# Workshop 1 小样本 self-supervision + few shot learning

Few shot: transfer, 小样本, 快速迁移
↳ 够用吗?
↳ 类内 diversity 大

第二项工作: 对比学习



**Does SSL Benefit FSL?**

- Self-supervised learning (SSL) aims to learn transferable presentation from a large set of unlabeled data:
  - Pretext task-based SSL
  - Contrastive learning
- Both SSL and FSL aim to transfer knowledge learned from the data of a set of seen tasks to a new unseen one:
  - Integrate SSL into FSL --> learn more transferable knowledge?

**Summary and New Results**

- New results for meta-learning:
  - With the PAC-Bayes framework, we derive a fast-rate generalization bound for FSL with DNN. For the first time, we rigorously demonstrate why an FSL model works well on novel tasks. (submitted to AAAI 2022 and ICLR 2022)
  - Large-scale multi-modal pre-training (like OpenAI CLIP and Wenlan-BriVL) may be a possible solution to FSL. The challenge changes to fine-tuning these pre-trained models for FSL.

**Sub-Conclusions**

- We have proposed a novel contrastive prototype learning with augmented embedding (CPLAE) model to address the lack of training data problem in FSL.
- This work shows for the first time that contrastive learning is effective under the supervised and few-shot learning setting.
- Different from existing embedding-based meta-learning methods, we introduce both data augmentation to form an augmented embedding space and a support set prototype centered loss to complement the conventional query centered loss.
- Extensive experiments on three widely used benchmarks demonstrate that our CPLAE achieves new state-of-the-art.

**Summary and New Results**

- Take-home message about few-shot learning (FSL):
  - Episode-level self-supervised learning (SSL) is more effective than instance-level SSL.
  - Contrastive prototype learning brings more benefits to FSL than classic contrastive learning.
  - Concatenating multiple data augmentations can significantly improve the SOTA methods.

# Workshop 2 川大 彭玺

cluster 表示形式? 约空间/列空间上 叫作对比学习

第二项工作: 聚类



**Background**
❖ Contrastive Learning
- MoCo
- SimCLR
- BYOL

**Observation**
❖ Label as representation[1,2]

| | Dog | Plane | Car |
|---|---|---|---|
| | 0.95 | 0.00 | 0.05 |
| | 0.07 | 0.01 | 0.92 |
| | 0.00 | 1.00 | 0.00 |

Cluster Assignment Probability = Soft Label = Instance Representation

Since each sample belongs to only one cluster, ideally, the rows tends to be one-hot.

**Observation** — Regarding the rows of the feature matrix as the soft labels of instances:
- $P(r_j|x_i)$ → the probability of sample $i$ belonging to cluster $j$
- rows → instance representations
- columns → cluster representations distributed over the dataset
- As a result, the instance- and cluster-level contrastive learning could be conducted in the row and column space of the feature matrix, respectively.

**Overview**

约
列 空间间 作
对比学习?

第二项工作: 多模态方法

**Improvement: Confidence-based Boosting**
❖ Two Improvements
- Use mixed augmentation (weak + strong)
- Select confident predictions as pseudo-labels to simultaneously guide the instance- and cluster-level contrastive learning

迁移? c额高数信息的样是(前面epoch找得比较自信的?)

**Confidence-based Boosting**
- Generate pseudo labels based on the confidence of predictions
  - Confidence definition
    $$y_i = g_c(f(x_i))$$
    $$conf_i = \max(y_i)$$
  - Top-$\gamma$ selection
    $$pred_i = \arg\max(y_i)$$
    $$n = \gamma \times N/M$$
    $$CONF_k = sort(\{conf_i | i \in [1,N], pred_i = k\})[n]$$
    $$conf_i \geq CONF_{pred_i}$$
  - Remove unconfident predictions
    $$conf_i < \alpha$$

**Observation** — Multi-view Learning
Either of JR and CR explicitly use cross view consistency which implicitly satisfies:
- Completeness of data
- Correspondence between views

数据是对齐的

**Observation** — Existing works
- POP: Some works have remarkable progress on this problem.
- PVP: Only PVC explicitly considers this challenging problem and its goal is to achieve the instance level alignment (IA), which is daunting and over-sufficient to the discriminative tasks such as classification and clustering.

