

---

# Bootstrap Your Own Latent

## A New Approach to Self-Supervised Learning

---

School of Industrial and Management Engineering, Korea University

Hansam Cho

# Contents

---

- ❖ Research Purpose
- ❖ Methods
- ❖ Experiments

# Research Purpose

---

- ❖ Bootstrap Your Own Latent: A New Approach to Self-Supervised Learning (2020, NeurIPS)
  - DeepMind 소속의 저자, 2022년 06월 06일 기준으로 1438회 인용

---

## Bootstrap Your Own Latent A New Approach to Self-Supervised Learning

---

Jean-Bastien Grill<sup>\*1</sup>, Florian Strub<sup>\*1</sup>, Florent Altché<sup>\*1</sup>, Corentin Tallec<sup>\*1</sup>, Pierre H. Richemond<sup>\*1,2</sup>

Elena Buchatskaya<sup>1</sup>, Carl Doersch<sup>1</sup>, Bernardo Avila Pires<sup>1</sup>, Zhaohan Daniel Guo<sup>1</sup>

Mohammad Gheshlaghi Azar<sup>1</sup>, Bilal Piot<sup>1</sup>, Koray Kavukcuoglu<sup>1</sup>, Rémi Munos<sup>1</sup>, Michal Valko<sup>1</sup>

<sup>1</sup>DeepMind

<sup>2</sup>Imperial College

[jbgrill, fstrub, altche, corentint, richemond]@google.com

### Abstract

We introduce Bootstrap Your Own Latent (BYOL), a new approach to self-supervised image representation learning. BYOL relies on two neural networks, referred to as *online* and *target* networks, that interact and learn from each other. From an augmented view of an image, we train the online network to predict the target network representation of the same image under a different augmented view. At the same time, we update the target network with a slow-moving average of the online network. While state-of-the-art methods rely on negative pairs, BYOL achieves a new state of the art *without them*. BYOL reaches 74.3% top-1 classification accuracy on ImageNet using a linear evaluation with a ResNet-50 architecture and 79.6% with a larger ResNet. We show that BYOL performs on par or better than the current state of the art on both transfer and semi-supervised benchmarks. Our implementation and pretrained models are given on GitHub.<sup>3</sup>

# Research Purpose

---

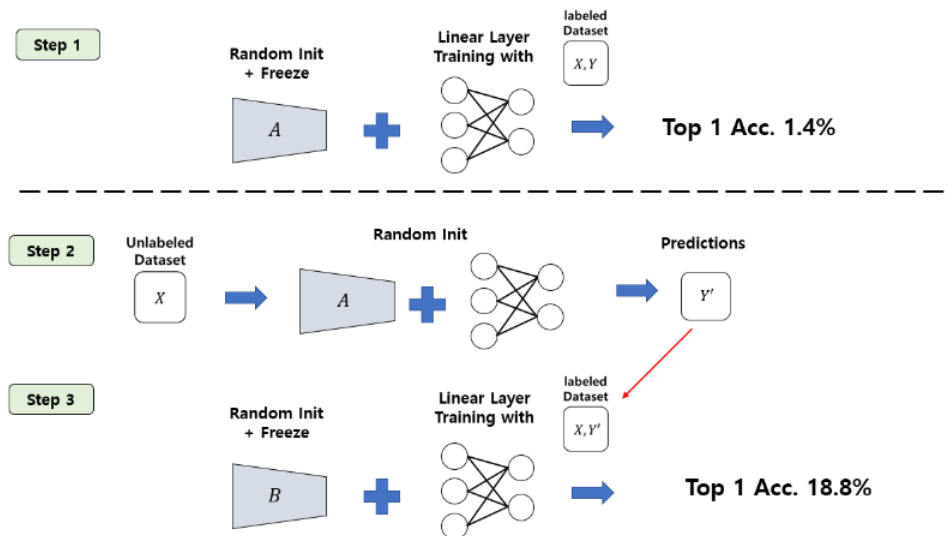
## ❖ Introduction

- 좋은 representation을 학습하는 것은 중요한 task이며 contrastive learning이 좋은 성능을 보임
- Contrastive learning은 negative pair의 정의와 augmentation 기법에 민감함
- Positive pair만을 활용하는 BYOL이라는 새로운 self-supervised learning 기법을 제안함

