# The story of the dark knowledge
## -- from model distillation to contrastive learning

Y.F. Liou

@NTU 2020/10/13

# 今天講題的兩大保證

- 保證講者沒有強大的數學背景
  - 請相信生物學家…
- 保證內容不會是全部正確的
  - 歡迎討論
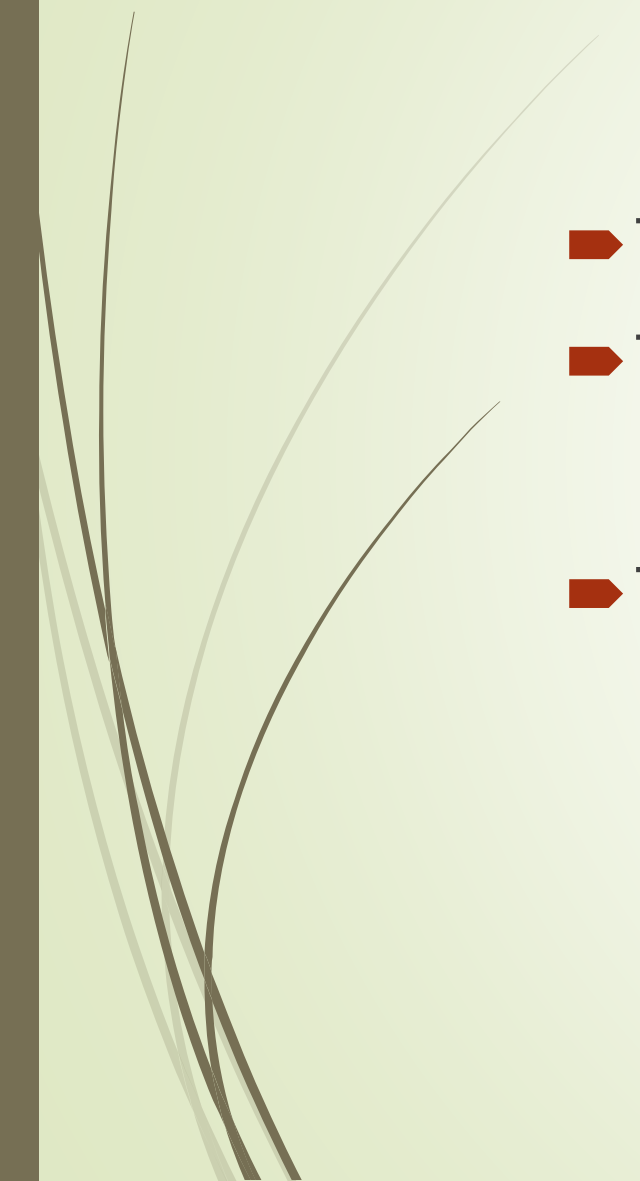
# Before we start

- 我們做個約定
  - 我做"民調"的時候，請各位盡量舉手
  - 希望每個到場的，可以盡量學到有趣的東西

- 我保證
  - 除非你想說話，不然我絕對不會拿麥克風讓你說話

# About me

- Manager of the AI contexture center of Coretronic Corp.

- BS., Life Science, Fu-Jen University

- MS., Biotechnology, National Taipei University of Technology

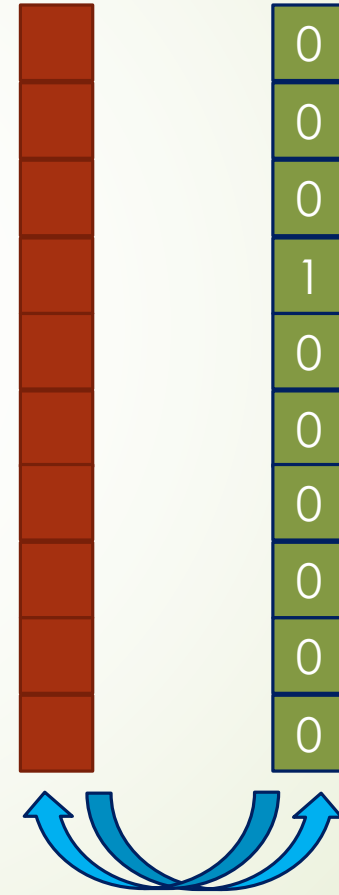- PhD., Bioinformatics and Systems Biology, National Chaio Tung University
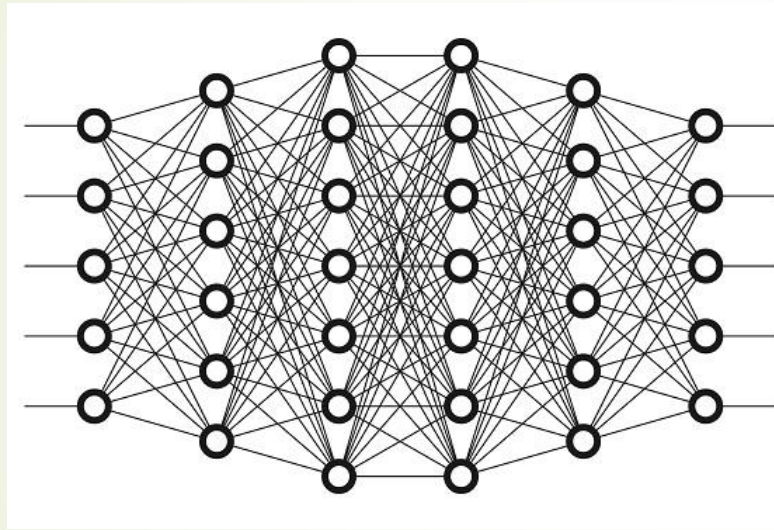
# The issues we will talk

- The basic training process of neural network
- The model distillation
  - Why and how will do the model distillation
- The contrastive learning

# How to training your neural network

# How to train your model – the gradient decent method



Cross entropy

Shannon's entropy equation:

$$H(X) = -\sum_{i=0}^{N-1} p_i \log_2 p_i$$

0
0
0
1
0
0
0
0
0
0

# The minimal risk theory of the traditional machine learning



Whtat happen to recent machine learning area?

