

Unsupervised Learning of Visual Features by Contrasting Cluster Assignments

[Caron et al., NeurIPS 2020]

-Summary-

Introduction (SwAV)

- Current contrastive learning strategy :
 - Compares pairs (positive, true) – (negative, true) to pull (positive, true) and push (negative, true) .
 - Problem : computing all the pairs on a large dataset (ex : ImageNet) is not practical
 - Alternatives:
 - ① : compute contrastive loss based on random batch (widely adopted)
 - ② : **clustering-based methods** discriminating between groups of images with similar features instead of individual images (Caron et al., 2018)
- This paper tries to improve their previous work (②) by changing clustering-based methods.

Introduction (SwAV)

- SwAV (**Sw**apping **A**ssignments between multiple **V**iews of the same image) :
 - **Goal : To learn a useful prototype $\{c_1, \dots, c_K\}$ that discriminate feature vectors $\{z_1, \dots, z_B\}$, which goes along with the concept of CL by using swapping assignments.**
 - Method :
 1. Image x_n is transformed into an augmented view x_{nt} by transformation $t \sim \mathcal{T}$.
 2. x_t is mapped into normalized feature vector $z_{nt} = \frac{f_\theta(x_{nt})}{\|f_\theta(x_{nt})\|_2} \in \mathbb{R}^D$ by model f_θ .
 3. Then, we compute a '**code**' q_{nt} by exploiting 'optimal transport' between z_{nt} and prototype vectors $\{c_1, \dots, c_K\} \subset \mathbb{R}^D$

Note : prototype matrix $C = [c_1, \dots, c_K] \in \mathbb{R}^{D \times K}$, code matrix $Q = [q_1, \dots, q_B] \in \mathbb{R}^{K \times B}$

(where D = dimension of feature z_{nt} , B = batch size, K = the number of polytopes)

Introduction (SwAV)

©: <https://www.youtube.com/watch?v=7QmsTleiRLs>

- Recall that SwAV clusters multiple samples to the prototypes.

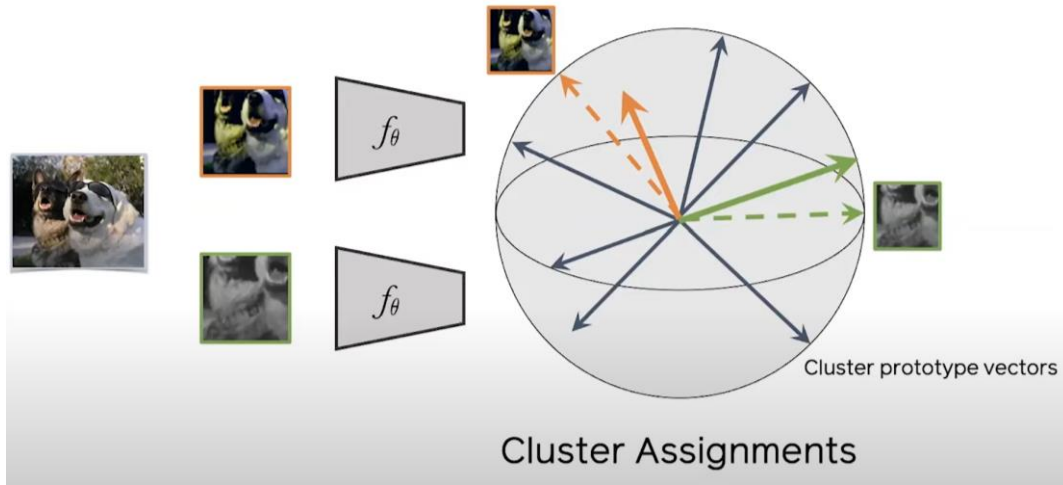


Illustration of clustering on SwAV

Note : why we use optimal transportation to get code q_{nt} ?

- q_{nt} should be probability vector (used as loss)
(achieved by constraint in optimal transportation)**
- Higher coupling value for positive pair / Smaller coupling value for negative pair**

Above objectives can be achieved by optimal transportation by setting $-C^T Z$ as the cost matrix.

- But how to exploit optimal transportation to compute the codes q_{nt} ?
 - Set code matrix Q as the coupling matrix in optimal transportation.
 - Set $-C^T Z$ as the cost matrix in optimal transportation. (where $Z = [z_1, \dots, z_B] \in \mathbb{R}^{D \times B}$)
 - Now, we want to minimize the coupling \times cost entirely (= solving optimal transportation)

Introduction (SwAV)

- Why we choose $-C^T Z = \begin{bmatrix} -c_1^T z_1 & \cdots & -c_B^T z_1 \\ \vdots & \ddots & \vdots \\ -c_K^T z_1 & \cdots & -c_K^T z_B \end{bmatrix} \in \mathbb{R}^{K \times B}$ as the cost matrix ?
 - This enforces the optimal transportation to assign high coupling value ($= q_{ij}$) for positive pairs.
Logic : Positive pair \rightarrow high similarity ($= c^T z$) \rightarrow low cost ($= -c^T z$) \rightarrow high coupling value assigned by optimal transportation algorithm
 - Also, this method assign low coupling value ($= q_{ij}$) for negative pair in the similar logic above.
- This optimization maximizes the similarity between the features and the prototypes.