**Improvement**
❖ Deep Clustering (JULE, DC, etc.)
- Two-stage
  - error accumulation during the alternation
- Offline
  - need entire dataset to perform clustering
  - limited application on large datasets
❖ Contrastive Learning (MoCo, SimCLR, BYOL)
- Instance-level
- General representation
  - off-the-shelf for downstream tasks

❖ Contrastive Clustering (Ours)
- One-stage
  - end-to-end
- Online 在线聚类
  - batch-wise optimization and cluster prediction
- Both instance- and cluster-level
  - dual contrastive learning
- Task-specified
  - clustering-oriented

**Completeness of data**
- Assumption: all samples will be present in all views.
- Partially Data-missing Problem (PDP): some samples are missed in some views.

**Correspondence between views**
- Assumption: data from different views must be strictly aligned.
- Partially View aligned Problem (PVP): portion of the correspondence...

**Motivation**
❖ Our basic idea
- Category level Alignment (CA), which embraces higher accessibility intuitively, an instance have a probability of 1/K or 1/N to be aligned in CA than IA.

**Key Idea**
- We reformulate CA as an Category level alignment problem which could be achieved by contrastive learning...

"最优传输问题"
instance level 对齐 → category level 对齐

随机 → 假阴性
正样本误对

# Conclusion

❖ Achieving online clustering by simultaneously conducting contrastive learning at both the instance- and cluster-level.

- By regarding the rows of the feature matrix as the soft labels of instances, rows and columns of the feature matrix correspond to instance and cluster representation, respectively.
- Select high-confident predictions as pseudo-labels to boosting the two-level contrastive learning.
- Mixed data augmentation (weak + strong).

# Conclusion

(a) Partially view-aligned (b) Categorial identification (c) Noisy correspondence FNPs (d) Category-level aligned

❖ Handle the Partially View-aligned Problem (PVP) by endowing contrastive learning with the robustness against noisy labels (noisy correspondence)

- Category- instead of instance-level alignment.
- Enable contrastive learning robust to false-negative pairs.
- Enrich the learning paradigm with noisy labels by treating the view correspondence as a special noisy label issue.

**Workshop 3 自监督 音视频理解.**

# Motivation

**Human baby learning**
Unsupervised, Weakly-supervised, One-shot learning

- When the baby learns something

Learning common representation inside one group

Distinguishing different characters among different categories

# Related work    信息迁移

- Supervised Pretraining + Finetune (2014)
- Unsupervised Pretraining + Finetune

IM.GENET  → Finetuning → Object Detection

Pretraining on ImageNet Classification

Semantic Segmentation

Fine-grained C...

# UESTC Varying-view action dataset

Three Capture settings

GitHub https://github.com/HRI-UE.../CFM-HRI-RGB-D-action-database

(a) Frame samples of 13 action categories in 8 fixed viewpoints and varying-view sequences.
(b) Temporal frame samples in the varying-view sequence of action a27.

**任意视角**
**汛别**

**Rotation    Category**

**辅助任务**

$X = (\quad, \quad); Y = 3$

Unsupervised Visual Representation Learning by Context Prediction. In ICCV 2015

S^4 L: Self-Supervised Semi-Supervised Lear... ICCV, 2019

**任意视角 → 特征学习 → 一致特征**
**1. 外向扩展式: GAN由一个视角生成更多视角 2. 内敛 View 预处理**
**→ 视角转化/统一.**

# Contrastive Learning

- Image discrimination

#1

#2

Input image

aug. views

ConvNets

features

push

## Arbitrary-view human action recognition via novel-view action generation
PR, 2021

- Motivation

A Common representation

## View-invariant Human Action Recognition via View Transformation Network (VTN)
TMM, 2021

- Motivation
- Larger quantity of possible views provides sufficient samples for self-supervised feature learning.
- Action features of different views have common representation inside one category.

**第二支工作 音视频理解**

## Vision-guided Music Source Separation via a Fine-grained Cycle-Separation Network
ACM MM, 2021

- Motivation
- Videos are composed of visual frames and sounds.
- Audio-vision two modals have a natural correspondence
- We try to build correct correspondence between vision and audio two modalities

Fig. 2. The overview of the proposed FCSN network.

# Related work

- MoCo (CVPR 2020) Kaiming He

MoCo

query — keys — from updating encoder

loss

feature of ...

encoder

momentum encoder slowly update
$\theta_k = m \cdot \theta_k + (1-m) \cdot \theta_q$

BYOL: we do not need negative pairs anymore (NeurIPS' 2020)

- an asymmetric design

+8.5  +2.0  +3.2  74.3

ImageNet-1K

MoCo  SimCLR  MoCo v2  BYOL

**Workshop 4  胡瀚  视觉自监督.**

## Yann LeCun's Cake Analogy

- "Pure" Reinforcement Learning (cherry)
  - The machine predicts a scalar reward given once in a while.
  - A few bits for some samples
- Supervised Learning (icing)
  - The machine predicts a category or a few numbers for each input
  - Predicting human-supplied data
  - 10→10,000 bits per sample
- Self-Supervised Learning (cake génoise)
  - The machine predicts any part of its input for any observed part.
  - Predicts future frames in videos
  - ...s of bits per sample

Credit by Yann L...

