# A Simple Framework for Contrastive Learning of Visual Representations

Ting Chen, Simon Kornblith, Mohammad Norouzi, Geoffrey Hinton

ICML (?)

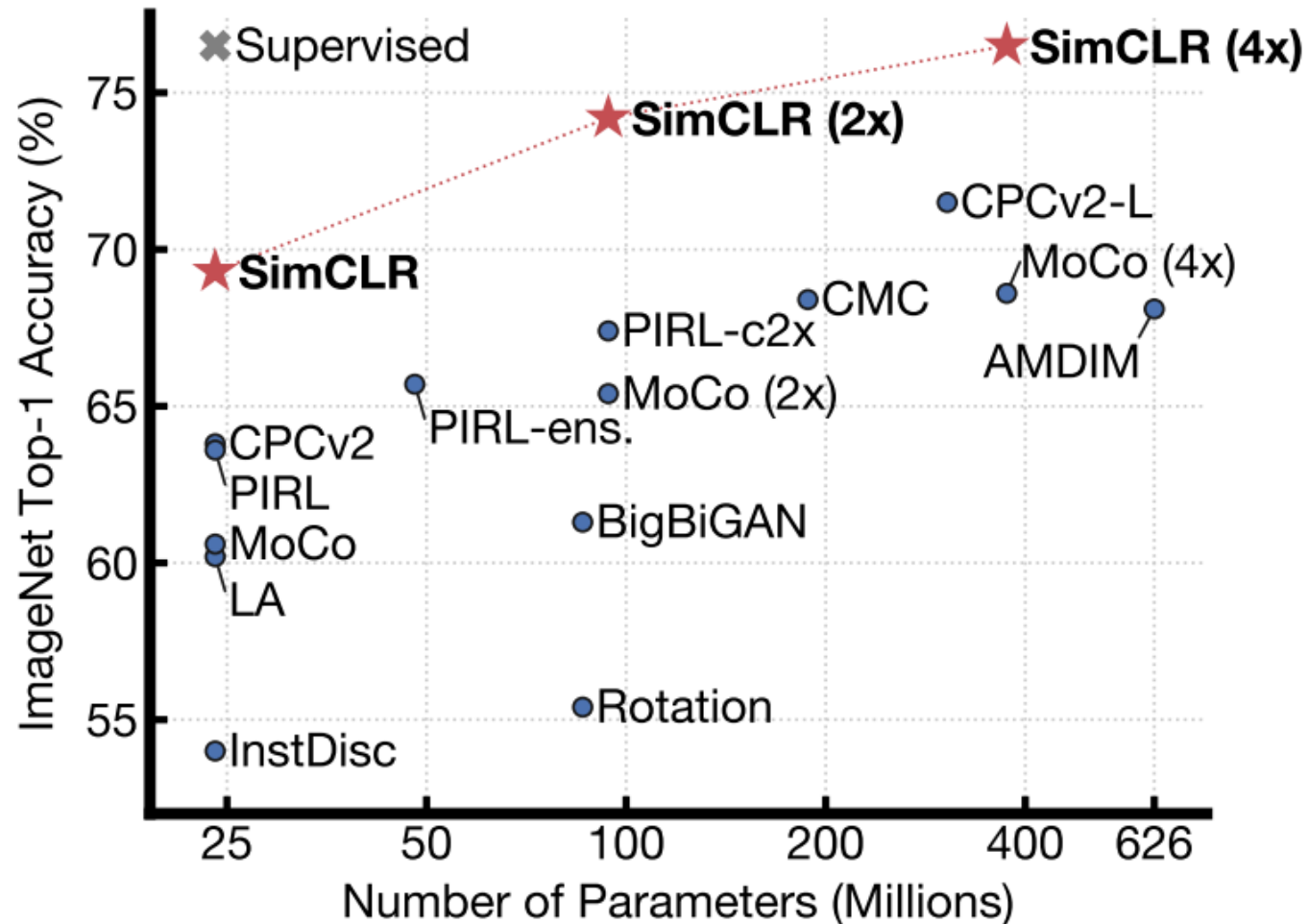2020. 04. 02

Presented by Junsoo Lee

# Overview:

- This paper simplify recently proposed contrastive self-supervised learning algorithms without requiring <u>specialized architectures</u> or a <u>memory bank</u>.

