# Cornell Bowers C·IS
## College of Computing and Information Science

# Deep Learning

## Week 05: Self-Supervised Vision Networks

# Overview

- Learning image representations
- Pre-training
- Self-supervised learning
- Contrastive Loss
  - Data Augmentation
  - SimCLR

Cornell Bowers C·IS

$f(x)$ = word embedding

**"orange"**
$x_1$

**"purple"**
$x_4$

**"puppies"**
$x_2$

**"apple"**
$x_5$

**"dogs"**
$x_3$

**Vector Space**

$f(x_1)$

$f(x_5)$ $f(x_4)$

$f(x_3)$

$f(x_2)$
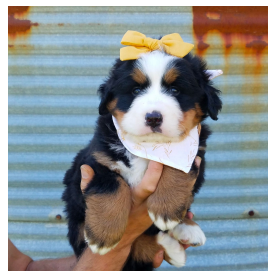
**Semantically different:**
puppy vs. cow

**Structurally similar:**
black and white
animal, grass

$x_1$

$x_2$

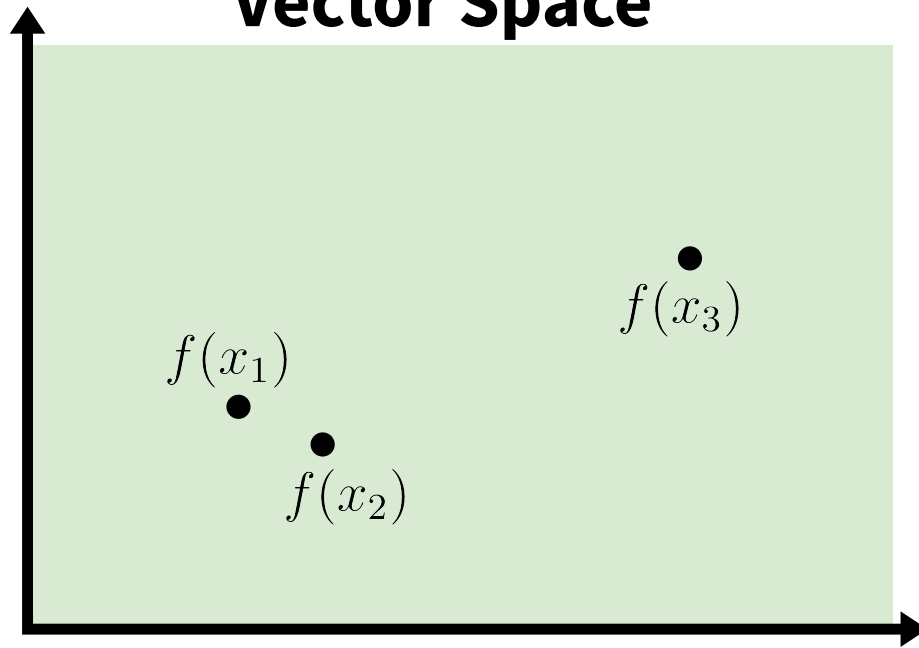**Structurally different:**
hands, different
backgrounds

**Semantically similar:**
Bernese puppies

$x_3$

$f(x)$ = raw pixels

**Vector Space**

$f(x_3)$

$f(x_1)$

$f(x_2)$

# Pixel-Space: Nearest Neighbors

- Dominated by shallow similarities
  - Background, etc.
- Poor semantic alignment



Cifar-10 Example
(http://cs231n.stanford.edu/slides/2023/lecture_13.pdf)

**Semantically different:**
puppy vs. cow

**Structurally similar:**
black and white
animal, grass

$x_1$

$f(x)$ = classification network

**Vector Space**

$x_2$

$f(x_1)$

$f(x_2)$

$f(x_3)$

**Structurally different:**
hands, different
backgrounds

**Semantically similar:**
Bernese puppies

$x_3$

# Image Classification



Image features!

flatten

0.9 "dog"

0.1 "cat"

feature extraction

classification

# Neural Net Features: Nearest Neighbors

- Image classification features work really well!
- Strong semantic alignment
- More robust to shallow variations



ImageNet Classification with Deep Convolutional Neural Networks by Krizhevsky et al.