# Is the traditional theory is the everything ?



The knowledge revolution is not first time in the human history.



Demo time ~~

# What is happening to overparameterization?

# The issues of over-parameter model would cause



The huge model need huge computing power. But several applications could not provide such computing resource.

# The ways to save the computing resources of neural network

- Physical ways
  - Neural network compiler
    - Operations on the chip
    - Pruning (maybe this can also be thought as a theoretical ways as mentioned below)
  - Quantization
    - The precision of the weights
- Theoretical ways
  - Distillation

- Since the distillation is a "theoretical way", there are many interesting results can be found in such researches.

# The model distillation -- Enhancing the dark knowledge

This would be the beginning of the contrastive learning

Usually,

The main logit will dominate most of output signals. If we want to enhancing the dark knowledge, softening the output would be a way …

# The basic form of "model distillation"

Teacher net

student net

Cross entropy

# When will the biologist be useful ?

## Distilling the Knowledge in a Neural Network

**Geoffrey Hinton**[*†]
Google Inc.
Mountain View
geoffhinton@google.com

**Oriol Vinyals**[†]
Google Inc.
Mountain View
vinyals@google.com

**Jeff Dean**
Google Inc.
Mountain View
jeff@google.com

## 1   Introduction

Many insects have a larval form that is optimized for extracting energy and nutrients from the environment and a completely different adult form that is optimized for the very different requirements of traveling and reproduction. In large-scale machine learning, we typically use very similar models

結論: 生物學家可能只有在introduction開頭有用，後面用處就沒有這麼明顯了。

# Distilling knowledge in a neural network

Teacher net

student net

Results can use:
1. Softmax with temperature
2. Logits (better)

金阿武　劉阿華　劉阿帆

Perturbation the logits would be better.
"Deep Model Compression: Distilling Knowledge from Noisy Teachers"
https://arxiv.org/pdf/1610.09650.pdf

results

Ground truth

loss

loss

Teacher nets are thought to lean much "empirical" information which would contain in the categories not activated.

Distilling the Knowledge in a Neural Network

Geoffrey Hinton[*†]
Google Inc.
Mountain View
geoffhinton@google.com

Oriol Vinyals[†]
Google Inc.
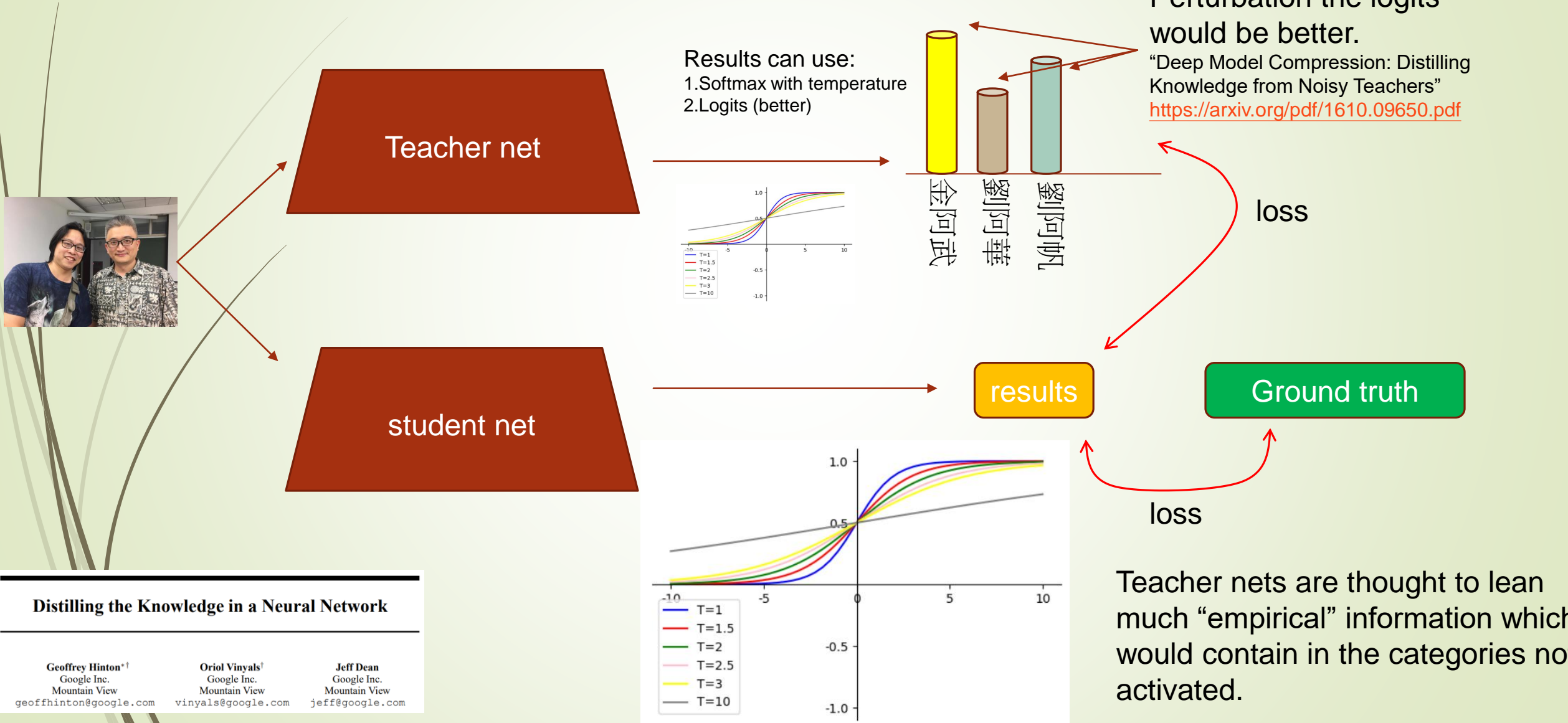Mountain View
vinyals@google.com