- Task: the contrastive prediction task.

- Purpose: to learn visual representation.

- Presentation: Studying the major components of proposed framework systematically.
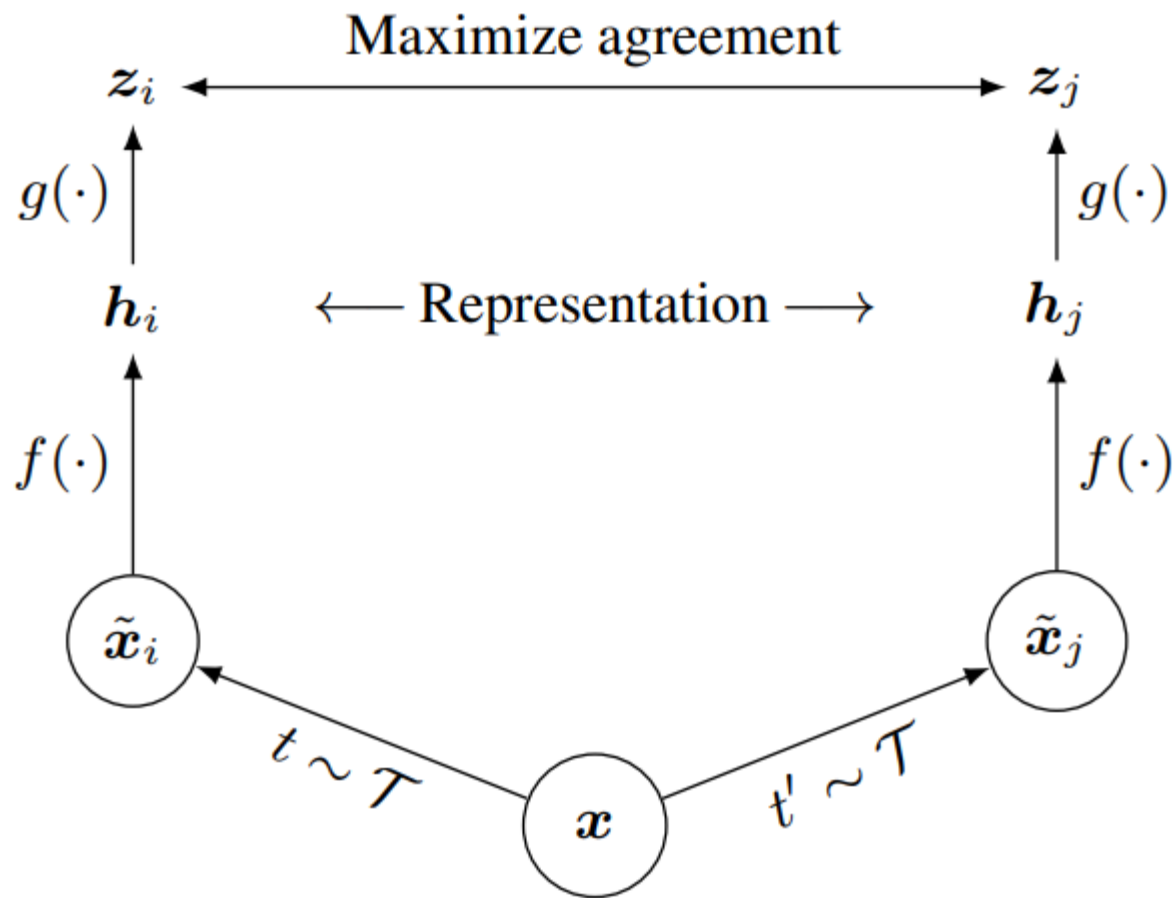
# Main Components:

- Composition of data augmentations plays a critical role in defining effective predictive tasks

- Introducing a learnable nonlinear transformation between the representation and the constrastive loss substantially improves the quality of the learned representations.

- Contrastive learning benefits from larger batch sizes and more training steps compared to the supervised learning.