**Pretraining:** Train a general purpose model on lots of data, and later customize it for more specific tasks



CV: **Imagenet**

NLP: **BERT**



Pre-training

**Vector Space**



Already have a very well-defined vector space

Image Pretraining

**First, train on a large, diverse dataset so that our model learns to extract robust image features**

**Vector Space**



**training**

# Pre-train

CAT

# Fine-tune

- ● Use image classification backbone as a feature extractor for other vision tasks
  - ○ E.g. Detection
  - ○ E.g. Instance segmentation

- ● Significantly accelerates training
  - ○ Random init requires much longer training

bbox AP: R50-FPN, GN



typical fine-tuning schedule

— random init
— w/ pre-train

iterations ($10^5$)

Ren, Shaoqing, et al. "Faster r-cnn: Towards real-time object detection with region proposal networks." Advances in neural information processing systems 28 (2015).
He, Kaiming, Ross Girshick, and Piotr Dollár. "Rethinking imagenet pre-training." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2019.

# Few Shot Learning

**Adapt to variations within known classes, with LIMITED labeled training data**

- We've only seen a few puppies and a few kittens, but a lot of other pretrained data

**Vector Space**

# What are potential problems with supervised pre-training?



Figure 1: Example images from the IMAGENET, the *retinal fundus photographs*, and the CHEXPERT datasets, respectively. The fundus photographs and chest x-rays have much higher resolution than the IMAGENET images, and are classified by looking for small local variations in tissue.

Raghu, Maithra, et al. "Transfusion: Understanding transfer learning for medical imaging." Advances in neural information processing systems 32 (2019).

Cornell Bowers C·IS

Transfer Learning

**Images may be out-of-distribution from the training data**

$x_1$   $x_2$   $x_3$

**Vector Space**

IMAGENET

$f(x_3) \bullet$

$f(x_1) \bullet$

$f(x_2) \bullet$

# Potential Problems?

- Classify diabetic retinopathy in retinal photographs
- Introduce simple architecture
  - Sequence of: Convolution, Batchnorm, ReLU (CBR)



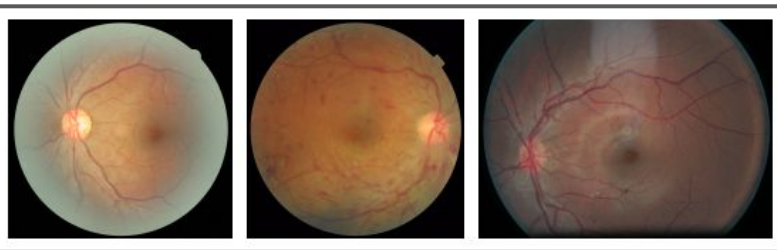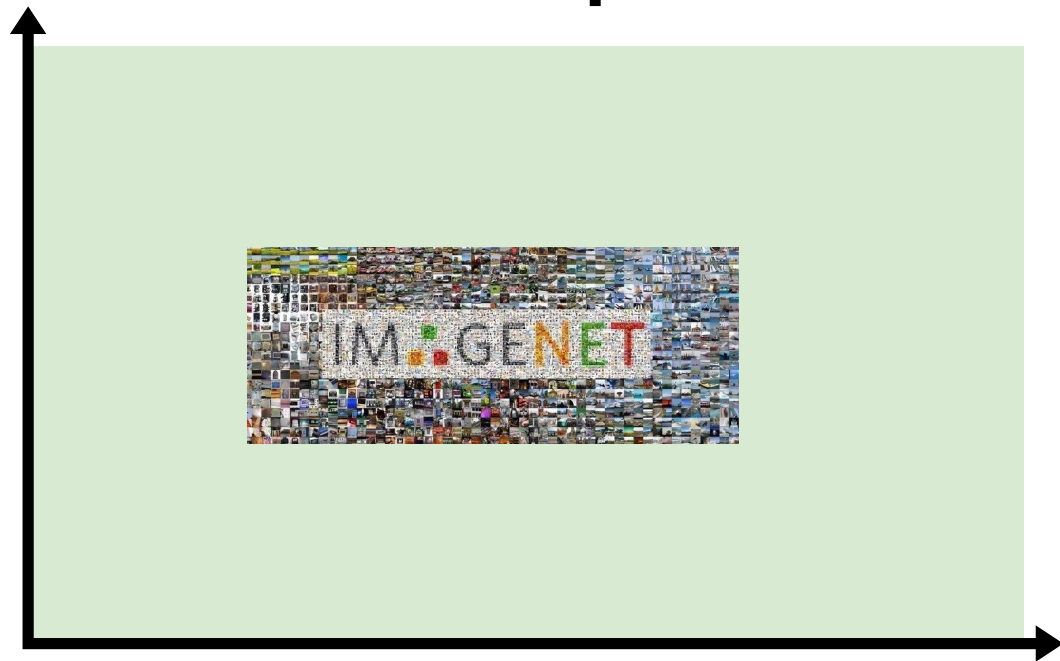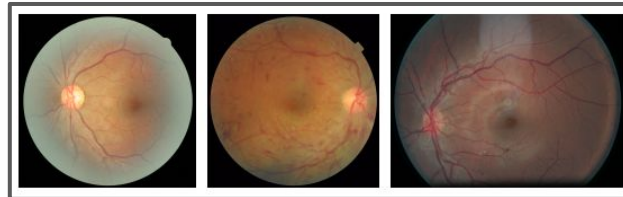| Dataset | Model Architecture | Random Init | Transfer | Parameters | IMAGENET Top5 |
|---------|--------------------|-------------|----------|------------|---------------|
| RETINA | Resnet-50 | 96.4% ± 0.05 | 96.7% ± 0.04 | 23570408 | 92.% ± 0.06 |
| RETINA | Inception-v3 | 96.6% ± 0.13 | 96.7% ± 0.05 | 22881424 | 93.9% |
| RETINA | CBR-LargeT | 96.2% ± 0.04 | 96.2% ± 0.04 | 8532480 | 77.5% ± 0.03 |
| RETINA | CBR-LargeW | 95.8% ± 0.04 | 95.8% ± 0.05 | 8432128 | 75.1% ± 0.3 |
| RETINA | CBR-Small | 95.7% ± 0.04 | 95.8% ± 0.01 | 2108672 | 67.6% ± 0.3 |
| RETINA | CBR-Tiny | 95.8% ± 0.03 | 95.8% ± 0.01 | 1076480 | 73.5% ± 0.05 |