## Why Self-Supervised Learning?

- Baby learns to see the world largely by observation

Photos courtesy of Emmanuel Dupoux

Credit by Yann LeCun

## SSL Opened Deep Learning

Science, 2006

Reducing the Dimensionality of Data with Neural Networks
G. E. Hinton and R. R. Salakhutdinov

## Burst of Deep Learning in Computer Vision

- Supervised learning using AlexNet (NeurIPS' 2012)

IMAGENET

ImageNet Challenge

## Supervised Pre-training + Fine-tuning

Pretraining on ImageNet Classification

SIM GENET

Finetuning

Object D... Semantic Segmentation Fine...

一种范式
12年
方法复兴

## How Did We Get Here?

- 2014.6 Exemplar Dosovitskiy et al. NIPS'2014
- 2018.5 Memory bank Wu et al, CVPR'2018
- 2018.12 Deep metric transfer MSRA
- 2019.11 MoCo FAIR
- For the first time, unsupervised pretraining outperform supervised pretraining on 7 down-stream tasks

Image #1 Image #2 Image #3
Pre-text task : Image discrimination

## How Did We Get Here?

by Andrew Zisserman

文本欢类

- Autoencoders — Hinton & Salakhutdinov Science 2006
- Denoising Autoencoders — Vincent et al. ICML 2008
- Exemplar networks — Dosovitskiy et al. NIPS 2014
- Co-Occurrence — Isola et al. ICLR Workshop 2016
- Egomotion — Agrawal et al. ICCV 2015; Jayaraman et al. ICCV 2015
- Context — Noroozi et al 2016
- Split-brain auto-encoders — Zhang et al. CVPR 2017

## Renaissance of Self-Supervised Learning

- Self-Supervised Pretraining + Finetuning

Momentum Contrast for Unsupervised Visual Representation Learning

Kaiming He    Haoqi Fan    Yuxin Wu    Saining Xie    Ross Girshick

Facebook AI Research (FAIR)

2019.11 MoCo
FAIR

有一个代表性的 contrastive learning

## MoCo (CVPR'2020)

- Large dictionary

Credit by Kaiming He

女开团子

发展变化

contrastive learning

query — loss — keys    dictionary
feature of ...

encoder    encoder

MoCo
query — loss — keys — from updating encoder
feature of ...
encoder    momentum encoder
slowly update
$\theta_k = m \cdot \theta_k + (1-m) \cdot \theta_q$

## Main Theme

- Improving ImageNet-1K linear evaluation (top-1 acc)



MoCo 60.6    SimCLR 69.1 (+8.5)    MoCo v2 71.1 (+2.0)    BYOL 74.3 (+3.2)    SwaV 75.3 (+1.0)

## A Finding by BYOL

- MoCo: we need larger dictionary size (more negative pairs)
- BYOL: we do not need negative pairs anymore
  - an asymmetric design

## PIC: a Single-Branch Method (NeurIPS'2020)

## SimCLR (ICML'2020)

- Simpler: no momentum, no memory (dictionary)
- Sufficient distance between pretext tasks and downstream tasks
  - a linear projection layer → a MLP layer
- Self-supervised learning benefit significantly for semi-supervised learning

ICML'20    → SimCLR V2  NIPS'20

## More Insights in SimCLR

- Self-supervised learning benefit more from larger models
- Self-supervised learning benefit significantly from larger training

模型越大 ← 数据越多了(半监督效果越好)

# New Trend:    改善下游任务性能

## Three Main Trends during the Last Year

- More study on BYOL, why it does not collapse
  - BYOL (NeurIPS'20), SimSiam (CVPR 2021)
- Pre-training good features for localization tasks
  - Pixel-level pre-training
    - PixPro DenseCL (CVPR'2021)
  - Object-level pre-training
    - SoCo (NeurIPS'2021)
  - Other pure contrastive learning
    - DINO (tech report)
- Self-supervised learning on Transformers
  - MoCo v3 (ICCV'2021), DINO (ICCV'2021)
  - SSL-Swin/MoBY/EsViT (tech report)
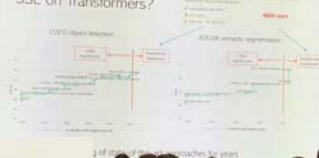
## Improvements on Pascal VOC object detection

- PixPro (CVPR'2021)

MoCo 55.9    SimCLR 56.3 (+0.4)    MoCo v2 57.6 (+1.3)    InfoMin 57.6    PixPro 60.2 (+2.6)

Totally 1.7% absolute improvements in 1 year!