# The performance of SimCLR (proposed method)

# The Constrastive Learning Framework:



$$\text{sim}(\boldsymbol{u}, \boldsymbol{v}) = \boldsymbol{u}^\top \boldsymbol{v} / \|\boldsymbol{u}\|\|\boldsymbol{v}\|$$

NT-Tent: normalized themperature-scaled cross entropy loss

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

where $\mathbb{1}_{[k\neq i]} \in \{0, 1\}$ is an indicator function evaluating to 1 iff $k \neq i$ and $\tau$ denotes a temperature parameter. The fi-

# The Constrastive Learning Framework:

**Algorithm 1** SimCLR's main learning algorithm.

---

**input:** batch size $N$, constant $\tau$, structure of $f, g, \mathcal{T}$.
**for** sampled minibatch $\{x_k\}_{k=1}^N$ **do**
  **for all** $k \in \{1, \ldots, N\}$ **do**
    draw two augmentation functions $t \sim \mathcal{T}, t' \sim \mathcal{T}$
    *# the first augmentation*
    $\tilde{x}_{2k-1} = t(x_k)$
    $h_{2k-1} = f(\tilde{x}_{2k-1})$              *# representation*
    $z_{2k-1} = g(h_{2k-1})$            *# projection*
    *# the second augmentation*
    $\tilde{x}_{2k} = t'(x_k)$
    $h_{2k} = f(\tilde{x}_{2k})$               *# representation*
    $z_{2k} = g(h_{2k})$               *# projection*
  **end for**
  **for all** $i \in \{1, \ldots, 2N\}$ and $j \in \{1, \ldots, 2N\}$ **do**
    $s_{i,j} = z_i^\top z_j / (\|z_i\| \|z_j\|)$     *# pairwise similarity*
  **end for**
  **define** $\ell(i,j)$ **as** $\ell(i,j) = -\log \dfrac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$
  $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$
  update networks $f$ and $g$ to minimize $\mathcal{L}$
**end for**
**return** encoder network $f(\cdot)$, and throw away $g(\cdot)$

---

# Evaluation Protocol

- Dataset: ImageNet                                        , CIFAR-10 -> Appendix

- Linear evaluation protocol:

  - A linear classifier is trained on top of the frozen base network, and test accuracy is used as a proxy for representation quality.

- Base network: ResNet-50

- Projection head: 2-layer MLP -> 128-dimensional
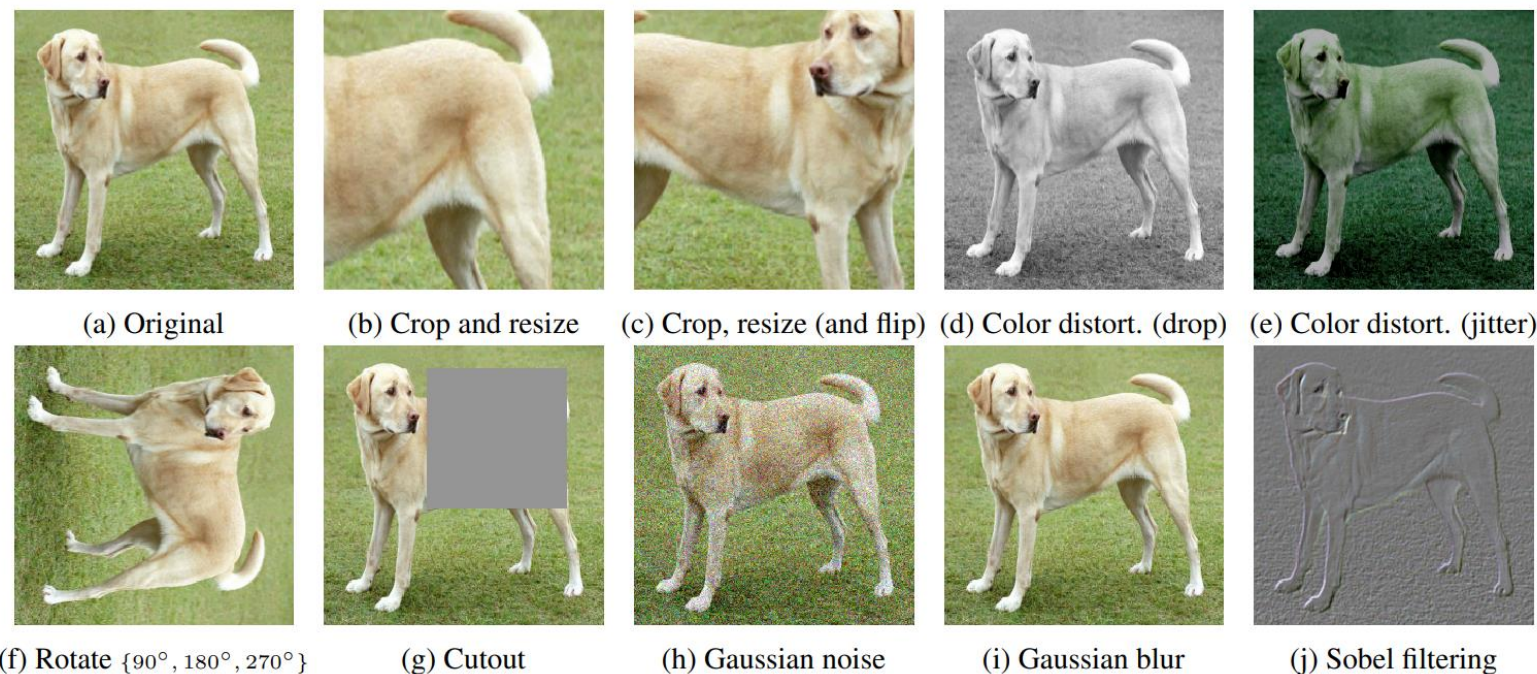
# Data Augmentations