Table 1: **Transfer learning and random initialization perform comparably across both standard IMAGENET architectures and simple, lightweight CNNs for AUCs from diagnosing moderate DR. Both sets of models also have similar AUCs, despite significant differences in size and complexity.** Model performance on DR diagnosis is also not closely correlated with IMAGENET performance, with the small models performing poorly on IMAGENET but very comparably on the medical task.

Raghu, Maithra, et al. "Transfusion: Understanding transfer learning for medical imaging." Advances in neural information processing systems 32 (2019).

# Not all images are labeled

- Particular problem for specialized domains (e.g. medicine)
  - Annotation is expensive!
- Much easier to collect unlabeled data
  - Similar to text!
- Can we still learn good image representations?

# Self-Supervised Learning

- Aim to learn from data without manual label annotation
  - Useful for specialized domains (e.g. medicine) with limited annotated data
- Self-supervised learning methods solve "pretext" tasks that produce good features for downstream tasks.
  - Learn with supervised learning objectives (e.g.,classification, regression)
  - Labels of these pretext tasks are generated automatically

Original text

Thank you ~~for inviting~~ me to your party ~~last~~ week.

Inputs

Thank you <X> me to your party <Y> week.

Targets

<X> for inviting <Y> last <Z>

Raffel, Colin, et al. "Exploring the limits of transfer learning with a unified text-to-text transformer." The Journal of Machine Learning Research 21.1 (2020): 5485-5551.

# Self-Supervised Learning: Rotation Prediction



| 90° rotation | 270° rotation | 180° rotation | 0° rotation | 270° rotation |

Figure 1: Images rotated by random multiples of 90 degrees (e.g., 0, 90, 180, or 270 degrees). The core intuition of our self-supervised feature learning approach is that if someone is not aware of the concepts of the objects depicted in the images, he cannot recognize the rotation that was applied to them.

Gidaris, Spyros, Praveer Singh, and Nikos Komodakis. "Unsupervised Representation Learning by Predicting Image Rotations." ICLR 2018. 2018.

# Rotation Prediction

- Self-supervised learning by rotating the input image
- Predict which rotation is applied
  - 4-way classification



Objectives:

Rotate 0 degrees — $g(X, y=0)$ — Rotated image: $X^0$ — ConvNet model F(.) — Maximize prob. $F^0(X^0)$ — Predict 0 degrees rotation (y=0)

Rotate 90 degrees — $g(X, y=1)$ — Rotated image: $X^1$ — ConvNet model F(.) — Maximize prob. $F^1(X^1)$ — Predict 90 degrees rotation (y=1)

Rotate 180 degrees — $g(X, y=2)$ — Rotated image: $X^2$ — ConvNet model F(.) — Maximize prob. $F^2(X^2)$ — Predict 180 degrees rotation (y=2)

Rotate 270 degrees — $g(X, y=3)$ — Rotated image: $X^3$ — ConvNet model F(.) — Maximize prob. $F^3(X^3)$ — Predict 270 degrees rotation (y=3)

Image $X$

Gidaris, Spyros, Praveer Singh, and Nikos Komodakis. "Unsupervised Representation Learning by Predicting Image Rotations." ICLR 2018. 2018.

