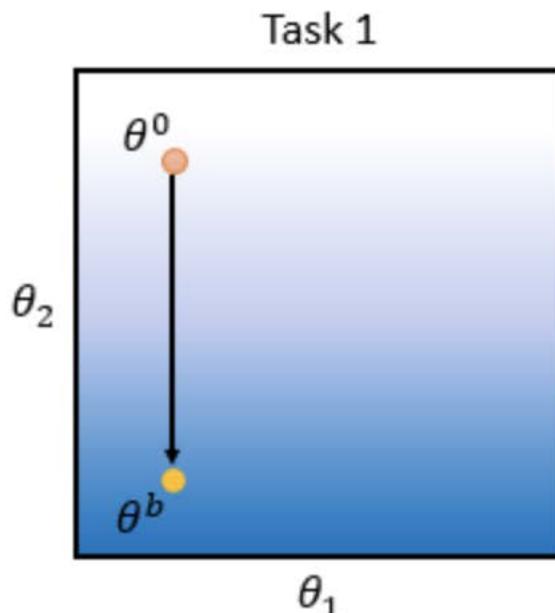
由上图可知，对于task1而言， θ^0 到 θ^b 的方向是损失函数值下降最快的方向，在 θ^b 处Task1可以取得最小损失值；然后将此模型用于Task2，在 θ^0 处Task2的损失函数值却很大，梯度方向是 θ^b 指向 θ^* 的方向。然而， θ^* 在Task1上的损失函数值很大，这解释了模型为什么会灾难性遗忘。

持续学习

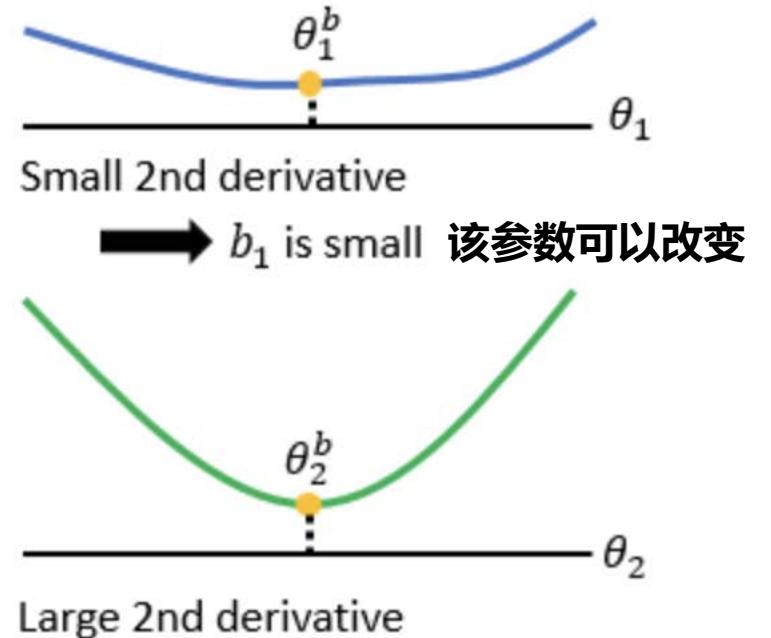
➤ Elastic Weight Consolidation (EWC)

EWC做的是什么呢？

通过引入正则项与参数 b_i 来衡量参数 θ_i 的重要性。当观察 θ_1^b 处的梯度很小时，表明参数变化对Task1模型影响小；当 θ_2^b 处的梯度很大时，表明参数变化对Task2模型影响大。



Each parameter has a
“guard” b_i



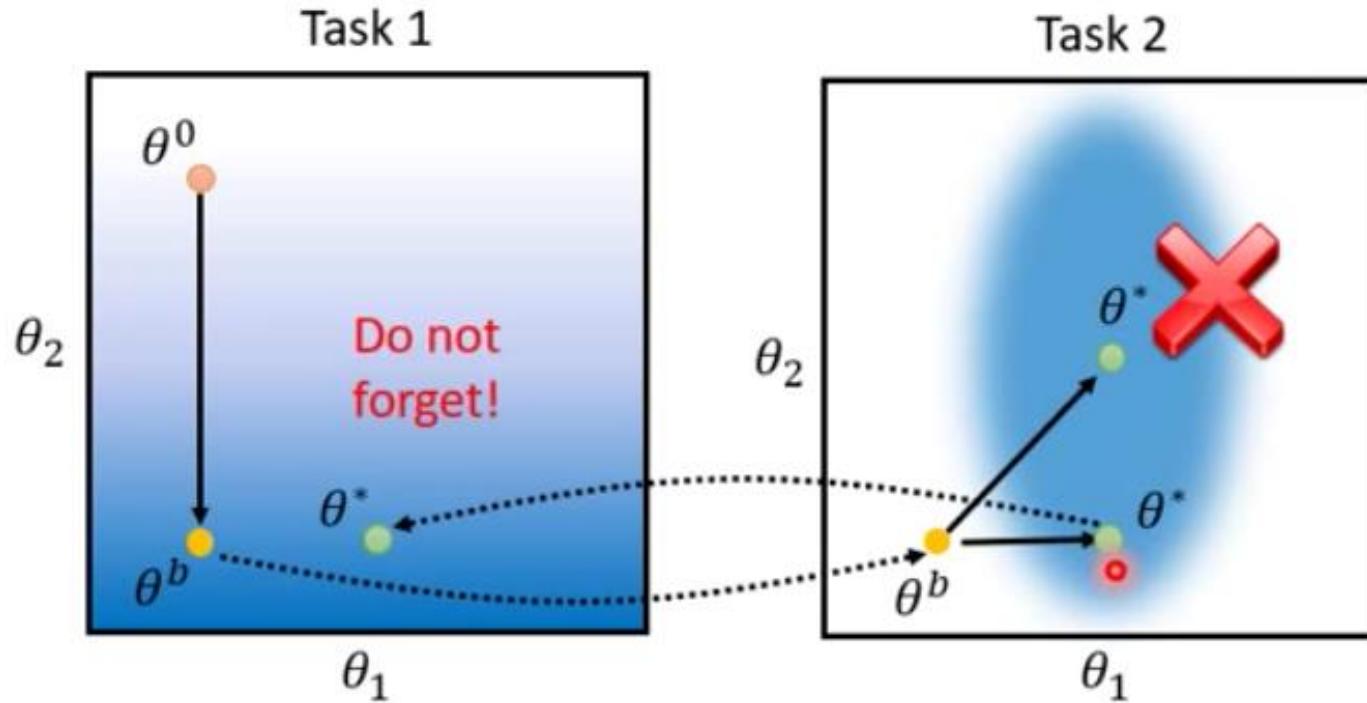
→ b_1 is small 该参数可以改变

→ b_2 is large 该参数不可以改变

持续学习

➤ Elastic Weight Consolidation (EWC)

即在解决新的任务时，需要**寻找最优的参数组合**，使其在所有任务都具有良好的。



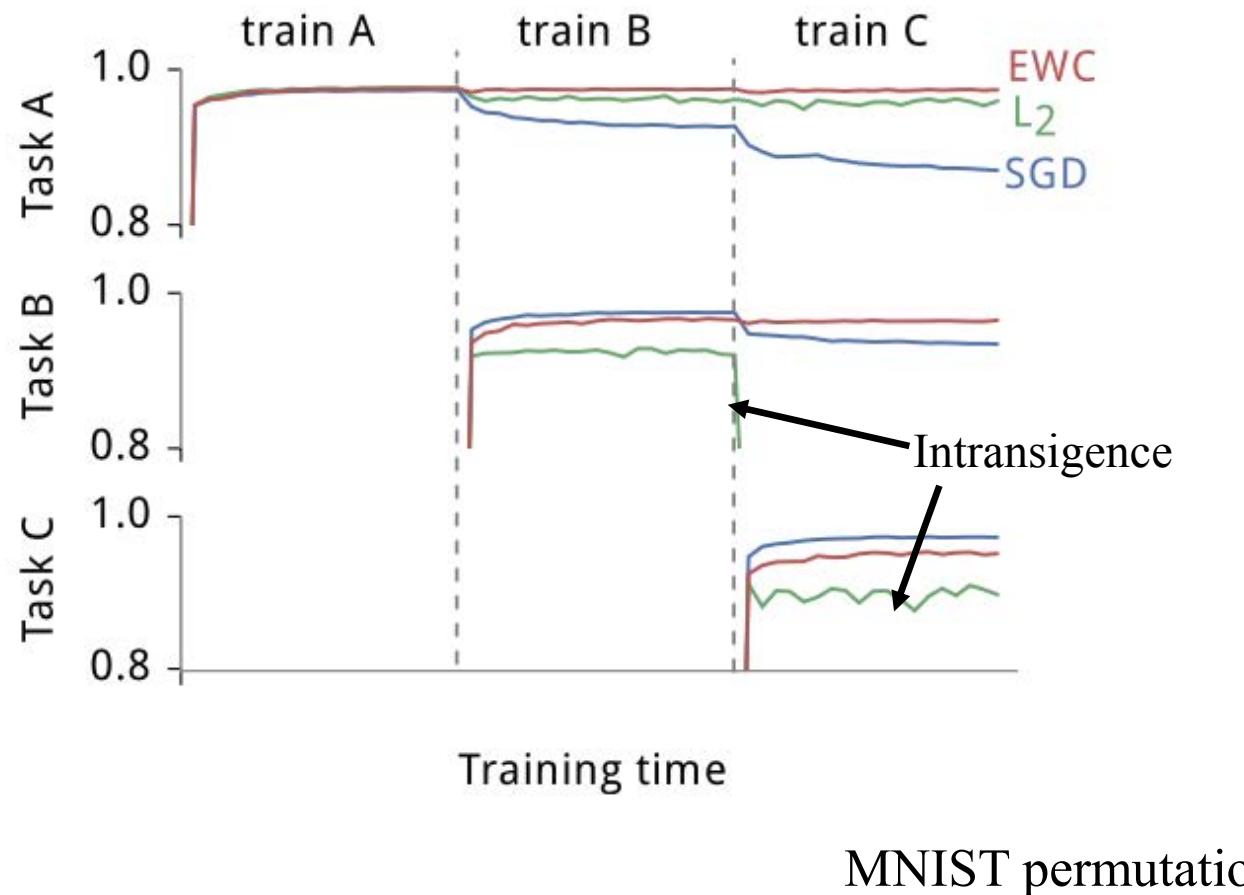
b_1 is small, while b_2 is large.