Introduction (SwAV)

- To prevent trivial solution such as giving all coupling values to smallest cost-valued item, we add entropy term $H(Q) := -\sum_{i,j} Q_{ij} \log Q_{ij}$.
- To sum up, our optimization problem is following :

$$\max_{Q \in \mathcal{Q}} \langle Q, C^T Z \rangle_F + \epsilon H(Q)$$

Note (Frobenius inner product) :

$$\langle A, B \rangle_F = \sum_{i,j} A_{ij} B_{ij} = \text{Tr}(A^T B)$$

where $\mathcal{Q} = \left\{ Q \in \mathbb{R}_+^{K \times B} \mid Q \cdot \mathbf{1}_B = \frac{1}{k} \mathbf{1}_K, Q^T \cdot \mathbf{1}_K = \frac{1}{B} \mathbf{1}_B \right\}$ (finite resource constraint in optimal transport.)

- Luckily, there is a good algorithm to solve this problem (**Sinkhorn-Knopp algorithm**) :
 - Solution : $Q^* = \text{diag}(u) \cdot K \cdot \text{diag}(v)$

$$\text{where } K = \exp\left(\frac{C^T Z}{\epsilon}\right), \quad u^{(k+1)} = \frac{r}{K V^{(k)}}, \quad v^{(k+1)} = \frac{c}{K^T u^{(k+1)}} \quad (\text{here, } r = \frac{1}{k} \mathbf{1}_K, c = \frac{1}{B} \mathbf{1}_B)$$

Introduction (SwAV)

- (Back to) Method :

Note : Overall loss adopting $L(z_t, z_s)$ can be rewritten as follows : (N : # of data)

$$-\frac{1}{N} \sum_{n=1}^N \sum_{s,t \sim \mathcal{T}} \left[\frac{1}{\tau} \mathbf{z}_{nt}^\top \mathbf{C} \mathbf{q}_{ns} + \frac{1}{\tau} \mathbf{z}_{ns}^\top \mathbf{C} \mathbf{q}_{nt} - \log \sum_{k=1}^K \exp \left(\frac{\mathbf{z}_{nt}^\top \mathbf{c}_k}{\tau} \right) - \log \sum_{k=1}^K \exp \left(\frac{\mathbf{z}_{ns}^\top \mathbf{c}_k}{\tau} \right) \right]$$

4. Calculate 'estimated similarity vector p_{nt} ' of z_{nt} based on the previous projection:

$$p_{nt}^{(k)} = \frac{\exp\left(\frac{1}{\tau} \mathbf{z}_{nt}^\top \mathbf{c}_k\right)}{\sum_{k'} \exp\left(\frac{1}{\tau} \mathbf{z}_{nt}^\top \mathbf{c}_{k'}\right)}$$

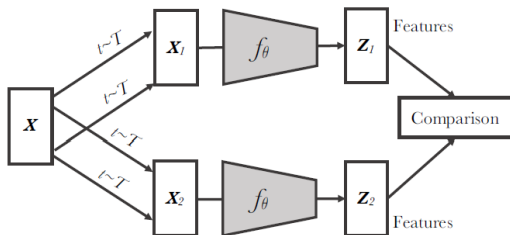
Softmax of columns of $\mathbf{z}_{nt}^\top \mathbf{C}$ with temperature τ

Why the loss is based on swapped prediction?

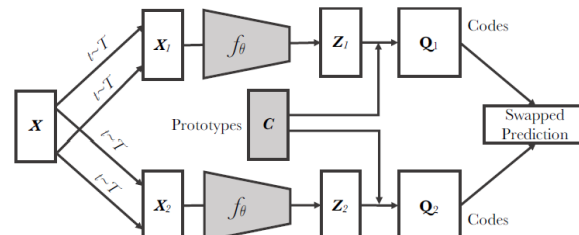
Idea : $\mathbf{z}_t, \mathbf{z}_s$ should be assigned into same prototype.