Jeff Dean
Google Inc.
Mountain View
jeff@google.com

# Temperature of activation function

$$q_i = \frac{exp(z_i/T)}{\sum_j exp(z_j/T)}$$

z = logits
T = Temperature
j = index of class
q = new logits

目標是把teacher 的label做soft。
可能的原因是teacher所做出來的答案不一定正確，而且內部有許多Dark knowledge。因此做soft後，希望student可以學到teacher所給予的所有資訊。

Here, "Temperature" enlarge the signal of dark knowledge.

```
>>> np.array([1.4, 2.8,.9])
array([1.4, 2.8, 0.9])
```

```
>>> sess.run(tf.math.softmax(a))
array([0.17662444, 0.71624742, 0.10712814])
>>> sess.run(tf.math.softmax(a * .3))
array([0.2956245 , 0.44992913, 0.25444637])
>>> sess.run(tf.math.softmax(a * .5))
array([0.26367459, 0.53097543, 0.20534998])
```

https://arxiv.org/pdf/1503.02531.pdf

# The Dark knowledge

- 這些非label的部分，在knowledge distillation的課題中稱呼"Dark knowledge"
- 影響Dark knowledge的表現，有可能也能影響整個模型
  - What if we control the dark knowledge directly?

0.7 0.1 0.2

0.7 0.2 0.1

這兩個結果對於最後使用SGD更新，回傳的loss是相同的。但是在dark knowledge部分表現非常不同。

The "labeling way" would be a key

# Label smoothing

The purpose of the label smoothing is design to overcome the "overconfidence" which would also cause the overfitting.

$$p_k = \frac{exp(x^T w_k)}{\sum_{l=1}^{L} exp(x^T w_l))}$$

Predictions as a function of activations in penultimate layer

Support :
The model give a softmax result.

$p\_k$: Likelihood the model assigns to the k-th class

$w\_k$: Weights and biases of the last layer

$x$: Vector containing the activations of the penultimate layer

| 0 | 1 | 0 | 0 |
|---|---|---|---|

| .05 | .85 | .05 | .05 |
|-----|-----|-----|-----|

$$y_k^{LS} = y_k(1 - \alpha) + \frac{\alpha}{K}$$

Applying label smoothing to hard targets

| 0 | 1 | 0 | 0 |
|---|---|---|---|

original

$\alpha$= .8

| .05 | .85 | .05 | .05 |
|-----|-----|-----|-----|

smoothing

The traditional method only considers the activated class

After label smoothing, all the logits are considered

# Considering the wrong or poor labeling

If there is a number like this



This number is 5, and it also look like 9.

Neural network

0.7

0.3

1 2 3 4 5 6 7 8 9 0

After feeding into a neural network, The answer of 5 and 9 will both give a penalty.
但實際上選定這兩個答案也有道理。

# Label smoothing helps to find the dense probability

smoothing

If we use the label smoothing, the penultimate layer representations give less spread results.



The "spreading" of the teacher is also containing some information. Using label smoothing (which improve the training in many tasks) in distillation work would hurt the training process.

**Rafael Müller,** **Simon Kornblith, Geoffrey Hinton**
Google Brain
Toronto
rafaelmuller@google.com

https://arxiv.org/pdf/1906.02629.pdf

# Considering the minor output

If there is a number like this

This number is 5, and it also look like 9.

Neural network

0.7

0.3

1 2 3 4 5 6 7 8 9
0

After feeding into a neural network, The answer of 5 and 9 will both give a penalty.
但實際上選定這兩個答案也有道理。

從模型預測的其他label輸出似乎也代表特殊意義。
• 是否表示"在這些次要的輸出中，模型告訴我們看到了甚麼"。
• 也許某些資訊也可以透過控制這些次要的輸出，反饋給模型。

| 0 | 1 | 0 | 0 |
|---|---|---|---|

| .05 | .85 | .05 | .05 |
|-----|-----|-----|-----|

Considering the minor outputs of logits is quite similar to consider the "dark knowledge" of model distillation.

# Teacher free knowledge distillation

- Label smoothing 跟knowledge distillation 都專注在minor logits

  - Label smoothing – 不依賴任何prior，直接設定超參數來看結果

  - Knowledge distillation – 依賴teacher 給予比重

  - 兩種想法結合以後，也許可以達到self-regularize。因為如果模型將錯誤的minor logits提高，也許就能更進一步讓模型來"理解"問題。

| Model | Baseline | Tf-KD$_{self}$ |
|-------|----------|----------------|
| MobileNetV2 | 68.38 | 70.96 (**+2.58**) |
| ShuffleNetV2 | 70.34 | 72.23 (**+1.89**) |
| ResNet18 | 75.87 | 77.10 (**+1.23**) |
| GoogLeNet | 78.72 | 80.17 (**+1.45**) |
| DenseNet121 | 79.04 | 80.26 (**+1.22**) |
| ResNeXt29 (8x64d) | 81.03 | 82.08 (**+1.05**) |

(a) Test accuracy (in %) on CIFAR100. We run 3 times and report the mean of the best results.