(可以调整 θ_1 ，但是不能调整 θ_2)

持续学习

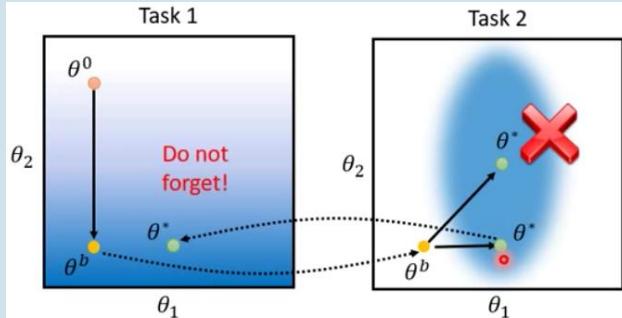
➤ Elastic Weight Consolidation (EWC)

图中显示的是EWC, L2正则化, 和随机梯度下降的实验结果。

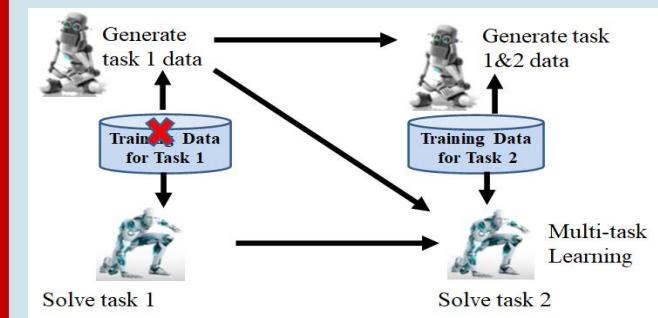


持续学习

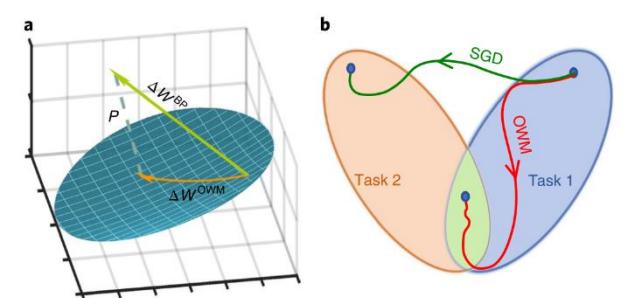
参数加权



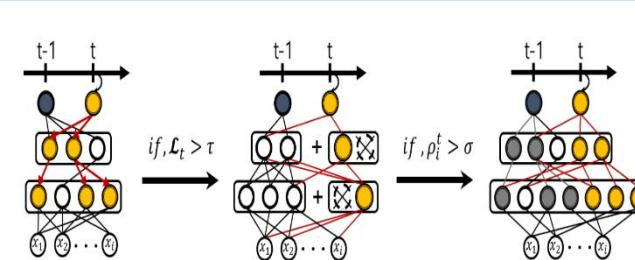
样本回放



约束梯度



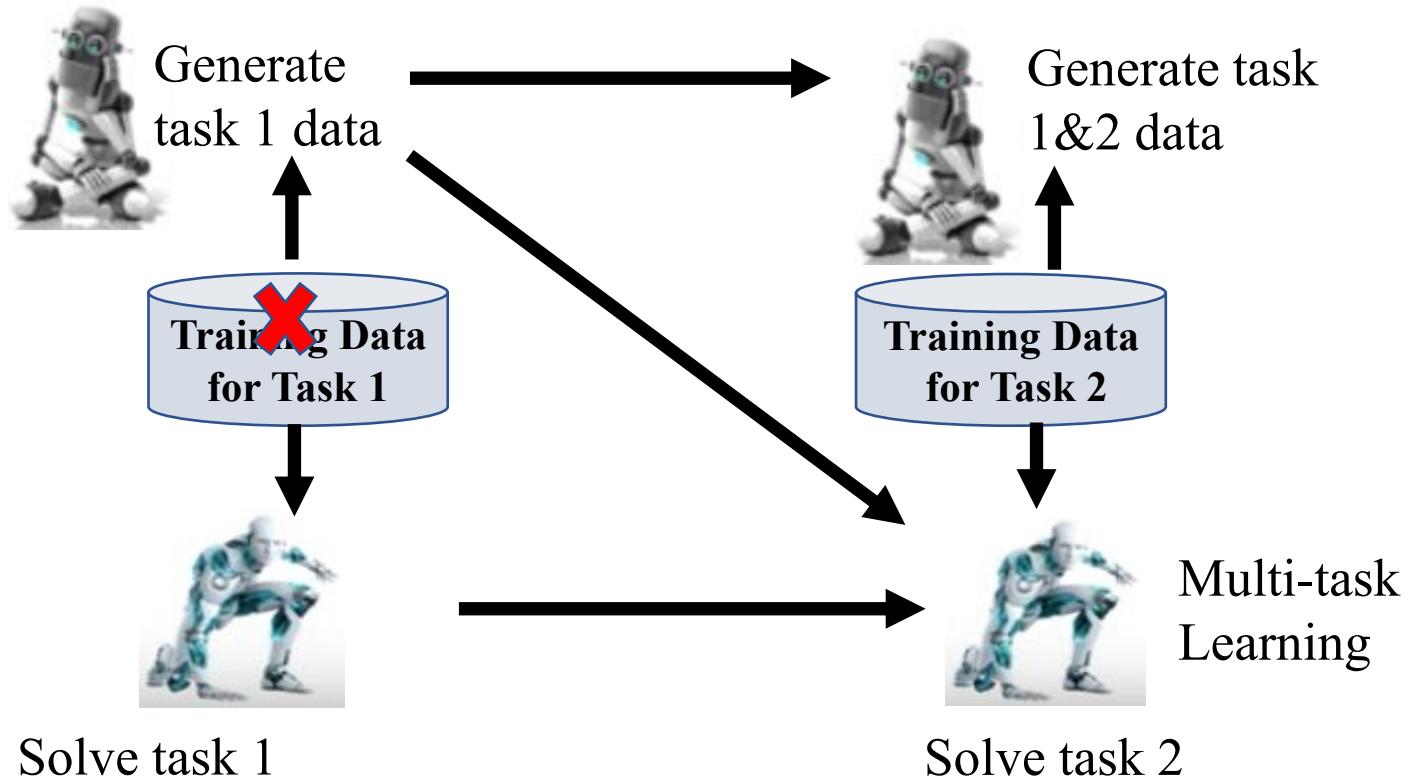
结构扩充



持续学习

➤ Generative replay methods

Motivation :既然我们无法存储所有的训练数据，那么就使用一个生成模型来学习已经训练过的数据的关键信息，当我们在解决下一个任务时，就用生成模型生成一些之前的数据来同时训练。这样既解决了数据存储的问题，又使用了前面所讲的同时在混合数据集上训练一个模型的想法。