# How to evaluate a self-supervised learning method?

- Don't care about the performance of the self-supervised learning task
  - E.g. Image rotation prediction
- Evaluate the learned feature encoder on downstream target tasks

# How to evaluate a self-supervised learning method?



lots of unlabeled data

self-supervised learning

feature extractor (e.g., a convnet)

supervised learning

evaluate on the target task

e.g. classification, detection

90°

conv    fc

**small amount** of labeled data on the target task
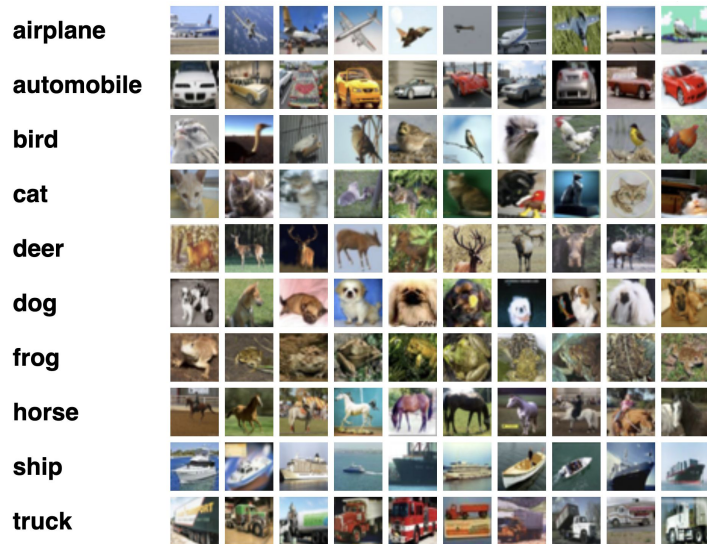
bird

conv    **linear classifier**

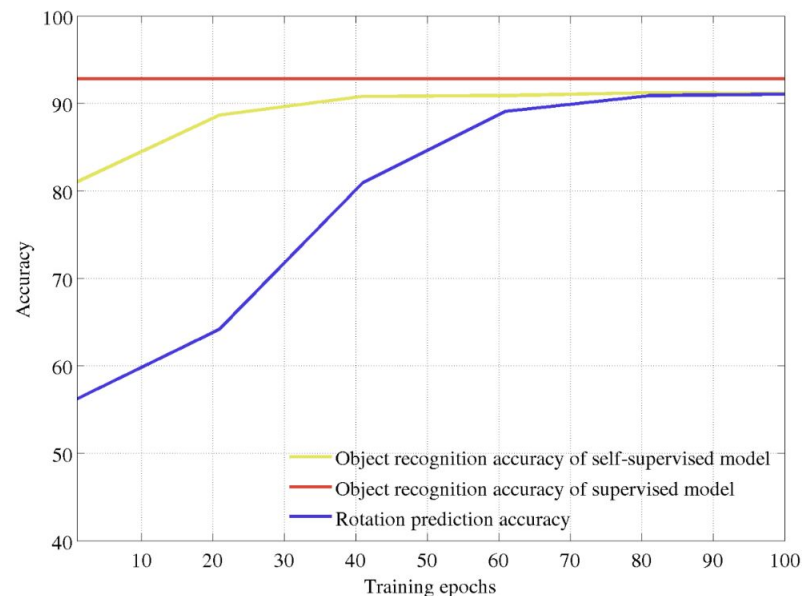1. Learn good feature extractors from self-supervised pretext tasks, e.g., predicting image rotations

2. Attach a shallow network on the feature extractor; train the shallow network on the target task with small amount of labeled data

# Self-Supervised Evaluation

- Downstream performance correlates with prefix task: rotation prediction



Cifar-10 Image Classification



Gidaris, Spyros, Praveer Singh, and Nikos Komodakis. "Unsupervised Representation Learning by Predicting Image Rotations." ICLR 2018. 2018.