Figure 4. Illustrations of the studied data augmentation operators. Each augmentation can transform data stochastically with some internal parameters (e.g. rotation degree, noise level). Note that we *only* test these operators in ablation, the *augmentation policy used to train our models* only includes *random crop (with flip and resize), color distortion,* and *Gaussian blur.* (Original image cc-by: Von.grzanka)

- Spatial/geometric:
  - Corp and resize
  - Rotation
  - Cutout

- Appearance:
  - Color distortion
  - Gaussian blur
  - Sobel filtering
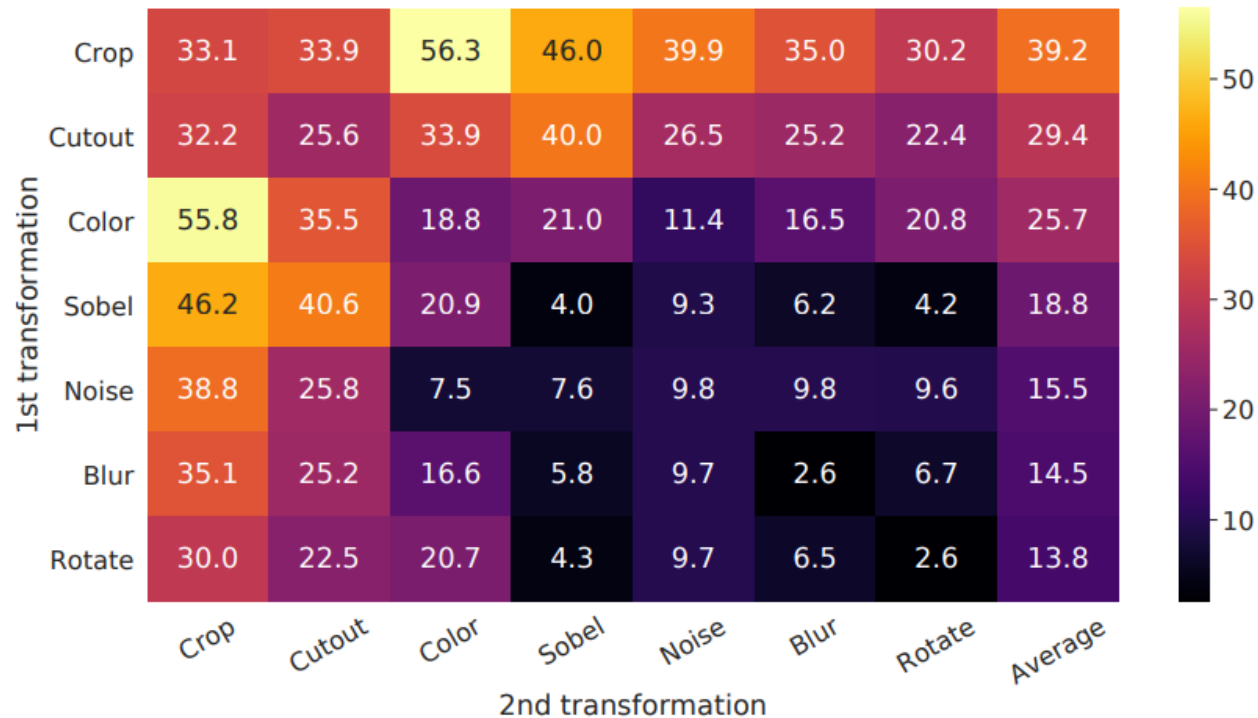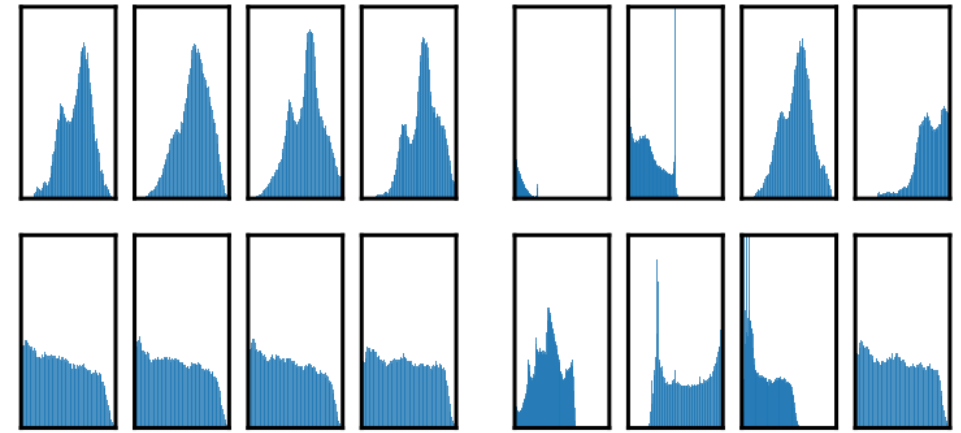
# Data Augmentations (1)