5. Calculate loss based on **swapped prediction** : (Note : $\mathbf{z}_t, \mathbf{z}_s$ are two different augmentations of x)

$$L(z_{nt}, z_{ns}) = l(z_{nt}, q_{ns}) + l(z_{ns}, q_{nt}), \quad l(z_{nt}, q_{ns}) = -\sum_k q_{ns}^{(k)} \log p_{nt}^{(k)}$$



Contrastive instance learning



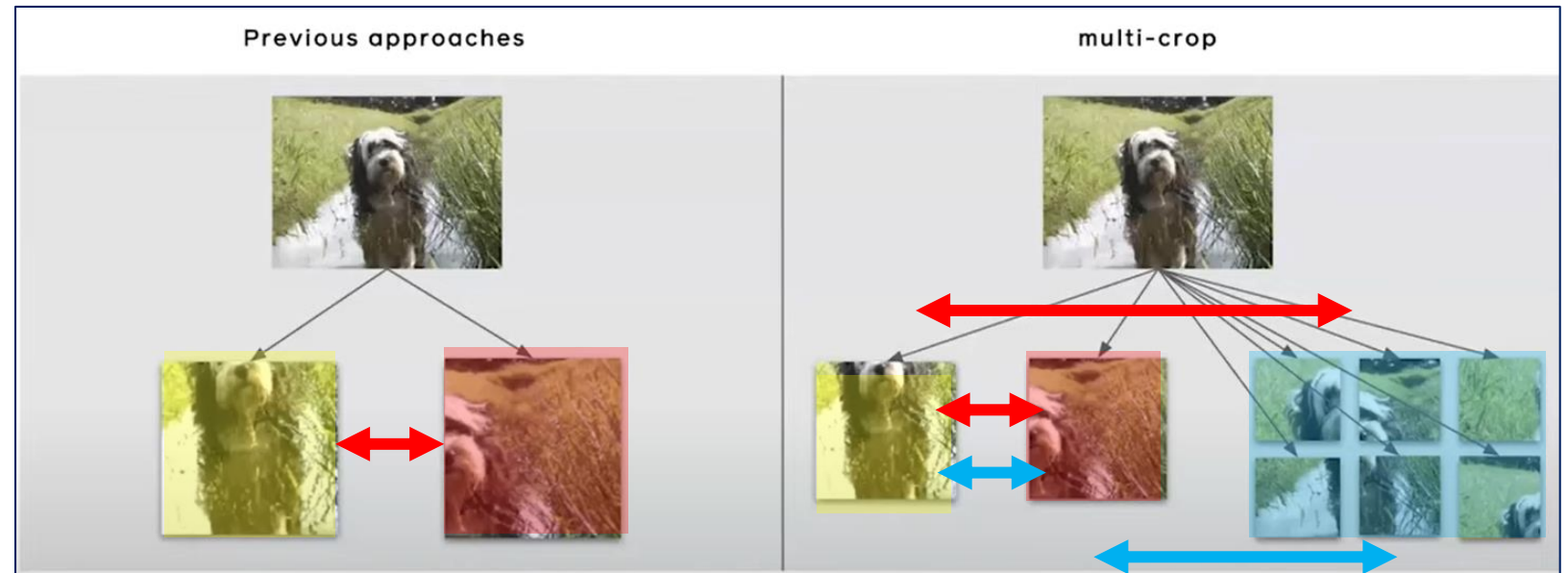
Swapping Assignments between Views

Comparison between CL instance learning vs. SwAV
(Note : Gray color implies trainable parameters)

Introduction (Multi-crop)

- It is widely known that comparing more views during CL training improves the resulting model. (Misra et al., 2019)
- However, most CL methods compare 'one' pair of transformations per image for saving comp. cost.
 - To exploit the advantage of multi view training, they propose to use V low resolution crops that cover only small parts of the image in addition to two standard resolution crops.

Illustration of multi-crop concept



Introduction (Multi-crop)

- Detailed description of multi-crop :
 - Recall the loss based on **swapped prediction** appeared in SwAV.

$$L(z_{t_1}, z_{t_2}) = l(z_{t_1}, q_{t_2}) + l(z_{t_2}, q_{t_1})$$

which adopts only two standard resolution crops z_{t_1}, z_{t_2} .

- We now sample V additional low resolution crops $z_{t_3}, \dots, z_{t_{V+2}}$ and generalize the loss :

$$L(z_{t_1}, z_{t_2}, \dots, z_{t_{V+2}}) = \sum_{i \in \{1, 2\}} \sum_{v=1}^{V+2} 1_{v \neq i} l(z_{t_v}, q_{t_i})$$

- Note : Using low resolution images ensures only a small increase in computational cost.