持续学习

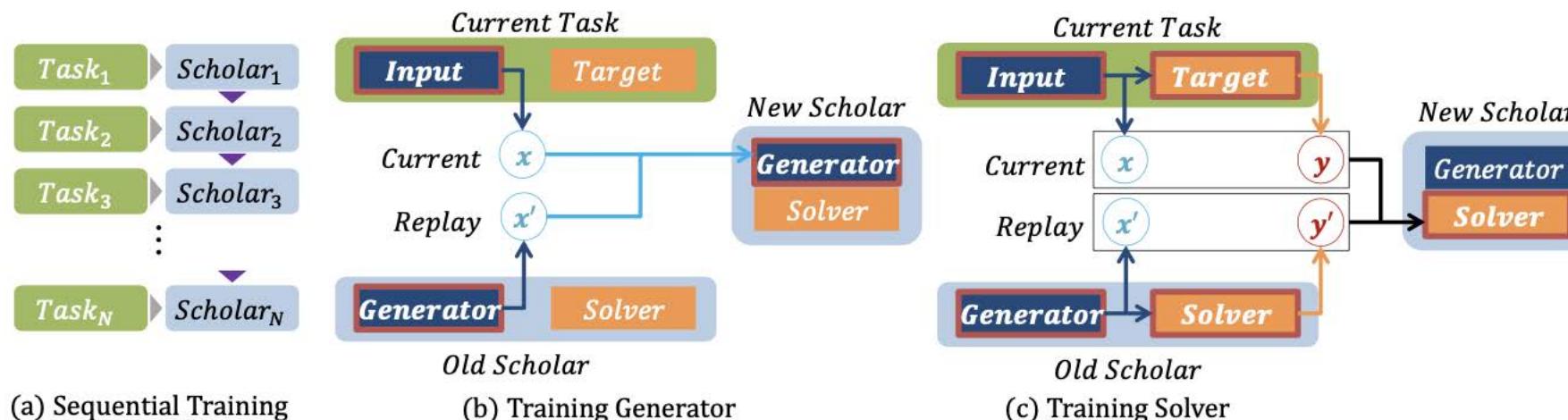
➤ Generative replay methods

Continual Learning with Deep Generative Replay

NIPS-2017

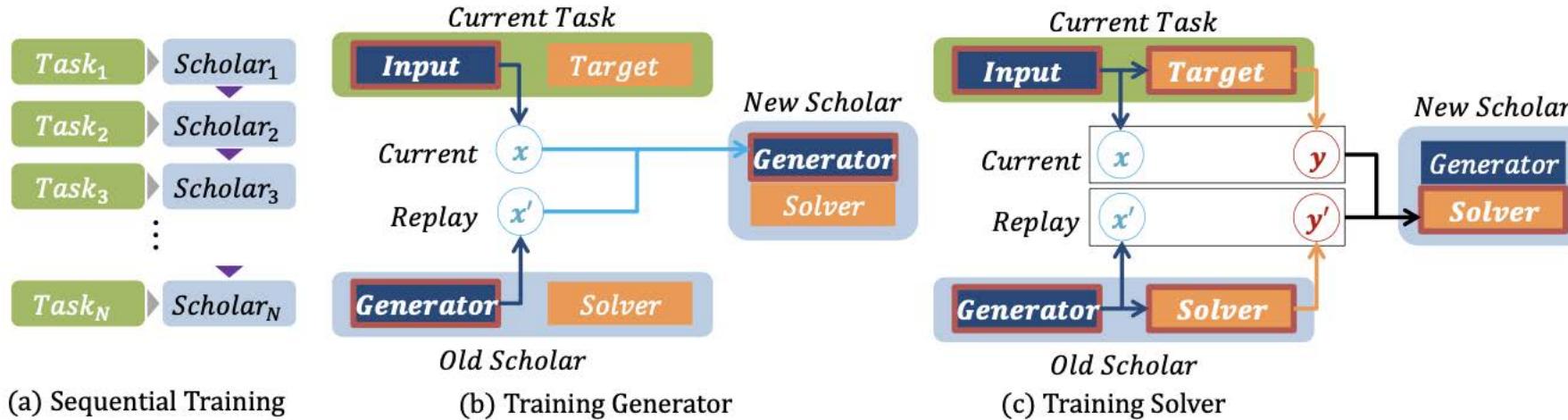
Basic idea: 基于“记忆回放”方法，利用对抗生成网络（GAN），提出了一种Continual learning框架。具体而言，利用GAN网络生成之前任务的样本加入到当前任务的训练中，以实现Continual learning。

该框架由“scholar”组成，它是一个元组 $\langle G, S \rangle$ ，其中 G 是一个生成器， S 用来解决任务的“solver”，下面将其考虑为分类器。它的学习过程：



持续学习

➤ Generative replay methods



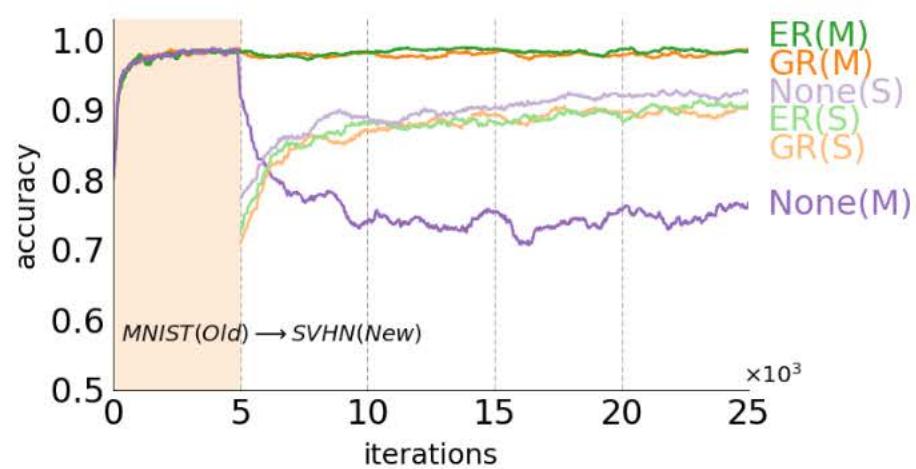
当新任务 $Task_t$ 到达时, $Scholar_{t-1}$ 首先回放之前任务的样本, 然后与当前任务的样本混合用于训练 $Scholar_t$ 的生成器。然后 $Scholar_{t-1}$ 的 $Solver_{t-1}$ 给回放的无标签样本打上标签, 并与当前任务带标签的数据混合用来训练 $Scholar_t$ 的 $Solver_t$ 。这两个过程一起使得当前的 $Scholar_t$ 既吸收了之前学者的知识, 同时也学会了当前任务的新知识, 这样其 $Solver_t$ 就有了解决所有任务的能力了。结合损失函数, 这个学习过程就更加清晰了:

$$L_{\text{train}}(\theta_i) = r \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim D_i} [L(S(\mathbf{x}; \theta_i), \mathbf{y})] + (1 - r) \mathbb{E}_{\mathbf{x}' \sim G_{i-1}} [L(S(\mathbf{x}'; \theta_i), S(\mathbf{x}'; \theta_{i-1}))]$$

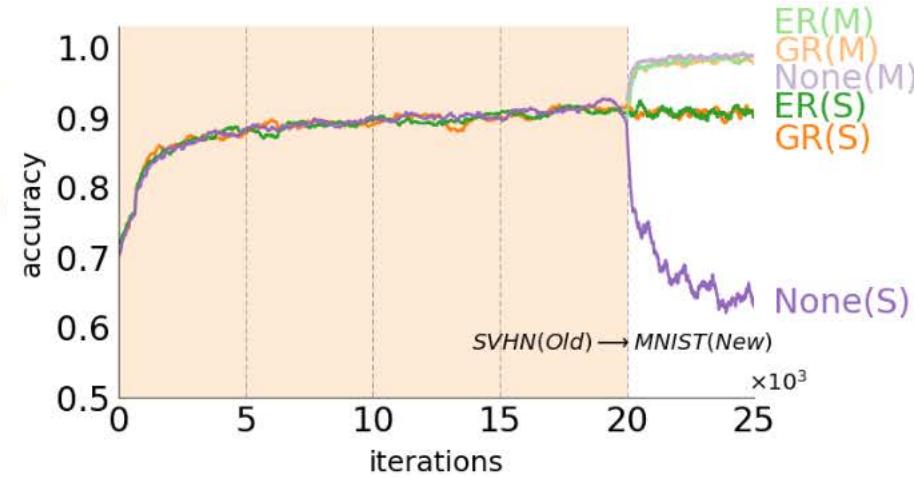
其中 r 用于权衡当前任务与之前任务的重要程度。

持续学习

➤ Generative replay methods



(a) MNIST to SVHN



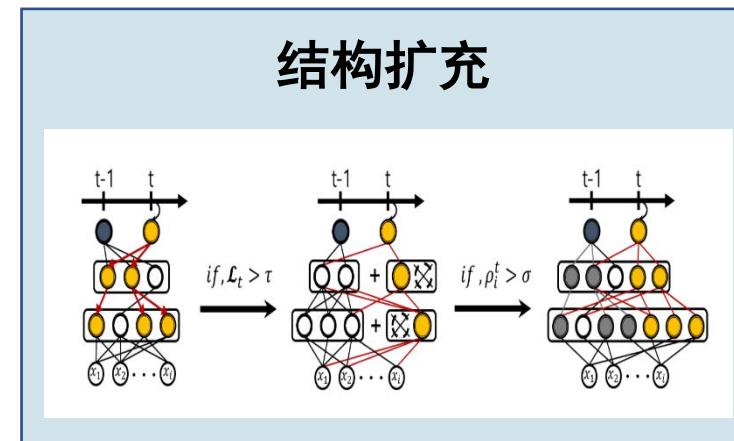
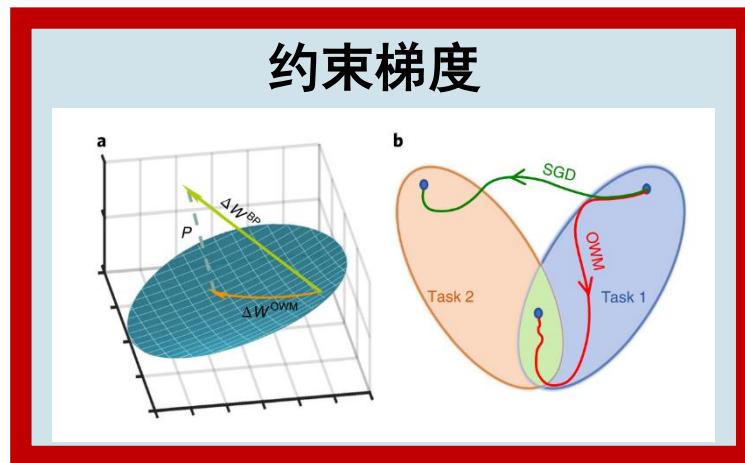
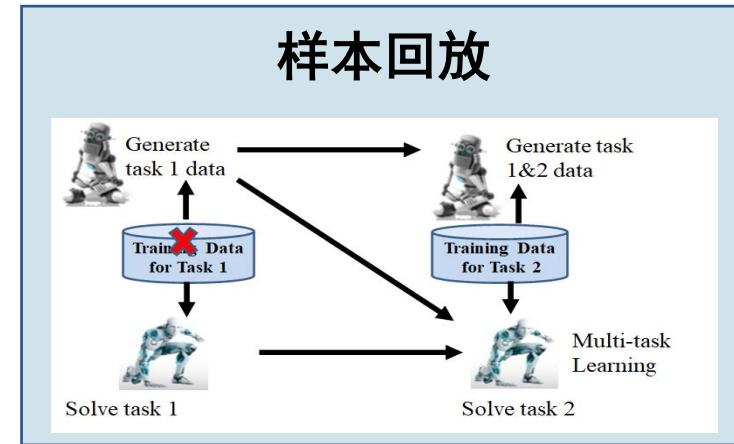
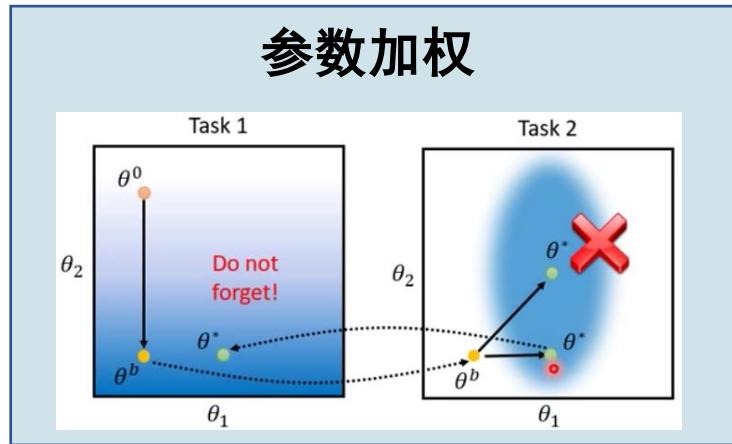
(b) SVHN to MNIST