使用不同溫度，可以凸顯出來的minor logits強度也不同。就看希望模型可以注意到多小的程度。

http://static.kancloud.cn/mikl_maple/python/1726331

# What if controlling the dark knowledge directly

# mixup

$$\tilde{x} = \lambda x_i + (1 - \lambda)x_j, \qquad \text{where } x_i, x_j \text{ are raw input vectors}$$

$$\tilde{y} = \lambda y_i + (1 - \lambda)y_j, \qquad \text{where } y_i, y_j \text{ are one-hot label encodings}$$

$\lambda$ + $1-\lambda$ →

| 1 | 0 |
|---|---|

| 0 | 1 |
|---|---|

| $\lambda$ | $1-\lambda$ |
|---|---|

Main purpose:
1. Data augmentation
2. Maybe ... Dark knowledge

# Manifold mixup



Layer k

Eligible layers (S)

$1-\lambda$     $+\lambda$

$\lambda$   $1-\lambda$

1   0

0   1

*Manifold Mixup* passes this sanity check (consult Appendix D for further details). While we found that using *Manifold Mixup* improves the robustness to single-step FGSM attack (especially over Input Mixup), we found that *Manifold Mixup* did not significantly improve robustness against stronger, multi-step attacks such as PGD (Madry et al., 2018).

# CutMix

label部分直接使用圖像所佔的比例，舉例:

M     $X_A \odot M$     $X_B \odot (1 - M)$

$$\tilde{x} = \mathbf{M} \odot x_A + (1 - \mathbf{M}) \odot x_B$$

使用Mask遮蔽後，把兩張圖片再合成一張

W

H

rx

ry

$HW - r_x r_y$

$r_x r_y$

$= 1 - (r_x r_y / HW)$

$= r_x r_y / HW$

如果Mask掉大小為$r_x r_y$，則圖A像原始比例剩下$HW - r_x r_y$，另外一張就是$r_x r_y$。
Label比例也調整為$(HW - r_x r_y)$跟$(r_x r_y)$

| | Mixup [48] | Cutout [3] | CutMix |
|---|---|---|---|
| Usage of full image region | ✔ | ✘ | ✔ |
| Regional dropout | ✘ | ✔ | ✔ |
| Mixed image & label | ✔ | ✘ | ✔ |

Table 2: Comparison among Mixup, Cutout, and CutMix.

內涵:
Dark knowledge部分，如果一開始無法從pretrain model中取得做蒸餾，也許我們可以自己製作這些dark knowledge。

# Extract the knowledge of unlabeled data

If you have lots of unlabeled data and you want to leverage the information inside them

# FixMatch



1. 先取出任意資料，並對該資料做 augmentation

2. 使用模型把augmented的資料放入模型中，如果模型有算出某個類別有高於特定的threshold(prediction的虛線)，就使用該圖，並且標上那個特定的類別。

依據接下來的任務來指定 pseudo-label

3. 同一個圖加上更多的noise做augmenting，並放到同一個模型當中，最後把算出來的結果進行cross entropy的運算

https://arxiv.org/pdf/2001.07685.pdf

# MixMatch

After various objections, the most important result will be get after average.

Consistency regularization



Figure 1: Diagram of the label guessing process used in MixMatch. Stochastic data augmentation is applied to an unlabeled image $K$ times, and each augmented image is fed through the classifier. Then, the average of these $K$ predictions is "sharpened" by adjusting the distribution's temperature. See algorithm 1 for a full description.

https://arxiv.org/pdf/1905.02249.pdf

# Unsupervised data augmentation (UDA)

Normal process to train the neural network



1. 把原始資料進行 augmentation(加噪音等等的處理)
2. 以原始圖透過模型預測出來的數值作為答案，augmentation的圖最後答案應該要跟原始圖的答案一模一樣。

https://arxiv.org/pdf/1904.12848.pdf

# AugMix

發揮你的想像力!!愛怎麼au就怎麼au



w1+w2+w3 = 1

Weights隨便設，總和為1即可。

w'1+w'2 = 1

*w1

*w2

*w3

*w'1

*w'2

⋮ ⋮ 重複n次得到n張

original

model

Loss 2 (*λ)
計算三者的JS divergence

Ground truth

$$M = (p_{\mathrm{orig}} + p_{\mathrm{augmix1}} + p_{\mathrm{augmix2}})/3$$

$$\mathrm{JS}(p_{\mathrm{orig}}; p_{\mathrm{augmix1}}; p_{\mathrm{augmix2}}) = \frac{1}{3}\Big(\mathrm{KL}[p_{\mathrm{orig}}\|M] + \mathrm{KL}[p_{\mathrm{augmix1}}\|M] + \mathrm{KL}[p_{\mathrm{augmix2}}\|M]\Big)$$

Loss 1 $\mathcal{L}(p_{\mathrm{orig}}, y) + \lambda\,\mathrm{JS}(p_{\mathrm{orig}}; p_{\mathrm{augmix1}}; p_{\mathrm{augmix2}})$

Since the unlabeled data could give "information"

Using "data" for regularization would be a way to control the overparameterization models …

# What is happening to overparameterization?



使用over-parameterization的model，基礎的machine learning concept多會不管用。
For example:
L1, L2 regularizations。在基礎的機器學習中是透過降低參數量達到generalized purpose。但是在over-parameter mode上，降低很多參數仍然會讓model處於over-parameterized的狀態。因此可能不會太管用。

解決方式:
目前最好的regularization method仍然是使用training data來達到目的。

https://lilianweng.github.io/lil-log/2019/03/14/are-deep-neural-networks-dramatically-overfitted.html

# Contrastive learning

Objection 3

target

Objection 1

Objection 2



(a) Predictive learning

(b) Contrastive learning

If we use different objection for the same target, we would get different presenting which are belong to the same objects

We need to tell the neural network that "these are the same thing".