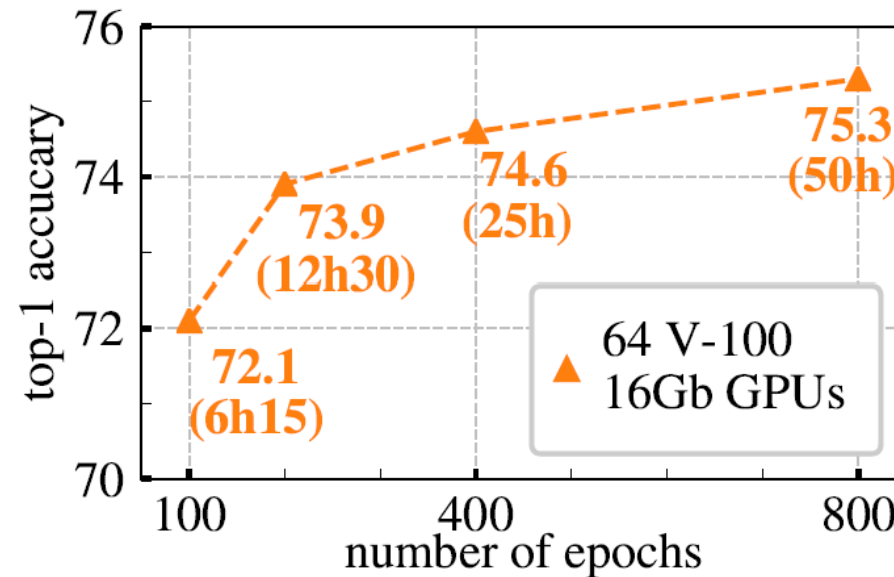
Experiments

Multi-crop (x)

Multi-crop (o)

- To verify downstream performance, they used linear evaluation protocol.

Method	Top-1		Δ
	2x224	2x160+4x96	
Supervised	76.5	76.0	-0.5
<i>Contrastive-instance approaches</i>			
SimCLR	68.2	70.6	+2.4
<i>Clustering-based approaches</i>			
SeLa-v2	67.2	71.8	+4.6
DeepCluster-v2	70.2	74.3	+4.1
SwAV	70.1	74.1	+4.0



Left : Comparison between clustering-based and contrastive instance methods and impact of multi-crop. (400 /200 epochs for unsupervised / supervised)

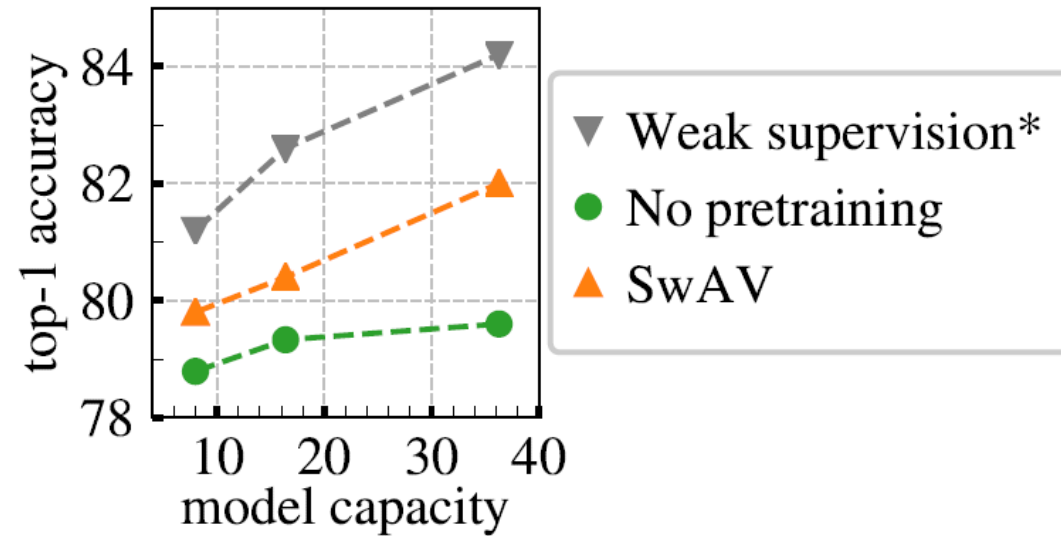
Right : SwAV (with multi-crop) downstream performance as a function of epochs

Note : DeepCluster-v2 is not online which makes it impractical for extremely large datasets (= requires much longer training compared to SwAV)

Experiments

- Further experiments demonstrate the SwAV's higher performance on representation learning.

Method	Frozen	Finetuned
Random	15.0	76.5
MoCo	-	77.3*
SimCLR	60.4	77.2
SwAV	66.5	77.8



Left : Top-1 acc on ImageNet for pretrained models on uncurated 1B Instagram images (linear evaluation on frozen features or finetuned features)

Right : Performance of finetuned models as we increase the capacity of ResNext

Note (weak supervision) : pretrained on a curated 1B Instagram images filtered with 1.5k hashtags similar to ImageNet classes.