# Discuss

We are provided this image without labels: what are some other tasks we can do with it?
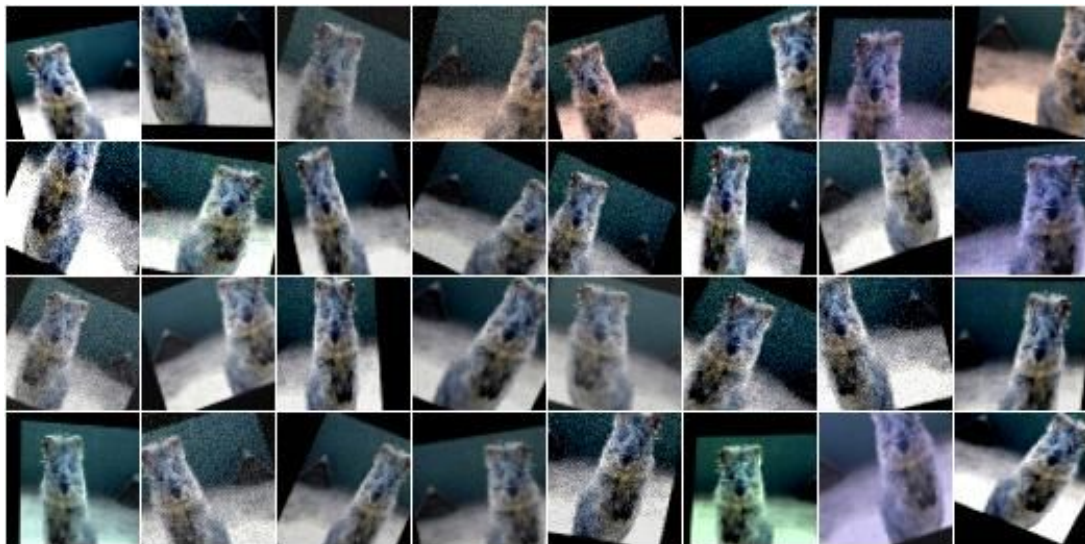
How can we perform self-supervised learning with images?

# Review: Image Augmentation

- Horizontal flips
- Rotate image
- Zoom/crop image
- Brighten/darken image
- Shift colors



Augmentation

https://imgaug.readthedocs.io/en/latest/index.html

Triplet loss function

Margin

$$\ell = \max(0, \|f(\mathbf{x_i}) - f(\mathbf{x}^+)\|^2 - \|f(\mathbf{x_i}) - f(\mathbf{x}^-)\|^2 + c)$$

Ensures loss is not negative

Model should map positive examples close together

Model should map negative examples far apart

# SimCLR: A Simple Contrastive Learning Framework for Images

- Sample two different augmentations of an image
- Apply a base encoder to each view of the image to extract an image feature
  - e.g. ResNet
- Apply an MLP projection head to generate final representations
  - Throw away projection head after training



Chen, Ting, et al. "A simple framework for contrastive learning of visual representations." International conference on machine learning. PMLR, 2020.

# SimCLR Augmentations



(a) Original   (b) Crop and resize   (c) Crop, resize (and flip)   (d) Color distort. (drop)   (e) Color distort. (jitter)

(f) Rotate $\{90°, 180°, 270°\}$   (g) Cutout   (h) Gaussian noise   (i) Gaussian blur   (j) Sobel filtering

Chen, Ting, et al. "A simple framework for contrastive learning of visual representations." International conference on machine learning. PMLR, 2020.

# SimCLR Loss

- Temperature-scaled cross-entropy loss

Model should map positive
examples close together

$$l = -\log\left(\frac{\exp(\mathrm{sim}(\mathbf{x}, \mathbf{x}^+)/\tau)}{\exp(\mathrm{sim}(\mathbf{x}, \mathbf{x}^+)/\tau) + \exp(\mathrm{sim}(\mathbf{x}, \mathbf{x}^-)/\tau)}\right)$$

Model should map negative
examples far apart

Chen, Ting, et al. "A simple framework for contrastive learning of visual representations." International conference on machine learning. PMLR, 2020.

# Discussion: Comparison of Loss Functions

Consider when our model has learned to push most negative examples to have low similarity with the anchor. What happens to the value of each loss?