# Mean teacher

Already proposing the concepts of self-supervise (consistency regularization)



Figure 2: The Mean Teacher method. The figure depicts a training batch with a single labeled example. Both the student and the teacher model evaluate the input applying noise ($\eta$, $\eta'$) within their computation. The softmax output of the student model is compared with the one-hot label using classification cost and with the teacher output using consistency cost. After the weights of the student model have been updated with gradient descent, the teacher model weights are updated as an exponential moving average of the student weights. Both model outputs can be used for prediction, but 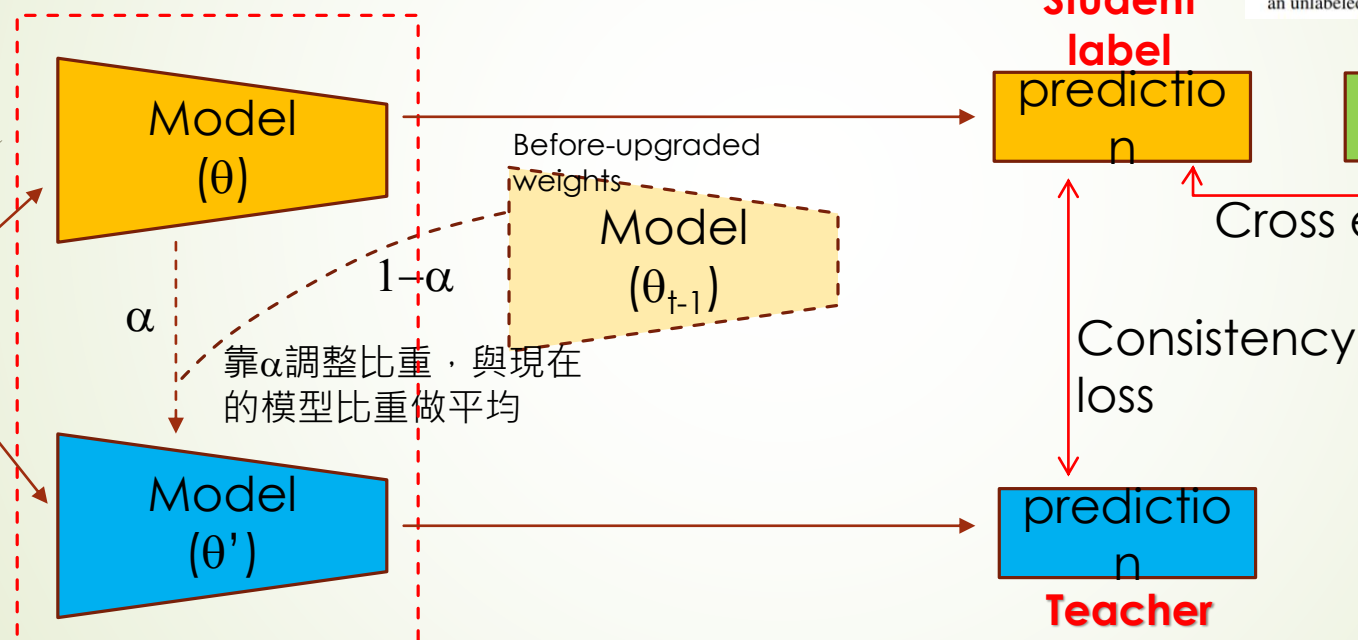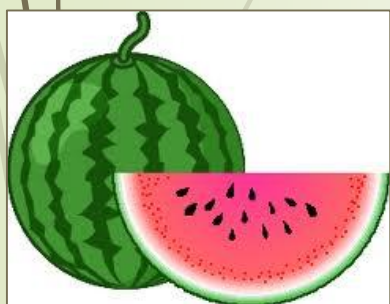at the end of the training the teacher prediction is more likely to be correct. A training step with an unlabeled example would be similar, except no classification cost would be applied.

https://arxiv.org/pdf/1703.01780.p

**Student label**

Model ($\theta$)

Before-upgraded weights

Model ($\theta_{t-1}$)

$\alpha$

$1-\alpha$

靠$\alpha$調整比重，與現在的模型比重做平均

prediction

Ground truth

Cross entropy

Consistency loss

Model ($\theta'$)

prediction

**Teacher label**

Applying the noise to the weights ($\eta$ and $\eta'$)
( Dropout would be a good

$\alpha$ would not to be large during initially training stage

The key concept:
1. Changing the "dark knowledge"
2. Extracting the ensemble characteristics.

