IMPERIAL

Unsupervised Training of Vision Transformers with Synthetic Negatives Nikolaos Giakoumoglou, Andreas Floros, Kleanthis Marios Papadopoulos, Tania Stathaki **Department of Electrical and Electronic Engineering, Imperial College London**

Introduction

This paper addresses the neglected potential of hard negative samples in self-supervised learning for vision transformers.

- Previous works explored synthetic hard negatives but not in the context of vision transformers
- We build on this observation and integrate synthetic hard negatives generated "on-the-fly" to improve vision transformer representation learning
- This simple yet effective technique notably improves the discriminative power of learned representations
- Our experiments show performance improvements for both **DeiT-S and Swin-T architectures**

"If intelligence is a cake, the bulk of the cake is unsupervised learning" - Yann LeCun

Preliminaries

Contrastive Learning

Contrastive learning aims to differentiate between similar and dissimilar data pairs. Given an image x, and two augmented views, we minimize the InfoNCE loss:

$$\begin{split} \mathcal{L}(\mathbf{q},\mathbf{k},\mathcal{Q}) = -\log\frac{\exp(\mathbf{q}^\top\cdot\mathbf{k}/\tau)}{\exp(\mathbf{q}^\top\cdot\mathbf{k}/\tau) + \sum\limits_{\mathbf{n}\in\mathcal{Q}}\exp(\mathbf{q}^\top\cdot\mathbf{n}/\tau)} \end{split}$$

where **q** and **k** are the representations of two augmented views, Q is a set of negative samples, and τ is a temperature parameter.

Methodology

Synthetic Hard Negatives Synthetic hard negatives provide challenging examples that help models learn more discriminative features. We generate them using an arbitrary function \mathcal{G} (see SynCo [2] for details):

where \mathcal{Q}^N are the top N hardest negatives and ξ controls the synthesis process.

The modified InfoNCE loss is given by:

The first sum represents memory-based negatives and the second sum represents synthetic negatives.

We develop our approach in PyTorch, building upon the implementation of MoBY [1] and SynCo [2]. We pretrain on ImageNet ILSVRC-2012 using a DeiT-Small or Swin-Tiny encoder. We refer to our method as SynBY.

Imperial College London

$$\mathbf{s} = \frac{\mathcal{G}(\mathbf{q}, \mathcal{Q}^N; \xi)}{\|\mathcal{G}(\mathbf{q}, \mathcal{Q}^N; \xi)\|_2}$$

$$\mathcal{L}(\mathbf{q},\mathbf{k},\mathcal{Q},\mathcal{S}) = -\log\frac{\exp(\mathbf{q}^{\top}\cdot\mathbf{k}/\tau)}{\exp(\mathbf{q}^{\top}\cdot\mathbf{k}/\tau) + Z}$$

where Z represents the negative samples:

$$Z = \sum_{\mathbf{n} \in \mathcal{Q}} \exp(\mathbf{q}^\top \cdot \mathbf{n} / \tau) + \sum_{\mathbf{s} \in \mathcal{S}} \exp(\mathbf{q}^\top \cdot \mathbf{s} / \tau)$$

Experiments

Linear Evaluation

Method
Supervised Supervised
DINO MoCo-v3
MoBY MoBY
SynBY (ours) SynBY (ours)

Attention Visualization



References

[1] Zhenda Xie, Yutong Lin, Zhuliang Yao, Zheng Zhang, Qi Dai, Yue Cao, and Han Hu. Self-supervised Learning with Swin Transformers, 2021.

[2] Nikolaos Giakoumoglou and Tania Stathaki. SynCo: Synthetic Hard Negatives for Contrastive Visual Representation Learning, 2025.



on ImageNet		
Architecture	Params (M)	Top-1 (%)
DeiT-S	22	79.8
Swin-T	29	81.3
DeiT-S	22	72.5
DeiT-S	22	72.5
DeiT-S	22	72.8
Swin-T	29	75.0
DeiT-S	22	73.0
Swin-T	29	75.2