*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.

- No single transformation suffices to learn good representations, even though the model can almost perfectly identify the positive pairs in the contrastive task.

- When composing augmentations, the contrastive task becomes harder, but the quality of representation improves dramatically.

10

# Data Augmentations (2)



Heatmap — 1st transformation (rows) vs 2nd transformation (columns):

| 1st \ 2nd | Crop | Cutout | Color | Sobel | Noise | Blur | Rotate | Average |
|---|---|---|---|---|---|---|---|---|
| Crop | 33.1 | 33.9 | 56.3 | 46.0 | 39.9 | 35.0 | 30.2 | 39.2 |
| Cutout | 32.2 | 25.6 | 33.9 | 40.0 | 26.5 | 25.2 | 22.4 | 29.4 |
| Color | 55.8 | 35.5 | 18.8 | 21.0 | 11.4 | 16.5 | 20.8 | 25.7 |
| Sobel | 46.2 | 40.6 | 20.9 | 4.0 | 9.3 | 6.2 | 4.2 | 18.8 |
| Noise | 38.8 | 25.8 | 7.5 | 7.6 | 9.8 | 9.8 | 9.6 | 15.5 |
| Blur | 35.1 | 25.2 | 16.6 | 5.8 | 9.7 | 2.6 | 6.7 | 14.5 |
| Rotate | 30.0 | 22.5 | 20.7 | 4.3 | 9.7 | 6.5 | 2.6 | 13.8 |



(a) Without color distortion.   (b) With color distortion.

*Figure 6.* Histograms of pixel intensities (over all channels) for different crops of two different images (i.e. two rows). The image for the first row is from Figure 4. All axes have the same range.

- Outstanding augmentations: random cropping, color distortion.

- Authors observe that when using only random cropping as data augmentation is that most patches from an image share a similar color distribution.

- It means that color histograms alone suffices to distinguish images; therefore, it is critical to compose cropping with color distortion in order to learn generalizable features.