# 自己跟自己比還不夠

透過跟別人比較，更了解自己的定位

Metric learning (meta-learning)

# Recurrent neural network

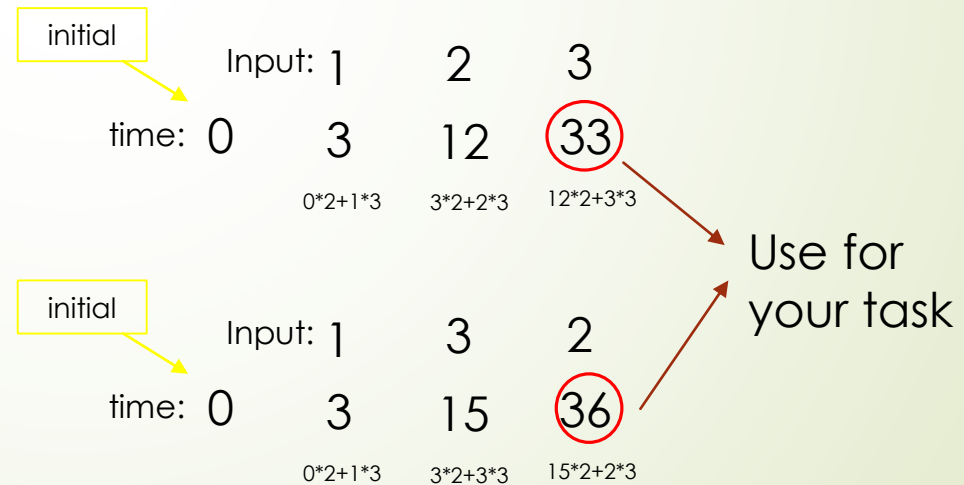- RNN( recurrent neural network )
  - Use the same tensor for sequence data
  - For example
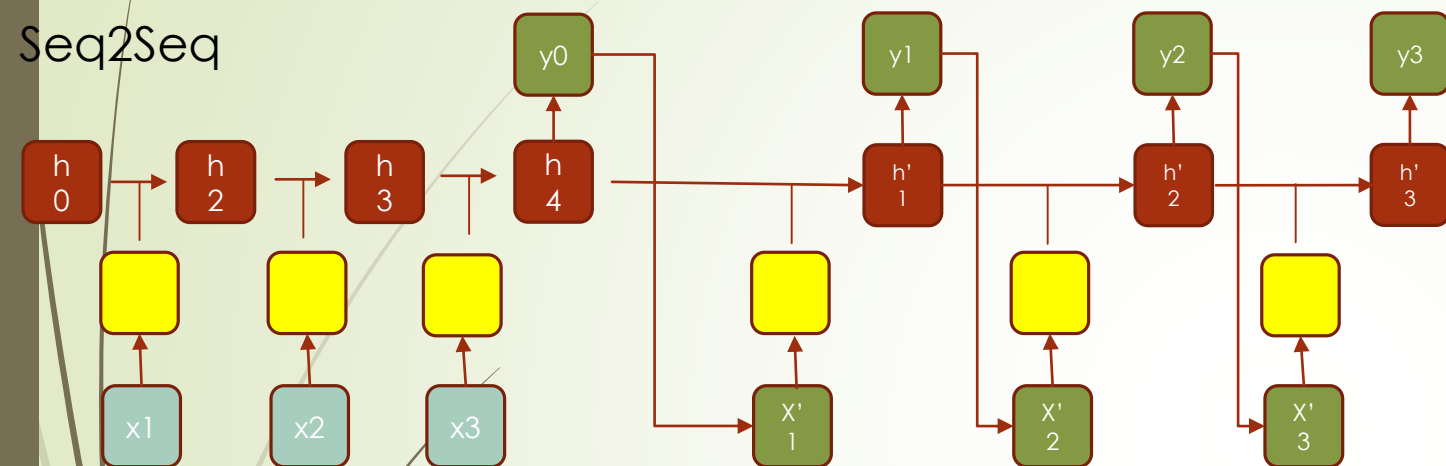    - time * 2
    - input * 3

Learning target

initial

Input: 1    2    3

time: 0    3    12    33

0*2+1*3    3*2+2*3    12*2+3*3

initial

Input: 1    3    2

time: 0    3    15    36

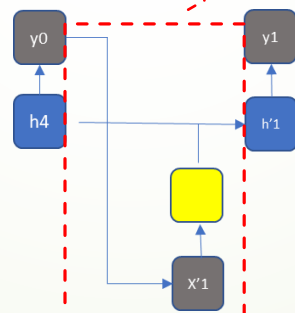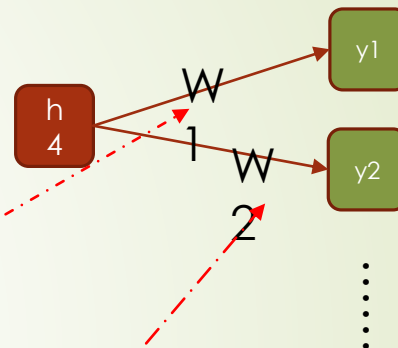0*2+1*3    3*2+3*3    15*2+2*3

Use for your task

2020/10/11

# Rethinking of RNN

Seq2Seq



1. h'系列全部都是由h4所得到，換句話說，
   h4早就隱含了h'系列的所有資訊
2. h'系列都是h4透過某些方式"打開"開關。
   例如h'1接收了x'1；h'2等於接收了x'1
   及x'2，以此類推。因此可以把這些"接
   收訊息"(例如x'1, x'1+x'2, …)當成"矩
   陣"，而這矩陣專門用來打開h4對於特
   定位置的開關。當h4接受到這些特定開
   關以後，就能把特定數值輸出即可。

Y1根本是自己
產生，只是透
過x'1打開自己

Y2是透過x'1和x'2的
訊息累加後做為開關。
但x'1和x'2也是由h4
自己產生

# Contrastive prediction code (CPC)

What is the "Contrastive"?
如果有"同一個"物件，我們用"不同角度"來看它，就能找出
不同處。

目前常用的方法就是加上noise或甚至直接乘以矩陣。

同一個h4，用不同的W處理，
同等于用不同角度在看h4。

因此把h4投射到不同位置，
等於在計算mutual
information。

假設從RNN部分取得
的資訊可以拆解(參
閱上一頁的投影片)



Figure 1: Overview of Contrastive Predictive Coding, the proposed representation learning approach.
Although this figure shows audio as input, we use the same setup for images, text and reinforcement
learning.