How does c affect Triplet Loss? (e.g. c = 0.01 vs. c = 1)
How does т affect contrastive loss? (e.g. т = 0.01 vs. т = 1)

**Triplet**: $$l = \max(0, \text{sim}(\mathbf{x}, \mathbf{x}^-) - \text{sim}(\mathbf{x}, \mathbf{x}^+) + c)$$

**SimCLR**: $$l = -\log\left(\frac{\exp(\text{sim}(\mathbf{x}, \mathbf{x}^+)/\tau)}{\exp(\text{sim}(\mathbf{x}, \mathbf{x}^+)/\tau) + \exp(\text{sim}(\mathbf{x}, \mathbf{x}^-)/\tau)}\right)$$

Chen, Ting, et al. "A simple framework for contrastive learning of visual representations." International conference on machine learning. PMLR, 2020.

# SimCLR Algorithm

- Use other images in the mini-batch as negatives
- L2 normalize representations
  - Use cosine similarity as the distance metric
- Compute temperature-scaled cross-entropy for all positive pairs

---

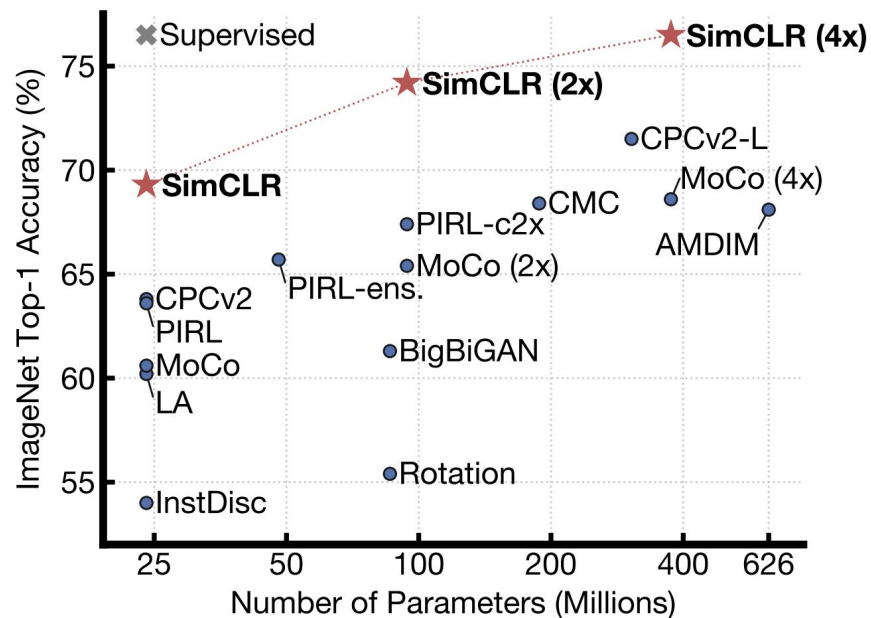**Algorithm 1** SimCLR's main learning algorithm.

---

**input:** batch size $N$, constant $\tau$, structure of $f$, $g$, $\mathcal{T}$.
**for** sampled minibatch $\{\boldsymbol{x}_k\}_{k=1}^N$ **do**
  **for all** $k \in \{1, \ldots, N\}$ **do**
    draw two augmentation functions $t \sim \mathcal{T}$, $t' \sim \mathcal{T}$
    # the first augmentation
    $\tilde{\boldsymbol{x}}_{2k-1} = t(\boldsymbol{x}_k)$
    $\boldsymbol{h}_{2k-1} = f(\tilde{\boldsymbol{x}}_{2k-1})$     # representation
    $\boldsymbol{z}_{2k-1} = g(\boldsymbol{h}_{2k-1})$     # projection
    # the second augmentation
    $\tilde{\boldsymbol{x}}_{2k} = t'(\boldsymbol{x}_k)$
    $\boldsymbol{h}_{2k} = f(\tilde{\boldsymbol{x}}_{2k})$     # representation
    $\boldsymbol{z}_{2k} = g(\boldsymbol{h}_{2k})$     # projection
  **end for**
  **for all** $i \in \{1, \ldots, 2N\}$ and $j \in \{1, \ldots, 2N\}$ **do**
    $s_{i,j} = \boldsymbol{z}_i^\top \boldsymbol{z}_j / (\|\boldsymbol{z}_i\| \|\boldsymbol{z}_j\|)$     # 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 $g(\cdot)$

---

Chen, Ting, et al. "A simple framework for contrastive learning of visual representations." International conference on machine learning. PMLR, 2020.

# SimCLR Results

- Train a linear classifier on features from SimCLR
- Approaches supervised performance!



Chen, Ting, et al. "A simple framework for contrastive learning of visual representations." International conference on machine learning. PMLR, 2020.