# Methods

## ❖ Motivation

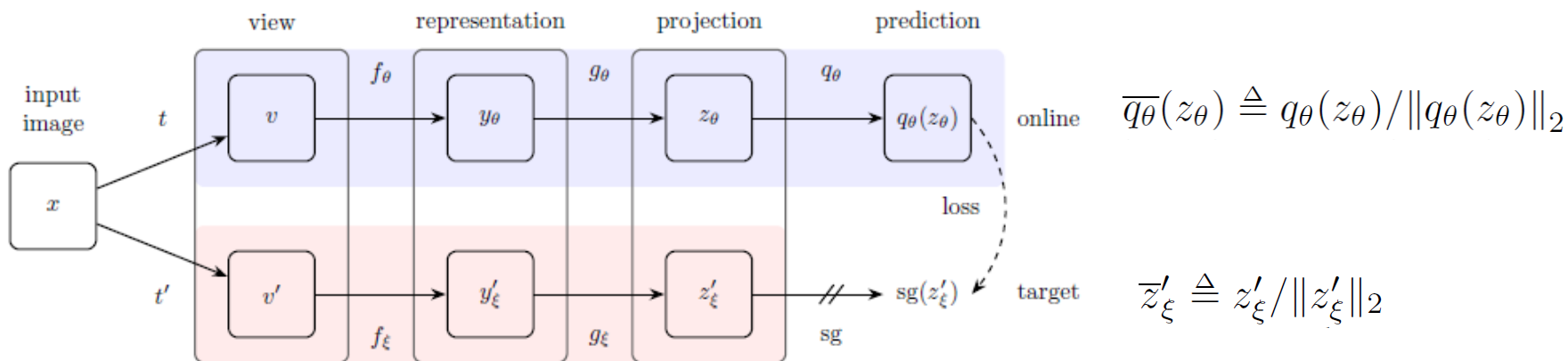
- ImageNet 데이터 활용한 classification 성능 비교
- 실험1: Random representation을 활용, linear classifier 학습 → Top1 Acc. 1.4%
- 실험2: Random representation을 예측하도록 encoder 학습 후 해당 encoder를 고정, linear classifier 학습 → Top1 Acc. 18.8%
- Representation을 예측하는 것 만으로도 좋은 representation 학습 가능



# Methods

## ❖ BYOL

- Online network: 실제 representation을 학습하는 모델
- Target network: Online network가 학습할 target값을 생성하는 모델
- Online network에만 predictor가 존재 → 모델의 비대칭성을 통해 서로 다른 representation 생성
- 서로 다른 augmentation ( $t, t'$ )을 활용해 두개의 이미지 생성 ( $v, v'$ ) 후 입력
- Online network와 target network의 output에 대해 정규화



# Methods

---

## ❖ BYOL

- MSE를 loss로 활용

$$\mathcal{L}_{\theta,\xi} \triangleq \|\overline{q_{\theta}}(z_{\theta}) - \overline{z'_{\xi}}\|_2^2 = 2 - 2 \cdot \frac{\langle q_{\theta}(z_{\theta}), z'_{\xi} \rangle}{\|q_{\theta}(z_{\theta})\|_2 \cdot \|z'_{\xi}\|_2}.$$

- 입력 이미지를 바꾼 경우에 대해서도 loss 계산, 두 loss를 더해서 최종 loss로 활용

$$\mathcal{L}_{\theta,\xi}^{\text{BYOL}} = \mathcal{L}_{\theta,\xi} + \tilde{\mathcal{L}}_{\theta,\xi}$$

- Online network만 학습을 진행하며 target network는 online network의 exponential moving average로 파라미터 업데이트

$$\theta \leftarrow \text{optimizer}(\theta, \nabla_{\theta} \mathcal{L}_{\theta,\xi}^{\text{BYOL}}, \eta) \quad \text{and} \quad \xi \leftarrow \tau \xi + (1 - \tau) \theta,$$