# SimCLR



representing

projection 只是為了訓練，最後可以丟掉

As close as possible (分子)

As far as possible (分母)

Use cos similarity in SimCLR

temperature

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\boldsymbol{z}_i, \boldsymbol{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\boldsymbol{z}_i, \boldsymbol{z}_k)/\tau)}, \quad (1)$$

Augmentation 1
Augmentation 2
Augmentation 3
Augmentation 4

Neural network

- Representation learning with contrastive cross entropy loss benefits from normalized embeddings and an appropriately adjusted temperature parameter.
- Contrastive learning benefits from larger batch sizes and longer training compared to its supervised counterpart. Like supervised learning, contrastive learning benefits from deeper and wider networks. https://arxiv.org/pdf/2002.05709.pdf

Tips
1. Because the sample size will be huge, LARS is recommended in such tasks.
2. Using layer normalization or global batch normalization.

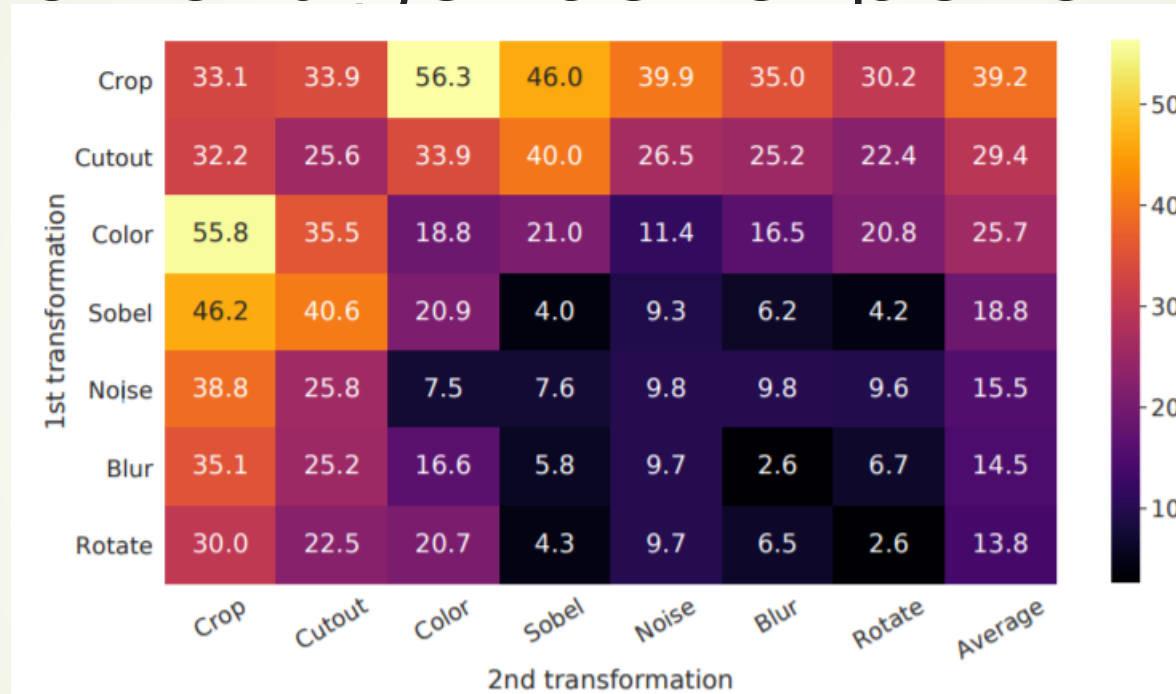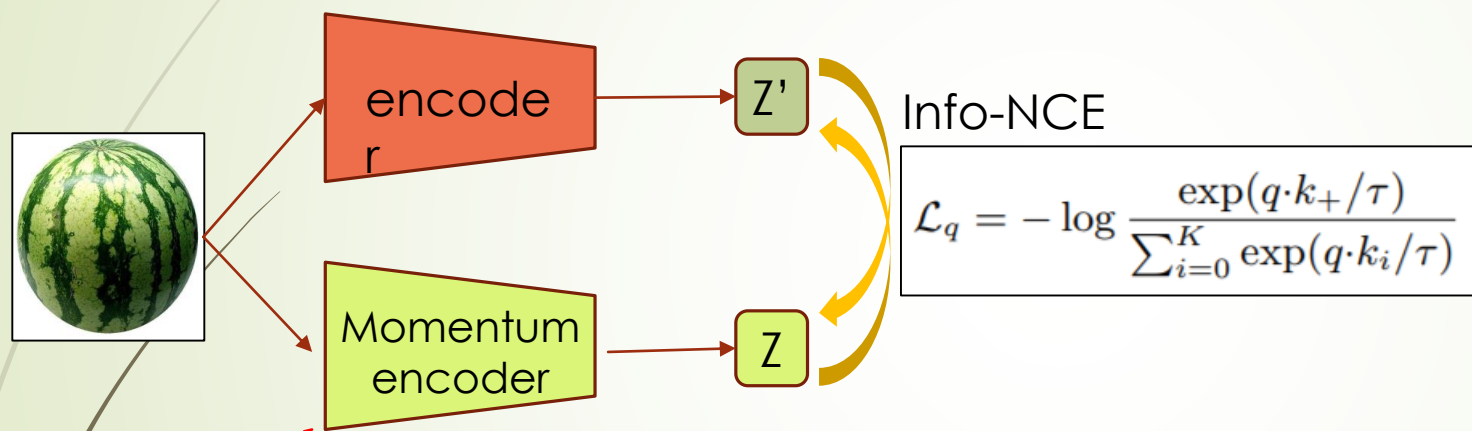# More in SimCLR – crop and color distortions get better performance



*Figure 5.* Linear evaluation (ImageNet top-1 accuracy) under individual or composition of data augmentations, applied only to one branch. For all columns but the last, diagonal entries correspond to single transformation, and off-diagonals correspond to composition of two transformations (applied sequentially). The last column reflects the average over the row.