# SimCLR Results

- Self-supervised vs. supervised ImageNet pre-training
- Evaluate transfer performance across 12 downstream classification datasets
  - Often outperforms supervised pre-training!

| | 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 |

Chen, Ting, et al. "A simple framework for contrastive learning of visual representations." International conference on machine learning. PMLR, 2020.

# Effect Of Projection Head

- Projects data to "augmentation-invariant" representation
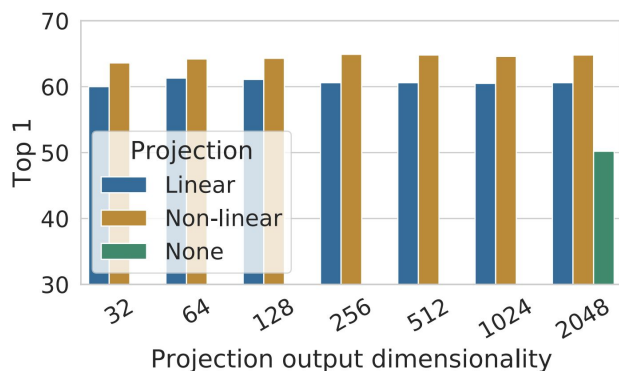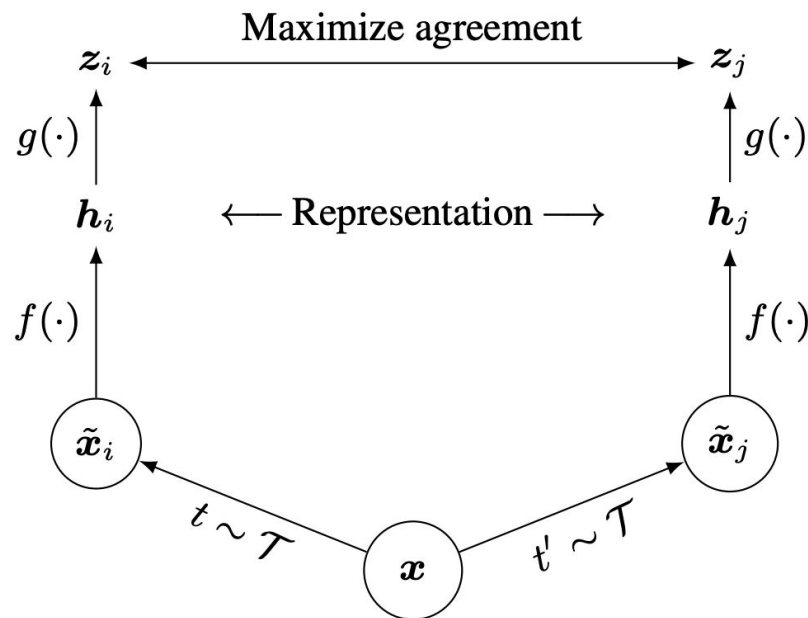  - Less useful features for downstream tasks



*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.

Chen, Ting, et al. "A simple framework for contrastive learning of visual representations." International conference on machine learning. PMLR, 2020.