GR: Generative Replay

ER: Replay actual past data paired with the predicted targets from the old solver network

None: A baseline of naively trained solver network

持续学习



持续学习

➤ 约束梯度

NIPS-2017

Basic idea: 提出了一种称为Gradient Episodic Memory (GEM) 的持续学习模型，该模型可以减轻遗忘，同时又可以将知识有益地转移到先前的任务中。

与EWC不同，GEM采用的是另一种不同的范式。在学习当前任务时，可以访问之前所有任务的部分数据，它们被收集在一个称为“episodic memory”的地方。这其中涉及到几个关键问题：

- (1) 收集的样本要有**代表性**，如何构建每个任务的coreset；
- (2) 如何利用“episodic memory”来解决遗忘问题。

希望新的task的训练过程中，对老的task的效果也有提升。

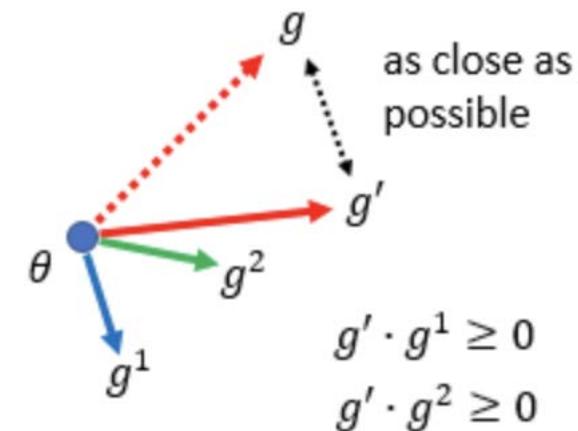
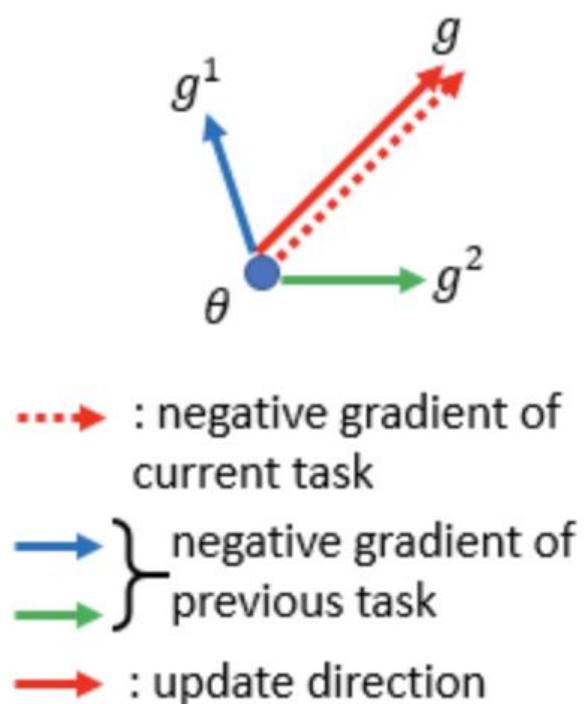
具体做法是，在使用 $\theta_i^{t_j}$ 的gradient $g_i^{t_j}$ 更新 $\theta_i^{t_j}$ 的时候，参考之前task中 θ_i 的gradient $g_i^{t_{j-1}}, g_i^{t_{j-2}}, g_i^{t_{j-3}}, \dots, g_i^{t_1}$ 。使用历史gradient修正当前 $g_i^{t_j}$ 为 g_i' ，使得：

$$g_i' \cdot g_i^{t_k} \geq 0, \forall 0 \leq k \leq j$$

意思就是修正后的方向不与之前的方向重冲突。

➤ 约束梯度

它的原理如下所示，红色虚线箭头 g 表示前面任务损失函数梯度下降最快方向，绿色箭头和蓝色箭头表示新任务损失函数梯度下降最快的方向。GEM所做的的是希望之前任务梯度下降的方向靠近新任务梯度下降的方向，如红色实线箭头 g' 所示，同时又要求 g 和 g' 不要离得太远，表示希望模型不要忘记已经学到的东西。



Need the data from previous tasks

持续学习

➤ 约束梯度

$$\text{Average Accuracy} = \frac{1}{T} \sum_{i=1}^T R_{T,i}$$

$$\text{Backward Transfer (BWT)} = \frac{1}{T-1} \sum_{i=1}^{T-1} (R_{T,i} - R_{i,i})$$

$$\text{Foreward Transfer (FWT)} = \frac{1}{T-1} \sum_{i=2}^T (R_{i-1,i} - \bar{b}_i)$$

- Accuracy 代表在 T 个任务训练完成后对所有任务预测的平均精度
- BWT 表示在 T 个任务训练完成后对过去每一个任务预测带来的平均表现增长
- FWT 表示模型在第 $i-1$ 个任务上训练后对于未来的第 i 个任务所带来的增益。。

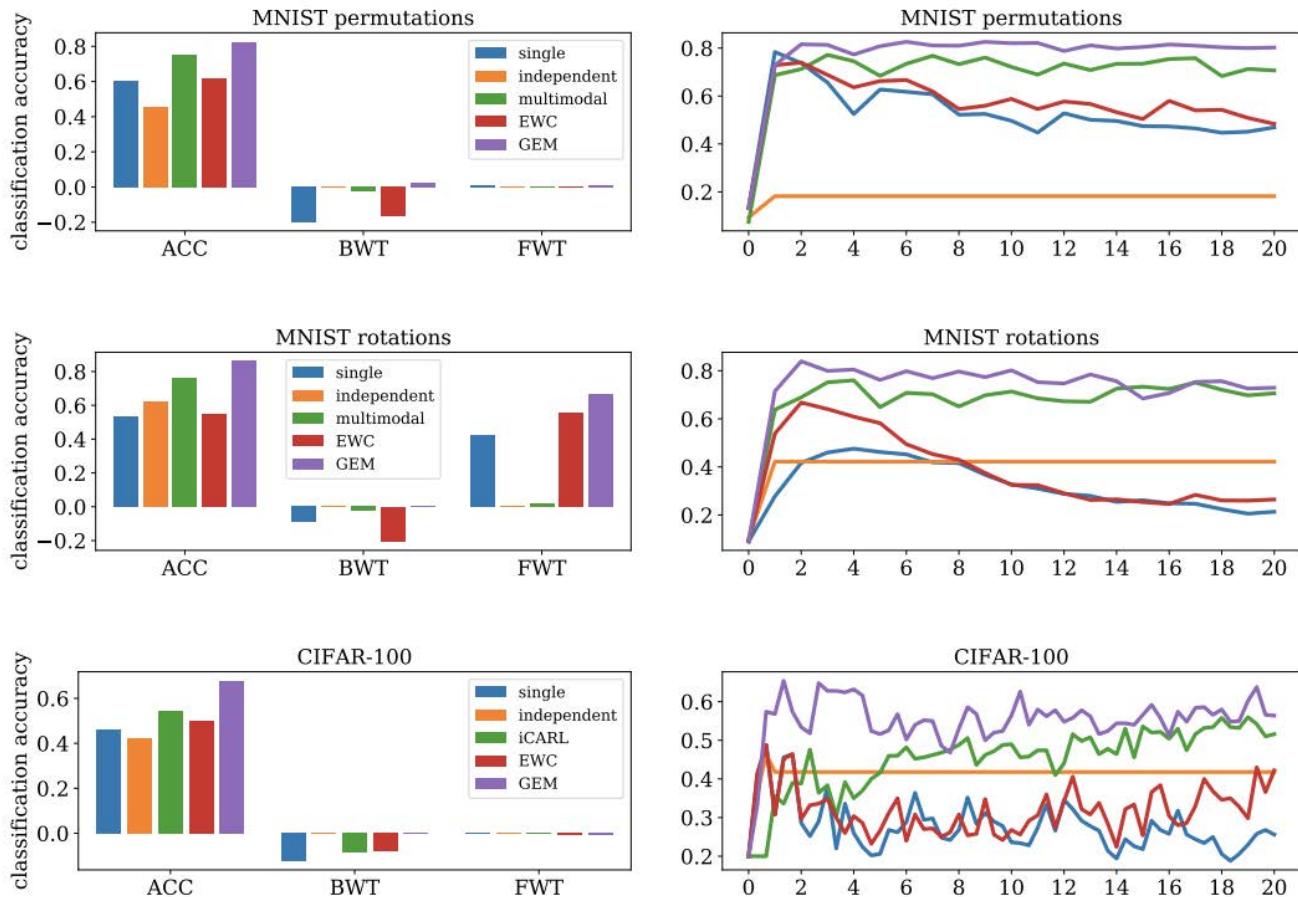
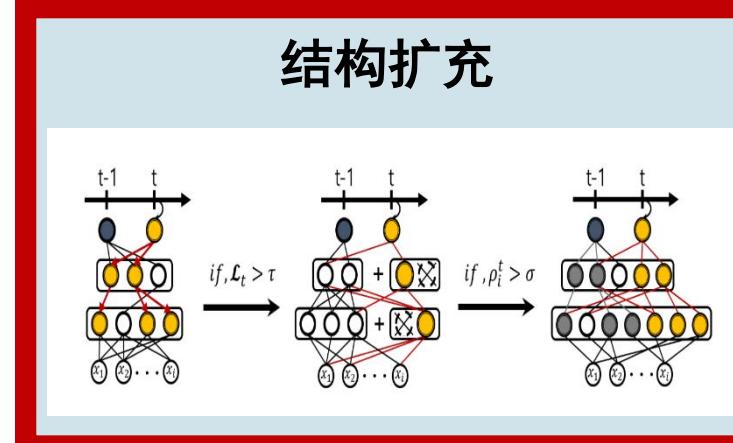
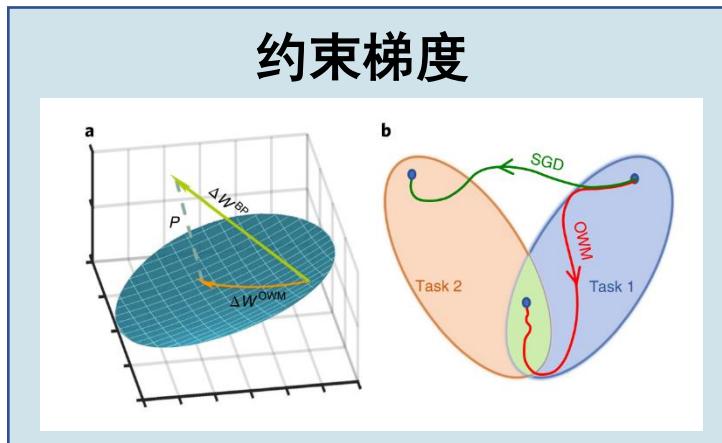
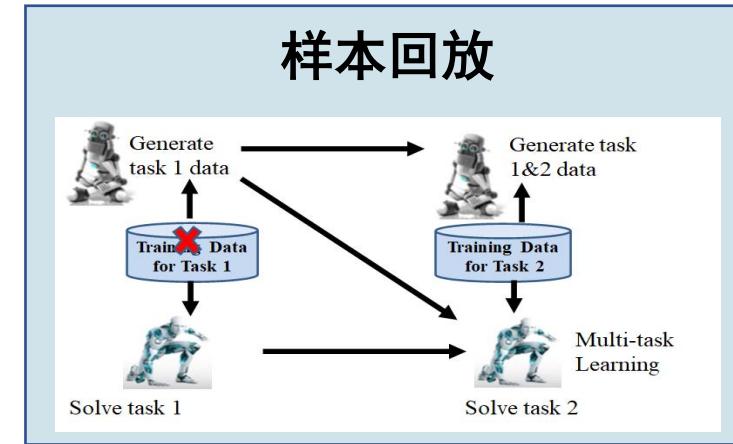
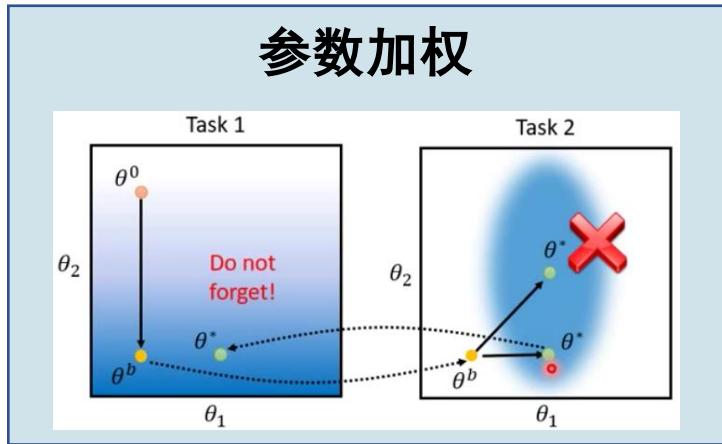