# MOCO

Knowledge can be transfer directly

encoder

Momentum encoder

Z'

Z

Info-NCE

$$\mathcal{L}_q = -\log \frac{\exp(q \cdot k_+/\tau)}{\sum_{i=0}^{K} \exp(q \cdot k_i/\tau)}$$

$$\theta_{\mathrm{Z}} \leftarrow \alpha\theta_{\mathrm{Z}} + (1-\alpha)\theta_{\mathrm{Z}'}$$

直接使用Z'比重更新Z。
其內涵為:
1. 兩個不同比重對同一件事情都能完整表達，表示這樣的表達是正確的。
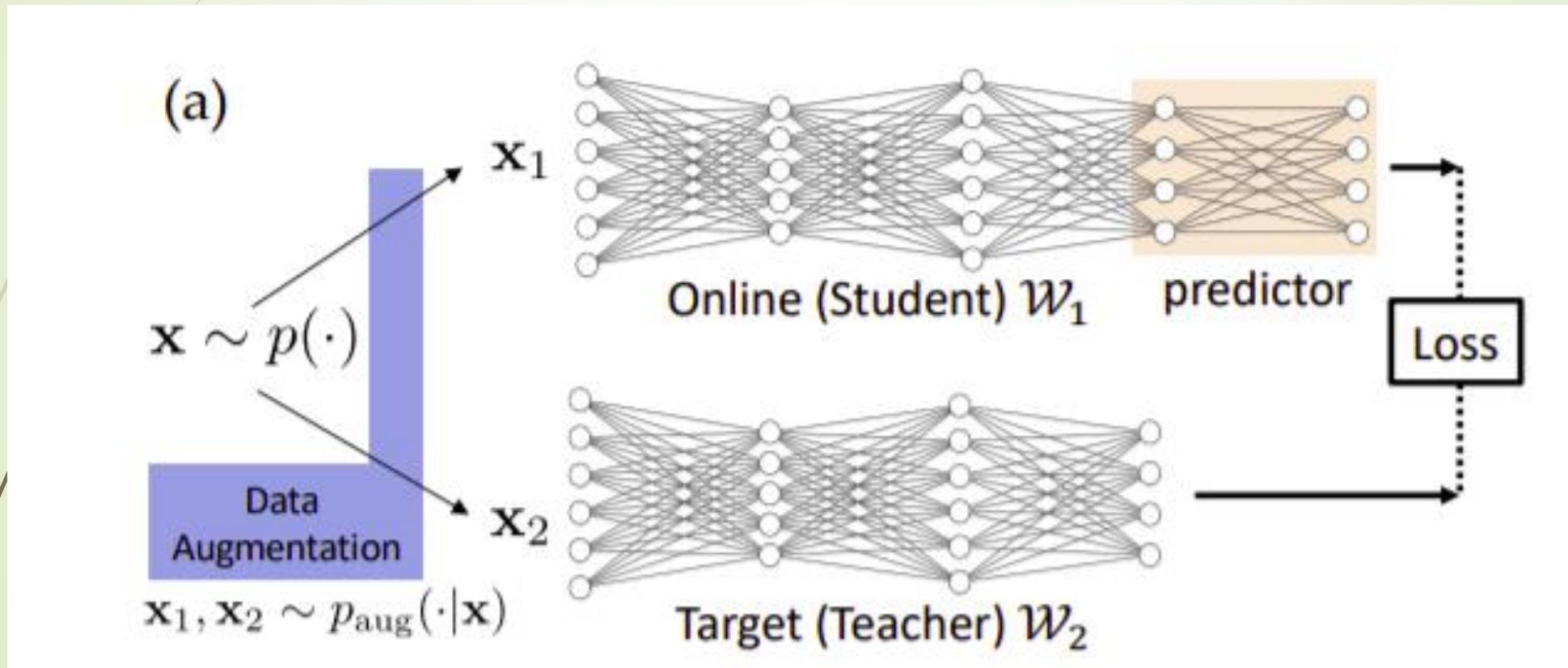2. 使用的info-NCE的並且搭配了temperature，表示在dark knowledge的部分模型也必須重視。

Tips:
- Batch normalization(BN)必須要進行shuffling。
因為使用info-NCE的狀況下，BN會洩漏正樣本跟負樣本之間的訊息造成模型只看BN輸出。

- α越大越好，即 momentum encoder更新幅度越小越好

## MOCO v2?
1. 補上projection (from SimCLR)
2. 補上augmentation (from SimCLR)

# Self-supervised learning could be thought as dual deep neural network

# Conclusion and suggestions

- Handling the dark knowledge would be a way to handling the overparameterization model

- In practical, we would need to learn some skill that the we might not notice in fashion way.

  - LARS ?

# Questions ?

If you like us, you can join our clubs: ( of course, they are all FREE)

1. 北Bio : https://www.facebook.com/groups/446434039038963

2. 台灣人工智慧社團: https://www.facebook.com/groups/525579498272187/

3. Deep learning 101 : https://www.facebook.com/TWMANDeepLearning/