11

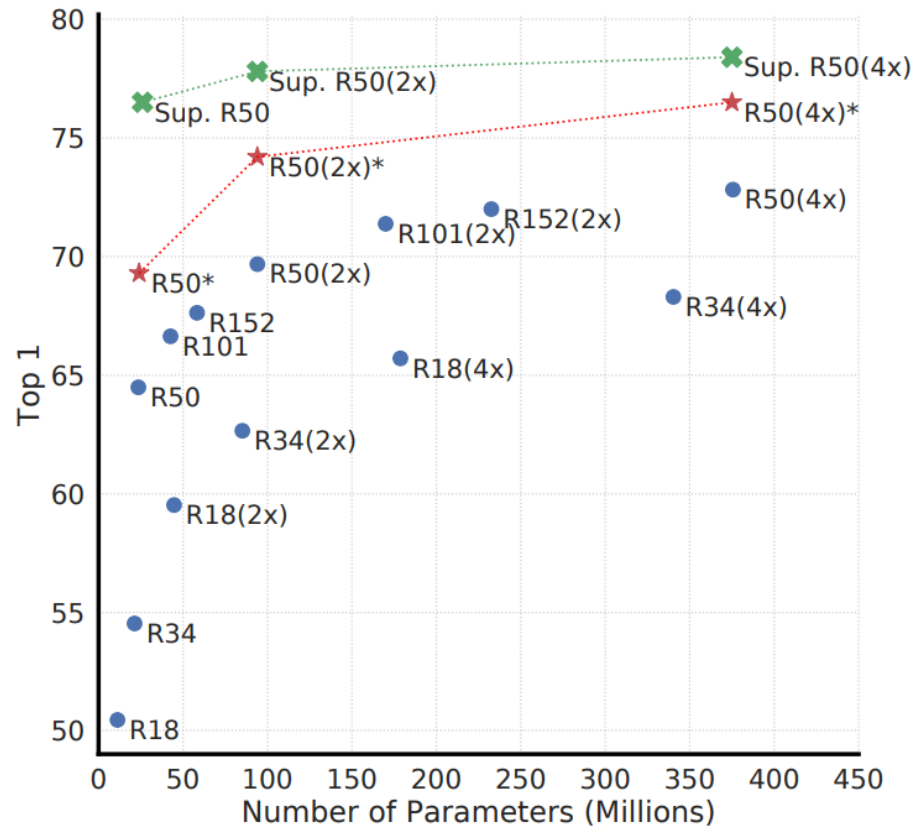# Contrastive learning benefits from bigger models.



Figure 7. Linear evaluation of models with varied depth and width. Models in blue dots are ours trained for 100 epochs, models in red stars are ours trained for 1000 epochs, and models in green crosses are supervised ResNets trained for 90 epochs[7] (He et al., 2016).

- The gap between supervised models and linear classifiers trained on unsupervised models shrinks as the model size increases.
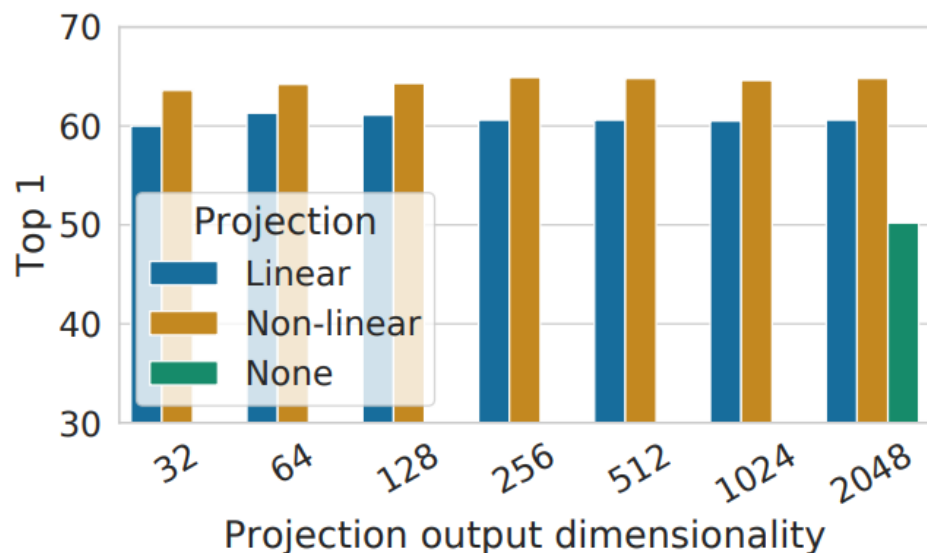
# A nonlinear projection head improve the quality.