Figure 1: Left: ACC, BWT, and FWT for all datasets and methods. Right: evolution of the test accuracy at the first task, as more tasks are learned.

持续学习

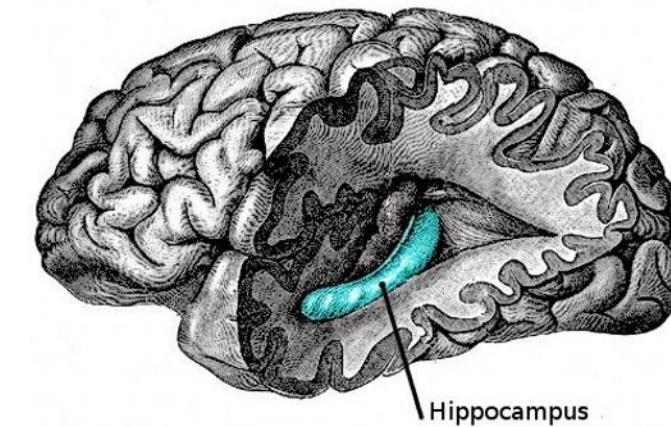
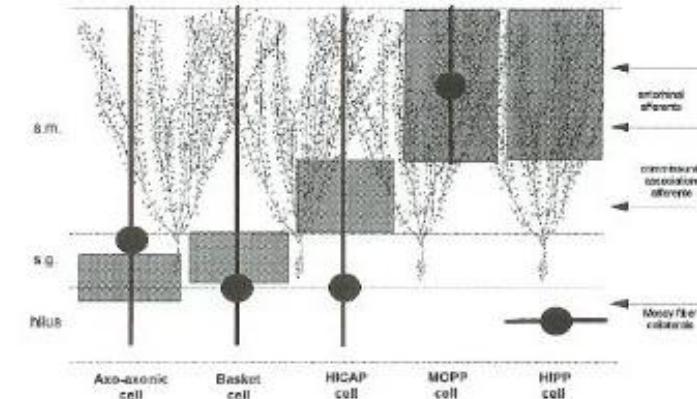


持续学习

➤ 结构扩充

Structural Plasticity via Neurogenesis

- Adult neurogenesis: generation of new neurons in adult brains throughout life, balanced by death of unused neurons (“use it or lose it”)
- In humans, it occurs primarily in the dentate gyrus of the hippocampus
- Increased neurogenesis is associated with better adaptation to new environments.

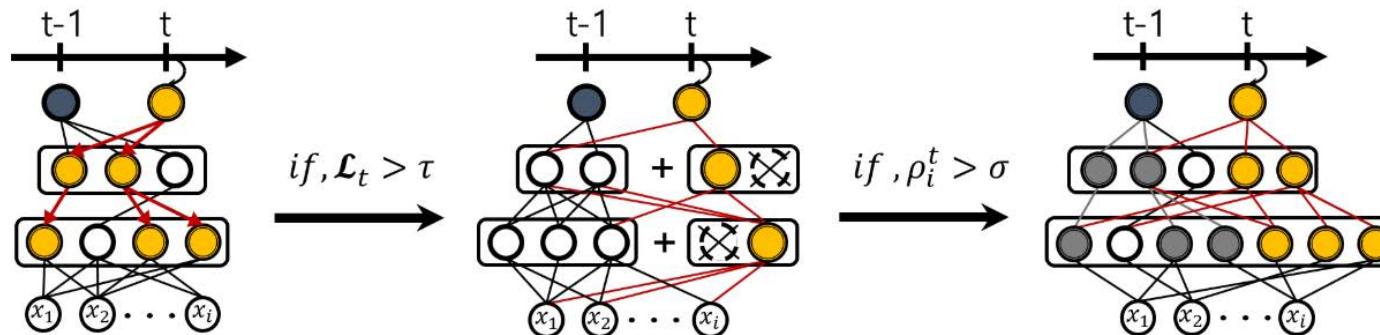


An inspiration for adaptive, expanding neural architecture methods.

持续学习

➤ 结构扩充

- Lifelong learning with dynamically expandable networks

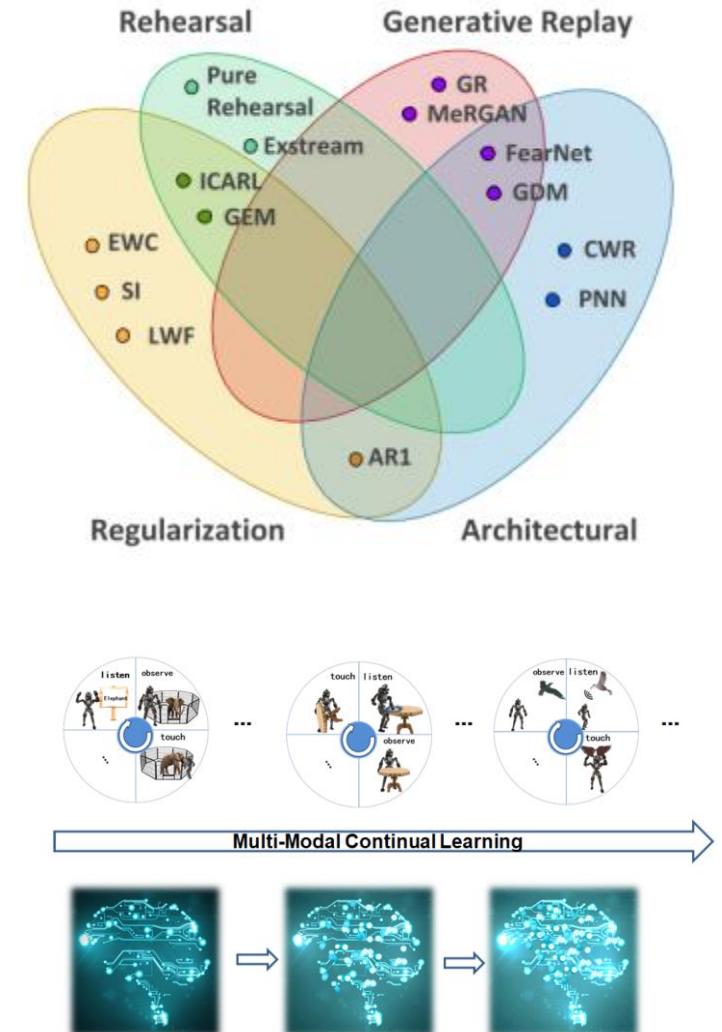


单一任务的网络结构不再是预先给定的，而是在训练任务的同时动态决定网络的容量，在必要时增加一定数量的神经元来扩充网络，相比于Progressive Neural Networks更灵活参数更少。