# Methods

## ❖ Intuitions on BYOL's behavior

- BYOL은 collapse\*를 막을 수 있는 장치가 없음 (\*collapse: 모든 이미지에 대해서 동일한 representation 생성)
- 하지만 collapse가 일어나는 경우가 불안정하기 때문에 collapse가 발생하지 않음
- Online network predictor가 optimal인 상황을 가정

$$q_\theta = q^\star \text{ with } q^\star \triangleq \arg \min_q \mathbb{E} \left[ \|q(z_\theta) - z'_\xi\|_2^2 \right], \quad \text{where } q^\star(z_\theta) = \mathbb{E}[z'_\xi | z_\theta],$$

- Online network update에 사용되는 gradient는 conditional variance에 영향을 받음

$$\nabla_\theta \mathbb{E} \left[ \|q^\star(z_\theta) - z'_\xi\|_2^2 \right] = \nabla_\theta \mathbb{E} \left[ \|\mathbb{E}[z'_\xi | z_\theta] - z'_\xi\|_2^2 \right] = \nabla_\theta \mathbb{E} \left[ \sum_i \text{Var}(z'_{\xi,i} | z_\theta) \right]$$

- Collapse 상황의 variance가 일반적인 상황의 variance에 비해 큼

$$\text{Var}(z'_\xi | z_\theta) \leq \text{Var}(z'_\xi | c)$$

- 따라서 학습과정 중 collapse 상황을 쉽게 벗어날 수 있음



# Experiments

## ❖ Linear evaluation on ImageNet

- BYOL은 online network encoder를 활용해 representation 생성
- Encoder를 고정 후 생성된 representation을 활용해 liner classifier 학습
- BYOL이 다른 방법론에 비해 좋은 성능을 보임

Method	Top-1	Top-5
Local Agg.	60.2	-
PIRL [35]	63.6	-
CPC v2 [32]	63.8	85.3
CMC [11]	66.2	87.0
SimCLR [8]	69.3	89.0
MoCo v2 [37]	71.1	-
InfoMin Aug. [12]	73.0	91.1
BYOL (ours)	<b>74.3</b>	<b>91.6</b>

(a) ResNet-50 encoder.

Method	Architecture	Param.	Top-1	Top-5
SimCLR [8]	ResNet-50 (2×)	94M	74.2	92.0
CMC [11]	ResNet-50 (2×)	94M	70.6	89.7
BYOL (ours)	ResNet-50 (2×)	94M	<b>77.4</b>	<b>93.6</b>
CPC v2 [32]	ResNet-161	305M	71.5	90.1
MoCo [9]	ResNet-50 (4×)	375M	68.6	-
SimCLR [8]	ResNet-50 (4×)	375M	76.5	93.2
BYOL (ours)	ResNet-50 (4×)	375M	<b>78.6</b>	<b>94.2</b>
BYOL (ours)	ResNet-200 (2×)	250M	<b>79.6</b>	<b>94.8</b>

(b) Other ResNet encoder architectures.

Table 1: Top-1 and top-5 accuracies (in %) under linear evaluation on ImageNet.

# Experiments

## ❖ Semi-supervised training on ImageNet

- 일부 데이터의 label 정보만 활용한 semi-supervised 실험, BYOL의 성능이 가장 우수

Method	Top-1		Top-5	
	1%	10%	1%	10%
Supervised [77]	25.4	56.4	48.4	80.4
InstDisc	-	-	39.2	77.4
PIRL [35]	-	-	57.2	83.8
SimCLR [8]	48.3	65.6	75.5	87.8
BYOL (ours)	<b>53.2</b>	<b>68.8</b>	<b>78.4</b>	<b>89.0</b>

(a) ResNet-50 encoder.

Method	Architecture	Param.	Top-1		Top-5	
			1%	10%	1%	10%
CPC v2 [32]	ResNet-161	305M	-	-	77.9	91.2
SimCLR [8]	ResNet-50 (2×)	94M	58.5	71.7	83.0	91.2
BYOL (ours)	ResNet-50 (2×)	94M	<b>62.2</b>	<b>73.5</b>	<b>84.1</b>	<b>91.7</b>
SimCLR [8]	ResNet-50 (4×)	375M	63.0	74.4	85.8	92.6
BYOL (ours)	ResNet-50 (4×)	375M	<b>69.1</b>	<b>75.7</b>	<b>87.9</b>	<b>92.5</b>
BYOL (ours)	ResNet-200 (2×)	250M	<b>71.2</b>	<b>77.7</b>	<b>89.5</b>	<b>93.7</b>

(b) Other ResNet encoder architectures.