# Impact of Batch Size
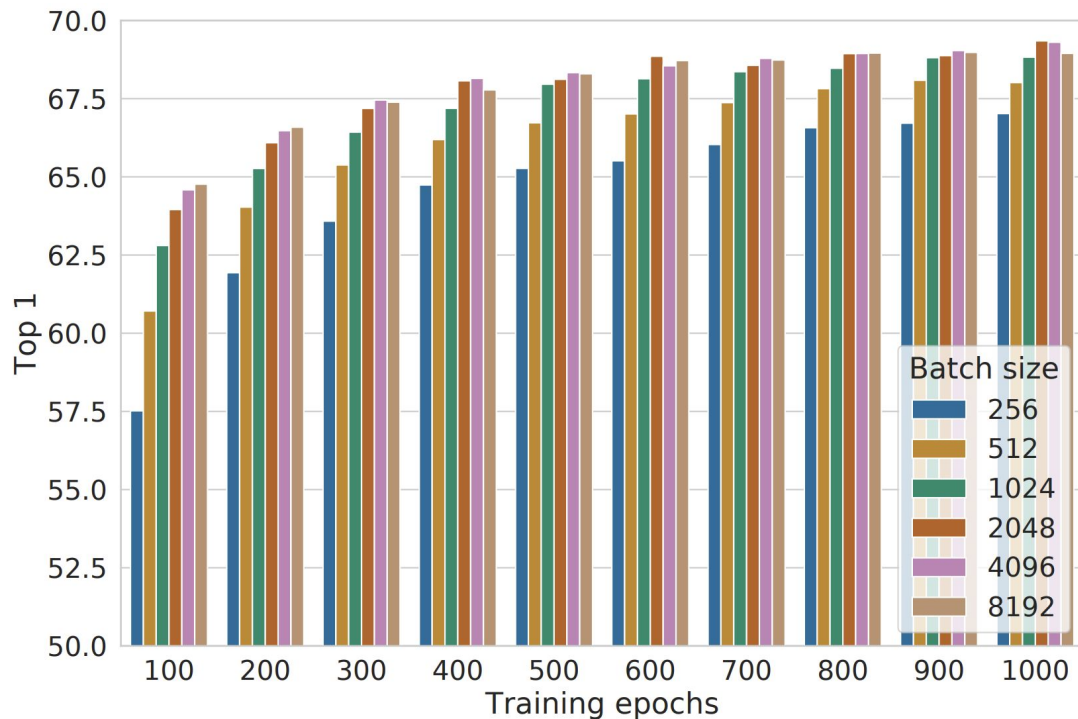
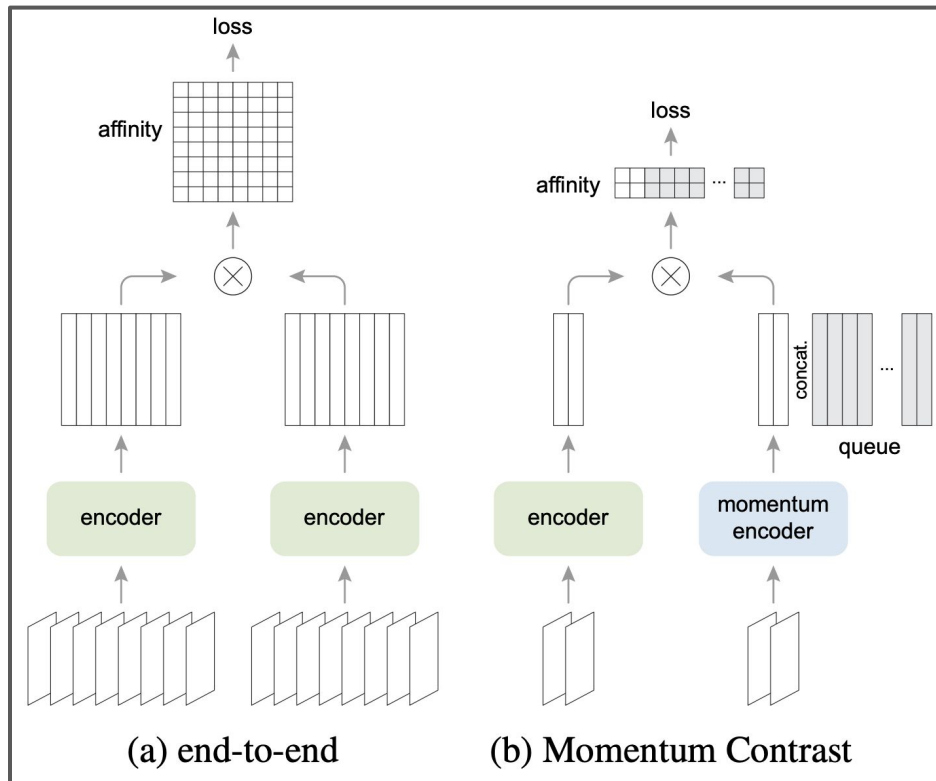- Requires large batches
  - Harder negatives!



*Figure 9.* Linear evaluation models (ResNet-50) trained with different batch size and epochs. Each bar is a single run from scratch.[10]

# Momentum Contrast (MoCo)

- Cache negative samples from earlier batches as you train

- Replace one encoder with an exponential moving average (EMA) of the model
  - Makes queued representations more stable

$$\theta_{\mathrm{k}} \leftarrow m\theta_{\mathrm{k}} + (1-m)\theta_{\mathrm{q}}$$



(a) end-to-end    (b) Momentum Contrast

Chen, Xinlei, et al. "Improved baselines with momentum contrastive learning." arXiv preprint arXiv:2003.04297 (2020).

# Choice of Data Augmentation

# Recap

- Supervised image classification pre-training produces strong image representations
  - Can efficiently transfer to other tasks
- Can apply self-supervised learning to images
  - Prefix tasks: rotation prediction, masked-image modeling, etc.
- Contrastive learning explicitly enforces similarity in representation space
  - Requires defining image augmentations