持续学习

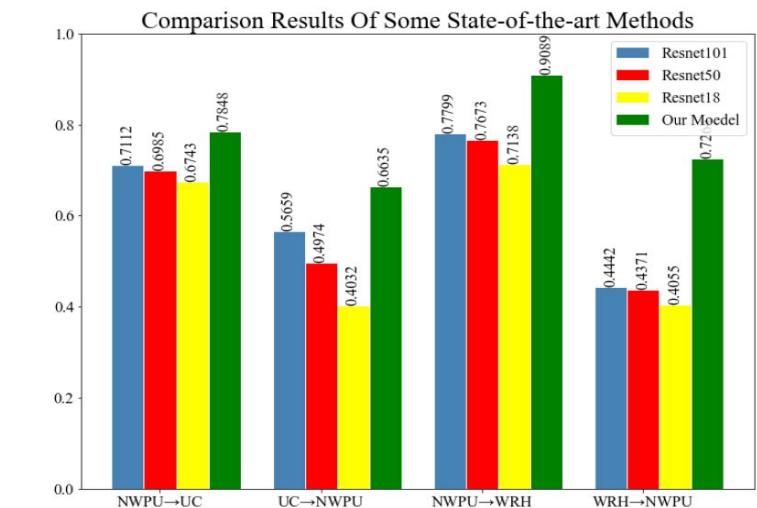
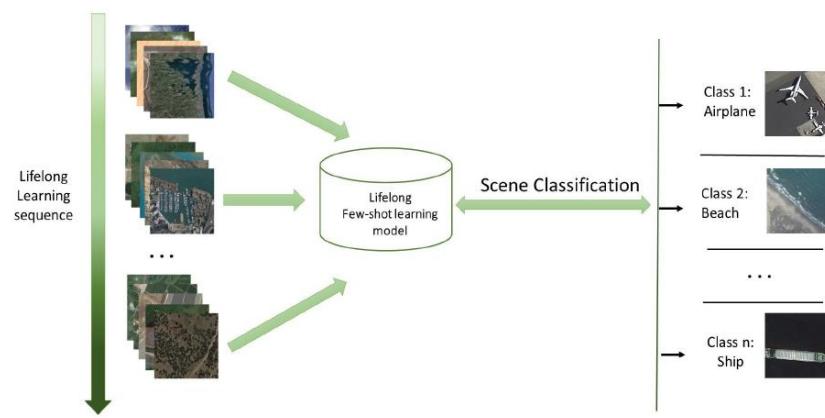
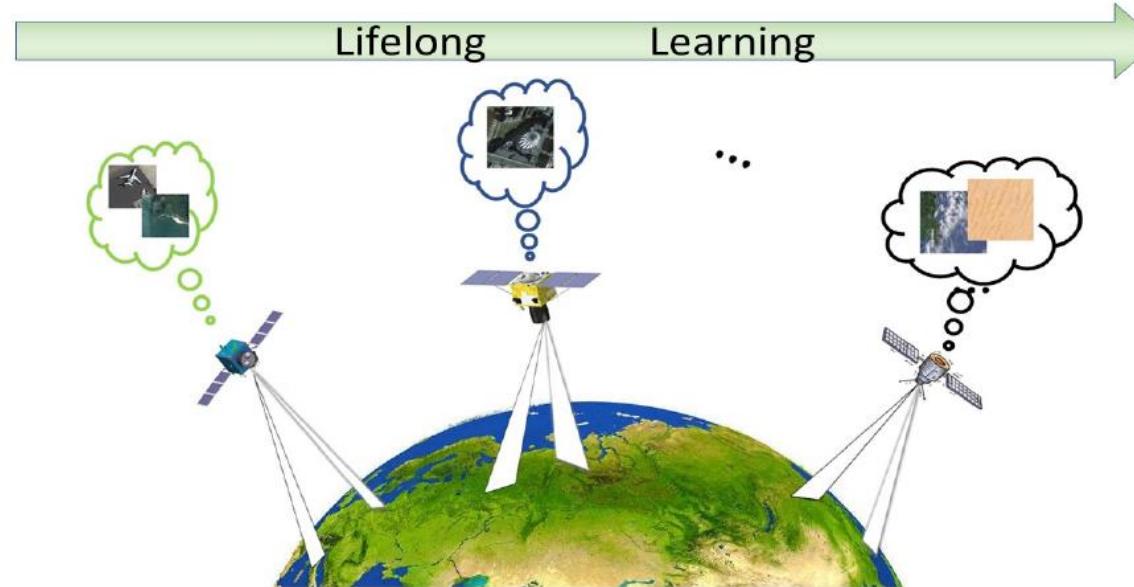
➤ 小结

- 一个完美的持续学习系统具有以下几个特性：首先，能够具有前向迁移的能力，即当前训练所得到的特征具有通用性和可塑性，能够容易地泛化到未来的任务；其次，具备后向迁移的能力，即在新任务上学到的特征具有稳定性，能够保留先前任务的知识；最后，模型具有扩展性，能够在未知任务，未知分类上进行扩展。
- 现有的Lifelong learning方法还有很多不足之处，对未来的展望主要体现在以下三个方面：
 - (1) 可迁移知识表示的研究，尤其是适用于任务类别数大且不一定完全相关的情况的知识表示；
 - (2) 如何处理数据多源异构问题，这里指的是针对来自不同的数据采集源，分布不同的异构数据；
 - (3) 结合其他先进技术，如GAN、强化学习、知识图谱、Spiking neural networks等，旨在提高特征学习性能和隐含因素捕捉能力以及对知识的组织和利用能力。



持续学习

应用



持续学习

➤ 应用

Toward Lifelong Learning for Industrial Defect Classification

A Proposed Framework

By Jingyu Zhang, Di Guo[✉], Yupei Wu, Xinying Xu[✉], and Huaping Liu[✉]

Automatic defect inspection is an important application for the development of smart factories in the era of Industry 4.0. It gathers data from production lines to train a model to automatically recognize certain types of defects. However, the defect types may vary in the production process, and it is difficult for the old model to adapt to new types of defects directly. Considering this problem, we propose an industrial defect classification framework based on lifelong learning, which continuously updates the defect classification model to adapt to different industrial scenarios as new defect appears. Specifically, a novel recursive gradient optimization (RGO) lifelong learning method is used to train the defect classification model, which only needs a fixed network capacity and does not need data replay. The proposed framework is evaluated on an experimental setup of six defect classification tasks. Extensive experiments in real scenarios are performed, demonstrating that the proposed framework can effectively relieve the catastrophic forgetting problem in lifelong learning compared with other state-of-the-art methods.

Digital Object Identifier 10.1109/RA.2023.3258743
Date of current version: xxxxx

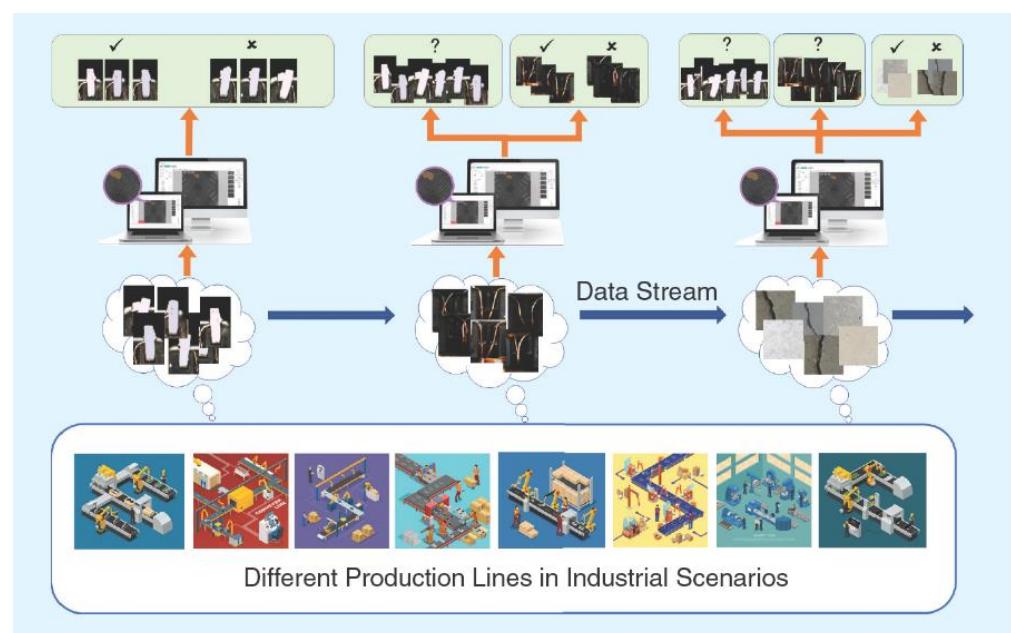
INTRODUCTION

Characterized by the increasing automation and development of smart factories, machine learning techniques are more and more employed in the era of Industry 4.0, among which machine defect inspection is a critical step to guarantee product quality in the process of industrial manufacturing (Figure 1). It utilizes the annotated data collected from production lines to train a defect recognition model, and it has been gradually replacing the traditional manual defect inspection that suffers from huge cost of labor and time.