Figure 8. Linear evaluation of representations with different projection heads $g(\cdot)$ and various dimensions of $z = g(h)$. The representation $h$ (before projection) is 2048-dimensional here.

| What to predict? | Random guess | Representation | |
|---|---|---|---|
| | | $h$ | $g(h)$ |
| Color vs grayscale | 80 | 99.3 | 97.4 |
| Rotation | 25 | 67.6 | 25.6 |
| Orig. vs corrupted | 50 | 99.5 | 59.6 |
| Orig. vs Sobel filtered | 50 | 96.6 | 56.3 |

Table 3. Accuracy of training additional MLPs on different representations to predict the transformation applied. Other than crop and color augmentation, we additionally and independently add rotation (one of $\{0°, 90°, 180°, 270°\}$), Gaussian noise, and Sobel filtering transformation during the pretraining for the last three rows. Both $h$ and $g(h)$ are of the same dimensionality, i.e. 2048.

- The importance of using the representation before the nonlinear projection is due to loss of information induced by the contrastive loss.

- g(.) can remove information that may be useful for the downstream task.

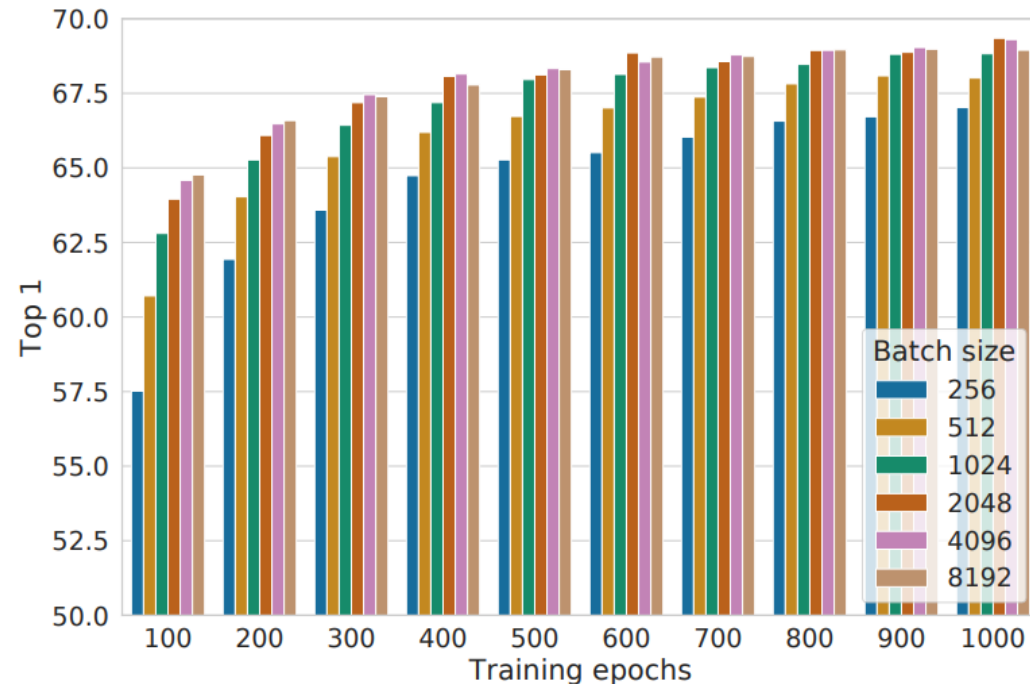# A nonlinear projection head improve the quality.



*Figure 9.* Linear evaluation models (ResNet-50) trained with different batch size and epochs. Each bar is a single run from scratch.

- In contrast to supervised learning, larger batch sizes provide more negative samples, facilitating convergence.

- Training longer also provides more negative examples, improving the results.