→像素级别预训练→目标级别预训练
Bert for transformer?
CV 〈 Contrast
     MM ?

## SSL on Transformers?

COCO object detection
ADE20k semantic segmentation

## Self-Supervised Learning + Transformer

- SSL can better leverage the model capacity
- Transformers have significantly stronger modeling power than CNNs

MoCo V3
DINO : segmentation
SSL-SWIN

# OPEN Questions:

## Open Crucial Questions

- Can SSL benefit from almost unlimited data?
- What is the relationship with multi-modality learning?
  - E.g., CLIP and DALL-E

## Take-Home Message

- Enjoy the "cake"
- Two trends:
  - Aligning pre-training to down-stream tasks
  - SSL + (Swin) Transformers
- Open critical questions
  - Can SSL benefit from almost unlimited data?
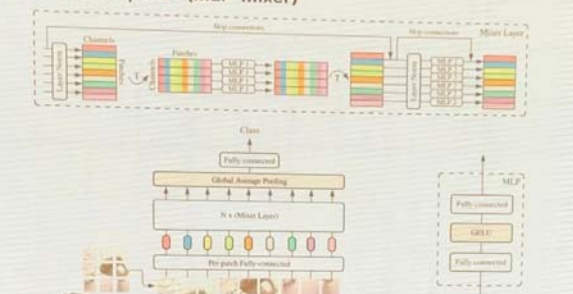  - What is the relationship with multi-modality learning?

## Overview

1. Self-supervised pre-training? (2019.01 "Revisiting SSL")

2. Semi-supervised pre-training? (2019.05 S⁴L)

3. How to measure? (2019.10 VTAB)

4. Large-scale supervised pre-training! (2019.12 BiT)

5. "Let the data speak" (2020.10 ViT, 2021.05 MLP-Mixer, 2021.06 Scaling ViT)

## 1. Self-supervised pre-training? - What's next

- We suggest looking at *architecture* instead of task, leading to new SOTA across all self-sup tasks we tried.
- Uncover key principles for self-sup architectures:
  - Widen the model and add skip connections.
  - These **help much more** than in supervised!
  - Used in all works since ours (CPC, MoCo, SimCLR, BYOL, …)
- Key problem in unsupervised learning: How to select best models? Need (few) labels!

## 2. Semi-supervised pre-training? - Impact

- Need few labels anyways? New model class: Self-supervised Semi-supervised Learning (S⁴L)
- Turn any self-sup task into a semi-sup method!
- We create solid ImageNet 1%, 10% baselines:
  - +20% over previous, used in all papers following ours.
  - This task been widely adopted in the future work: CPC, MoCo, SimCLR, BYOL, …
- "MOAM" gets SOTA on ImageNet with 10% labels:
- Somewhat finicky:
  - Multiple losses may counteract
  - How to balance batch content? BN?

| | Top-5 | Top-1 |
|---|---|---|
| MOAM full (proposed) | 91.23 | 73.21 |
| ResNet50v2 (4×wider) | 81.29 | 58.15 |
| VAE + Bayesian SVM | 64.76 | 48.41 |
| Mean Teacher | 90.89 | |
| UDA | 88.52 | 68.66 |
| CPCv2 | 84.88 | 64.03 |

## 5. "Let the data speak" (MLP-Mixer)

## 5. "Let the data speak" (ViT)

A surprisingly simple and efficient adaptation of Transformer to raw images

## 4. Large-scale supervised pre-training! (BigTransfer/BiT)
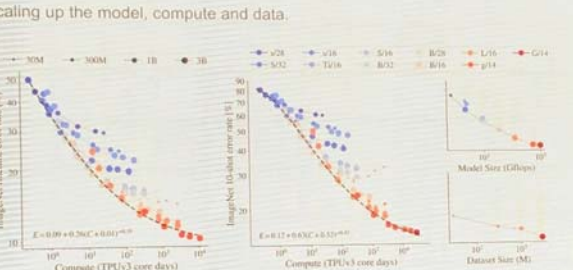
Simple, noisily supervised pre-training gives huge wins!

1. Be patient     2. Scale up *everything*     3. Profit!

Larger models are great at Few-shot learning!

Larger models are very robust (ObjectNet)

## 5. "Let the data speak" (Scaling ViT)

Scaling up the model, compute and data.

## 5. "Let the data speak" (Scaling ViT)

Scaling up the model: What shape can we fit?