Generally, the trained model mainly focuses on solving a specific kind of defect. In real industrial manufacturing processes, the production lines often vary due to different production requirements and it is inevitable that new types of defects could appear. Thus, it is difficult for the old model to adapt to all of these new scenarios properly. To solve this problem, we can either train several models for each new independent defect or train a generalized model for all defect types. However, the former approach may generate too many models and require a huge deployment load in production lines. Moreover, for some kinds of defects, the number of data are scarce, which brings difficulties to train an effective model. For the latter approach, it requires to train the model with all the data every time a new



FIGURE 1. A representative machine defect inspection system in industrial production lines.



持续学习

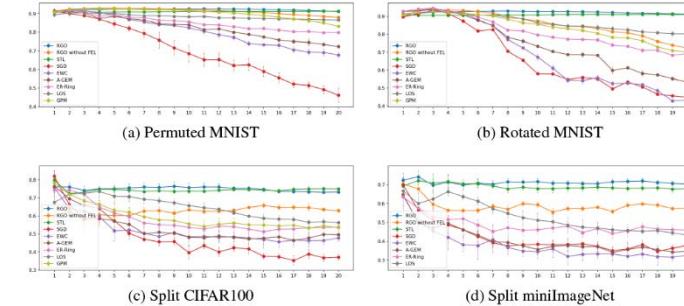
➤ 应用

$$L_k(\theta) = \frac{1}{n_k} \sum_{i=1}^{n_k} l(f(\theta; k, x_{k,i}), y_{k,i})$$

$$F_k(\theta) = \sum_{j=1}^{k-1} L_j(\theta) \approx \sum_{j=1}^{k-1} [L_j(\theta_j^*) + \frac{1}{2}(\theta - \theta_j^*)^T H_j(\theta - \theta_j^*)]$$

$$F_k^{RLL}(\theta) := \frac{1}{2}(\theta - \theta_{k-1}^*)^T \left(\sum_{j=1}^{k-1} H_j \right) (\theta - \theta_{k-1}^*)$$

$$\theta_k^* : \quad \min_{\theta} F_k^{RLL}(\theta), \quad \text{subject to } \nabla L_k(\theta) = 0$$



Theorem 2 (upper bound). Denote $\hat{\sigma}_m(\cdot)$ as the symbol for maximum eigenvalue and η_m as the maximum single-step learning rate, the recursive least loss has an upper bound:

$$F_k^{RLL}(\theta_k^*) \leq \frac{1}{2} n_k \eta_m \hat{\sigma}_m(P \bar{H}) L_k(\theta_{k-1}^*) \quad (8)$$

where $\bar{H} = \sum_{j=1}^{k-1} H_j$ is defined as the sum of the Hessian matrices of all old tasks.

- Continual learning with recursive gradient optimization, ICLR, 2022

持续学习

➤ 应用

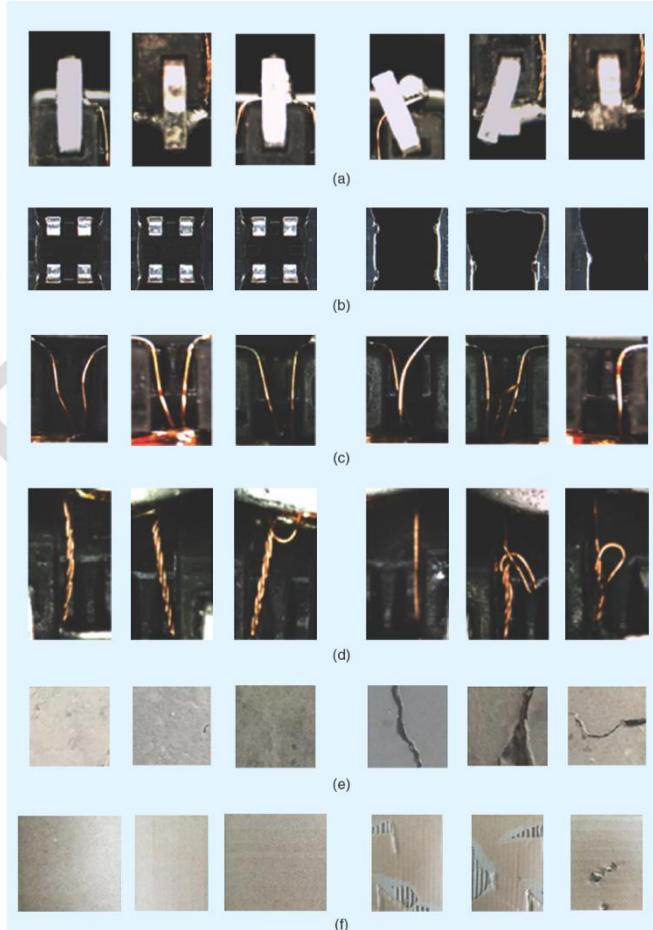


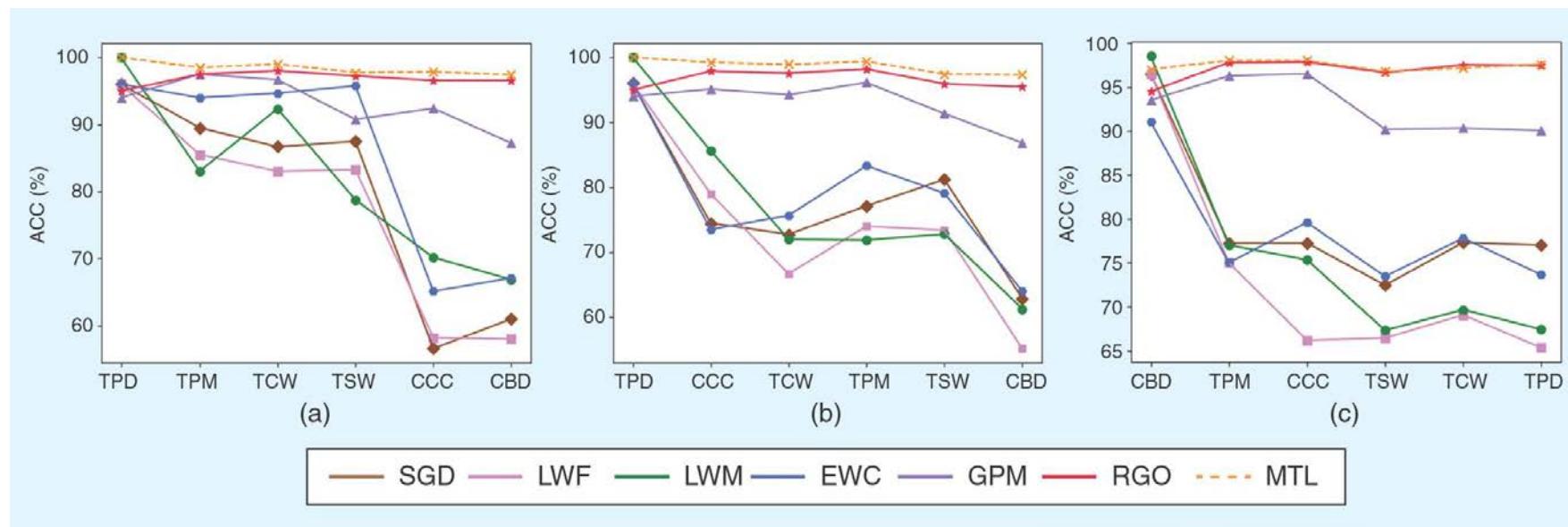
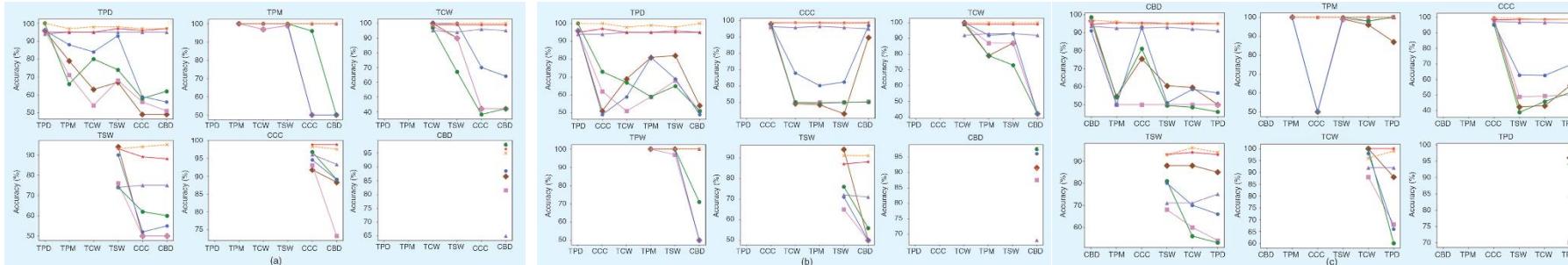
FIGURE 6. Six different kinds of industrial defect data sets: (a) TPD data set, (b) TPM data set, (c) TCW data set, (d) TSW data set, (e) CCC data set, and (f) CBD data set. The images on the left denote the normal samples and the right are the defective samples.

TABLE 1. The number of samples of six industrial defect data sets.