# Comparison with State-of-the-art.

Unsupervised setting:

| Method | Architecture | Param. | Top 1 | Top 5 |
|---|---|---|---|---|
| *Methods using ResNet-50:* | | | | |
| Local Agg. | ResNet-50 | 24 | 60.2 | - |
| MoCo | ResNet-50 | 24 | 60.6 | - |
| PIRL | ResNet-50 | 24 | 63.6 | - |
| CPC v2 | ResNet-50 | 24 | 63.8 | 85.3 |
| SimCLR (ours) | ResNet-50 | 24 | **69.3** | **89.0** |
| *Methods using other architectures:* | | | | |
| Rotation | RevNet-50 ($4\times$) | 86 | 55.4 | - |
| BigBiGAN | RevNet-50 ($4\times$) | 86 | 61.3 | 81.9 |
| AMDIM | Custom-ResNet | 626 | 68.1 | - |
| CMC | ResNet-50 ($2\times$) | 188 | 68.4 | 88.2 |
| MoCo | ResNet-50 ($4\times$) | 375 | 68.6 | - |
| CPC v2 | ResNet-161 ($*$) | 305 | 71.5 | 90.1 |
| SimCLR (ours) | ResNet-50 ($2\times$) | 94 | 74.2 | 92.0 |
| SimCLR (ours) | ResNet-50 ($4\times$) | 375 | **76.5** | **93.2** |

*Table 6.* ImageNet accuracies of linear classifiers trained on representations learned with different self-supervised methods.

# Comparison with State-of-the-art.

Few labels
(semi-supervised) setting:

| Method | Architecture | Label fraction | |
| --- | --- | --- | --- |
| | | 1% | 10% |
| | | Top 5 | |
| Supervised baseline | ResNet-50 | 48.4 | 80.4 |
| *Methods using other label-propagation:* | | | |
| Pseudo-label | ResNet-50 | 51.6 | 82.4 |
| VAT+Entropy Min. | ResNet-50 | 47.0 | 83.4 |
| UDA (w. RandAug) | ResNet-50 | - | 88.5 |
| FixMatch (w. RandAug) | ResNet-50 | - | 89.1 |
| S4L (Rot+VAT+En. M.) | ResNet-50 ($4\times$) | - | 91.2 |
| *Methods using representation learning only:* | | | |
| InstDisc | ResNet-50 | 39.2 | 77.4 |
| BigBiGAN | RevNet-50 ($4\times$) | 55.2 | 78.8 |
| PIRL | ResNet-50 | 57.2 | 83.8 |
| CPC v2 | ResNet-161($*$) | 77.9 | 91.2 |
| SimCLR (ours) | ResNet-50 | 75.5 | 87.8 |
| SimCLR (ours) | ResNet-50 ($2\times$) | 83.0 | 91.2 |
| SimCLR (ours) | ResNet-50 ($4\times$) | **85.8** | **92.6** |

*Table 7.* ImageNet accuracy of models trained with few labels.

# Comparison with State-of-the-art.

Transfer learning
(Pretrained?) setting:

| | Food | CIFAR10 | CIFAR100 | Birdsnap | SUN397 | Cars | Aircraft | VOC2007 | DTD | Pets | Caltech-101 | Flowers |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Linear evaluation:* | | | | | | | | | | | | |
| SimCLR (ours) | **76.9** | **95.3** | 80.2 | 48.4 | **65.9** | 60.0 | 61.2 | **84.2** | 78.9 | 89.2 | **93.9** | **95.0** |
| Supervised | 75.2 | **95.7** | **81.2** | **56.4** | 64.9 | **68.8** | **63.8** | 83.8 | **78.7** | **92.3** | **94.1** | 94.2 |
| *Fine-tuned:* | | | | | | | | | | | | |
| SimCLR (ours) | **89.4** | **98.6** | **89.0** | **78.2** | **68.1** | **92.1** | **87.0** | **86.6** | 77.8 | 92.1 | **94.1** | 97.6 |
| Supervised | 88.7 | 98.3 | **88.7** | **77.8** | 67.0 | 91.4 | **88.0** | 86.5 | **78.8** | **93.2** | **94.2** | **98.0** |
| Random init | 88.3 | 96.0 | 81.9 | **77.0** | 53.7 | 91.3 | 84.8 | 69.4 | 64.1 | 82.7 | 72.5 | 92.5 |

*Table 8.* Comparison of transfer learning performance of our self-supervised approach with supervised baselines across 12 natural image classification datasets, for ResNet-50 ($4\times$) models pretrained on ImageNet. Results not significantly worse than the best ($p > 0.05$, permutation test) are shown in bold. See Appendix B.8 for experimental details and results with standard ResNet-50.