## ❖ Transfer to other classification tasks

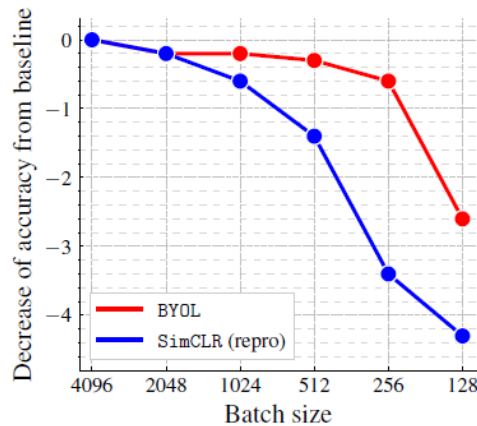
- 다양한 데이터에 적용할 수 있는 일반적인 representation을 생성할 수 있는지에 대한 실험
- Supervised learning을 통해 생성된 representation보다 좋은 성능을 보이기도 함

Method	Food101	CIFAR10	CIFAR100	Birdsnap	SUN397	Cars	Aircraft	VOC2007	DTD	Pets	Caltech-101	Flowers
<i>Linear evaluation:</i>												
BYOL (ours)	<b>75.3</b>	<b>91.3</b>	<b>78.4</b>	<b>57.2</b>	<b>62.2</b>	<b>67.8</b>	<b>60.6</b>	<b>82.5</b>	<b>75.5</b>	<b>90.4</b>	<b>94.2</b>	<b>96.1</b>
SimCLR (repro)	72.8	90.5	74.4	42.4	60.6	49.3	49.8	81.4	<b>75.7</b>	84.6	89.3	92.6
SimCLR [8]	68.4	90.6	71.6	37.4	58.8	50.3	50.3	80.5	74.5	83.6	90.3	91.2
Supervised-IN [8]	72.3	<b>93.6</b>	78.3	53.7	61.9	66.7	<b>61.0</b>	<b>82.8</b>	74.9	<b>91.5</b>	<b>94.5</b>	94.7
<i>Fine-tuned:</i>												
BYOL (ours)	<b>88.5</b>	<b>97.8</b>	<b>86.1</b>	<b>76.3</b>	<b>63.7</b>	<b>91.6</b>	<b>88.1</b>	<b>85.4</b>	<b>76.2</b>	<b>91.7</b>	<b>93.8</b>	<b>97.0</b>
SimCLR (repro)	87.5	97.4	85.3	75.0	63.9	91.4	87.6	84.5	75.4	89.4	91.7	96.6
SimCLR [8]	88.2	97.7	85.9	75.9	63.5	91.3	88.1	84.1	73.2	89.2	92.1	97.0
Supervised-IN [8]	88.3	97.5	<b>86.4</b>	<b>75.8</b>	<b>64.3</b>	<b>92.1</b>	86.0	85.0	74.6	<b>92.1</b>	<b>93.3</b>	<b>97.6</b>
Random init [8]	86.9	95.9	80.2	76.1	53.6	91.4	85.9	67.3	64.8	81.5	72.6	92.0

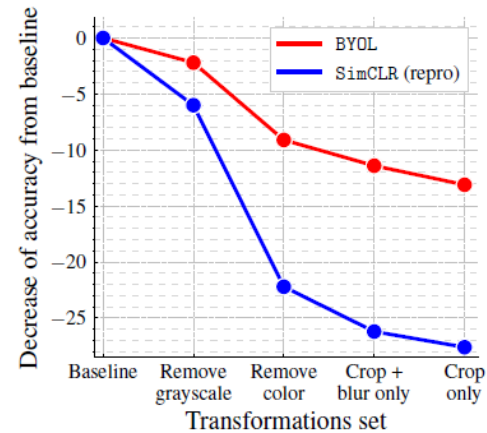
# Experiments

## ❖ Building intuitions with ablations

- SimCLR에 비해 batch size 크기, augmentation 기법에 대해 강건한 성능을 보임



(a) Impact of batch size



(b) Impact of progressively removing transformations

*Thank You*