	NUMBER OF TRAIN SET SAMPLES/PIECE		NUMBER OF TEST SET SAMPLES/PIECE	
	NORMAL SAMPLE	DEFECT SAMPLE	NORMAL SAMPLE	DEFECT SAMPLE
TPD	304	51	49	51
TPM	70	70	50	50
TCW	248	42	58	42
TSW	300	50	50	50
CCC	400	400	200	200
CBD	200	200	100	100

持续学习

应用



持续学习

应用

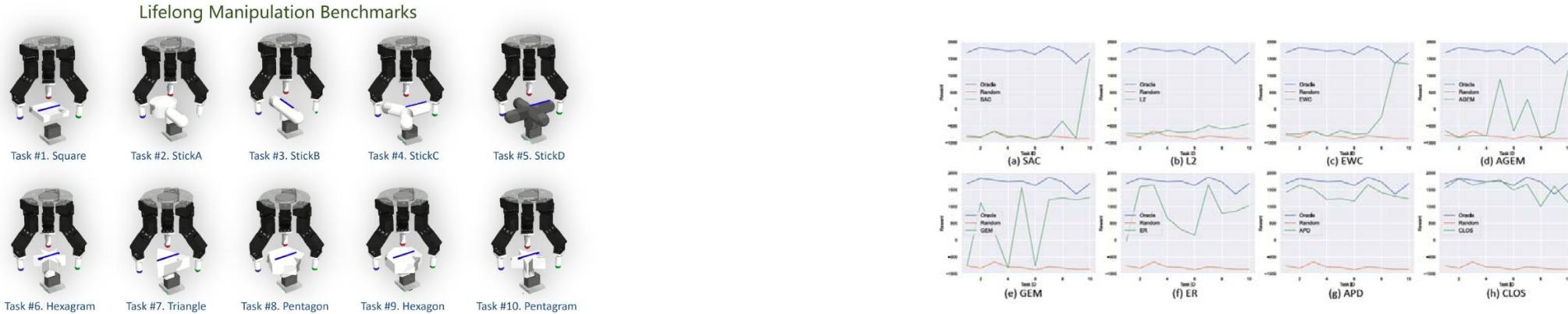


Figure 1: 10 manipulation tasks for lifelong manipulation benchmarks

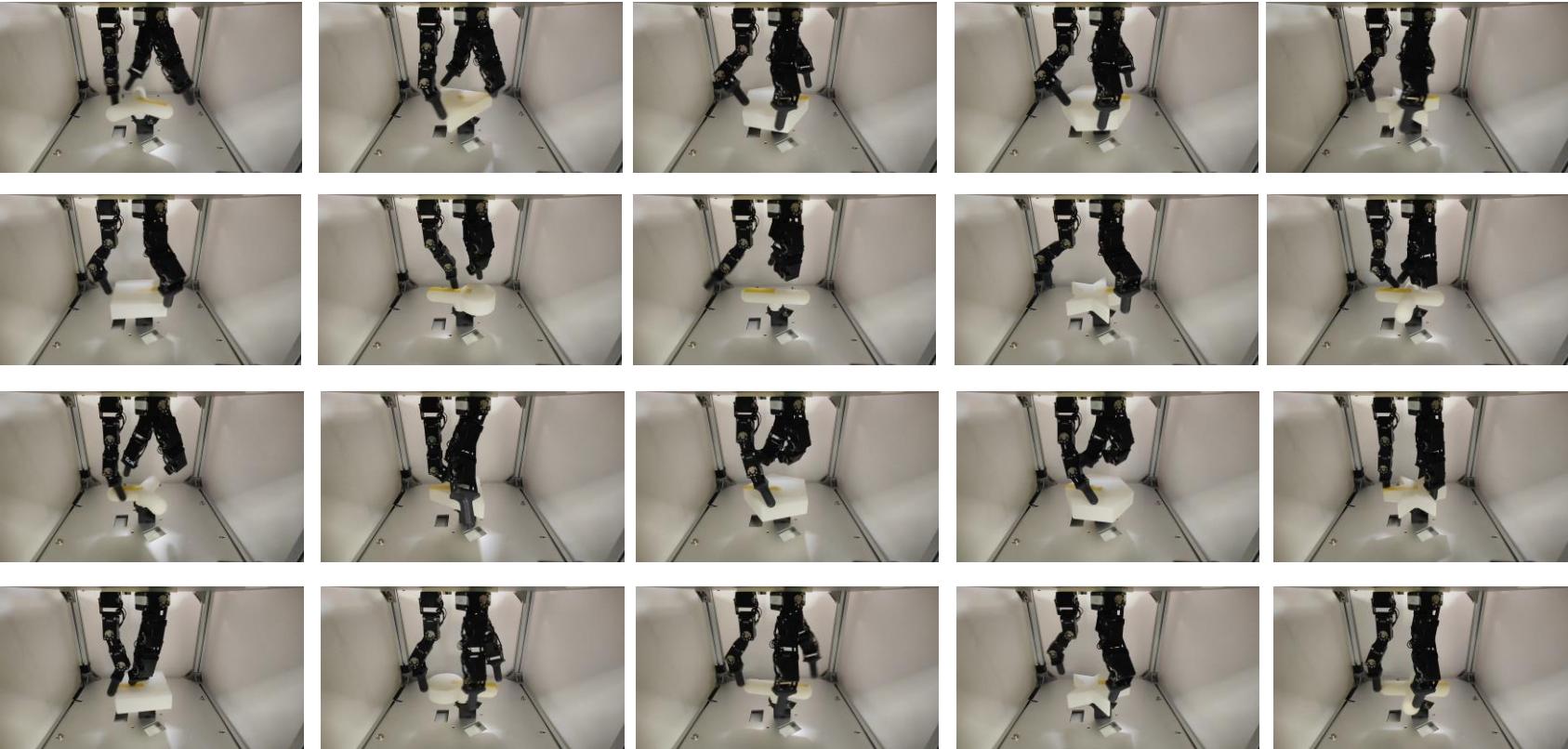
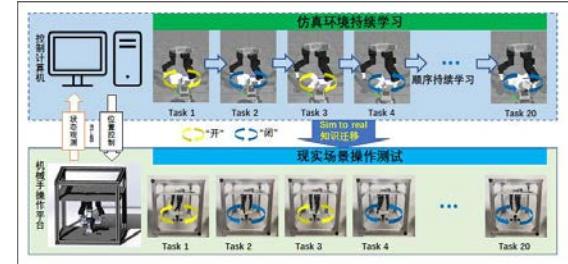
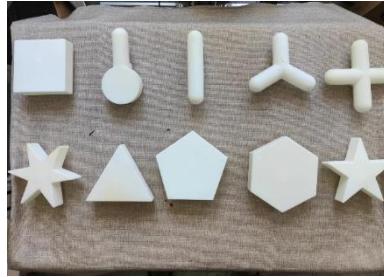
- RL tasks would intensify the weakness of lifelong learning algorithms. EWC, GEM, and AGEM, which are notably effective in supervised lifelong learning, have a poor performance in RL.
- expansion-based methods and memory-based methods have a good performance on lifelong RL methods. Future work could consider incorporating these two methods as a part of its algorithm for a better performance. Regularization-based methods and gradient-based methods still suffer in RL tasks.
- Existing lifelong learning methods requires memorizing lots of data, which makes it impractical in robotic applications. They require memorizing samples or models from previous tasks. However, for applications on real robots, it will not be possible for the robot to remember too many samples or will not allow a substantial increase in the model size. Existing lifelong learning methods still do not have a perfect solution to it.

Evaluations of the gap between supervised and reinforcement lifelong learning on robotic manipulation tasks, in: Proc. of Annual Conference on Robot Learning (CoRL), 2021

持续学习

应用

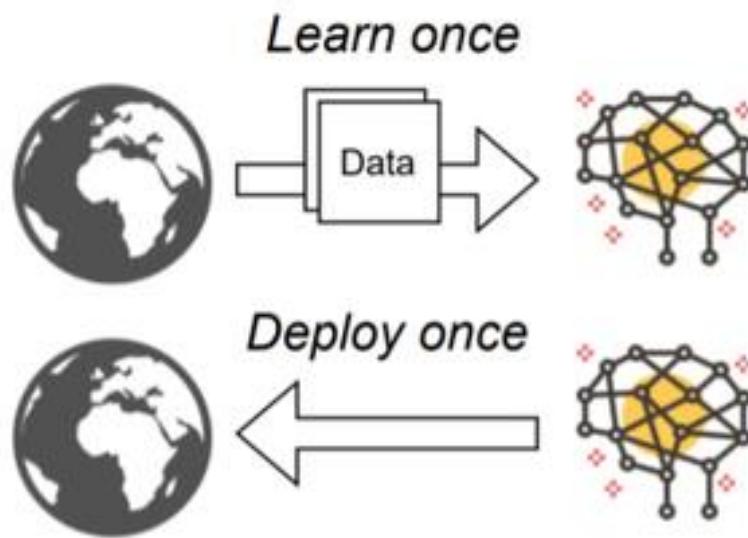
实现20种灵巧操作任务的持续学习。



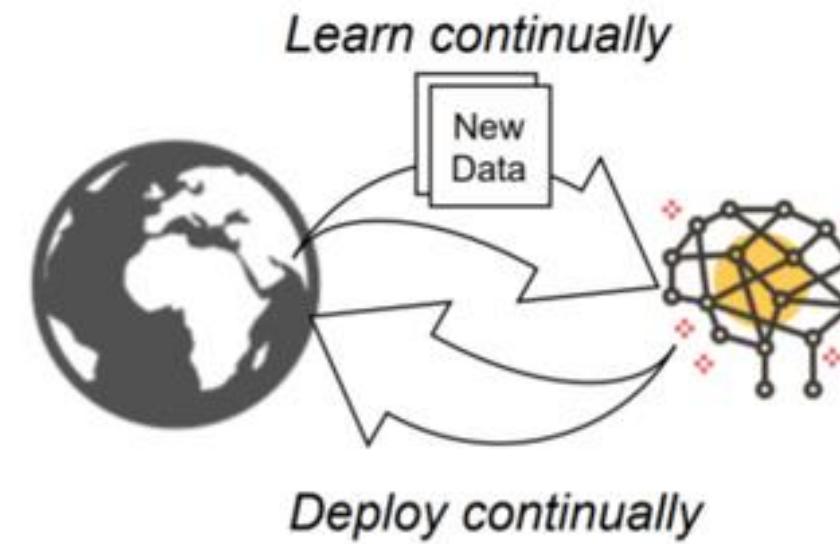
持续学习

➤ 总结

Static ML



Adaptive ML



选题报告

- ~20个队
- 汇报6分钟，交流3分钟
- 4月2日，FIT 1-515, 18:00
- PPT形式
 - 对任务的理解
 - 研究目标
 - 研究内容
 - 计划安排
 - 人员分工（针对多人团队）