

Cornell Bowers C-IS

College of Computing and Information Science

Deep Learning

Week 05: Self-Supervised Vision Networks

Logistics

- HW3 is out
- We have a feedback form (due Friday, March 8th)
- Project proposal due Thursday, March 7th

Project

- Aim: to get hands on experience with implementing modern deep learning methods
- To be completed in groups of 2-3
- Find a recent deep learning research paper
- Reproduce a specific result from the paper
 - Need to implement yourself!
- It's ok if open-source implementations exist
 - But you can't use them!

Project Proposal

A page long project proposal due **March 7**. It should contain the following:

- Paper selection:
 - Title, authors, and publication venue of the chosen paper
 - Brief summary of the chosen paper
 - Brief justification of why you choose this paper
- Result Selection
 - Tell us which result you want to replicate
- Re-implementation Plan
 - Describe architecture, method, and metrics
 - Details about how much compute and time is required to replicate results
- **Detailed instructions on canvas**

x_1

The

tastiest

fruits

are

oranges

x_2

The

best

animals

are

puppies

x_3

I

enjoy

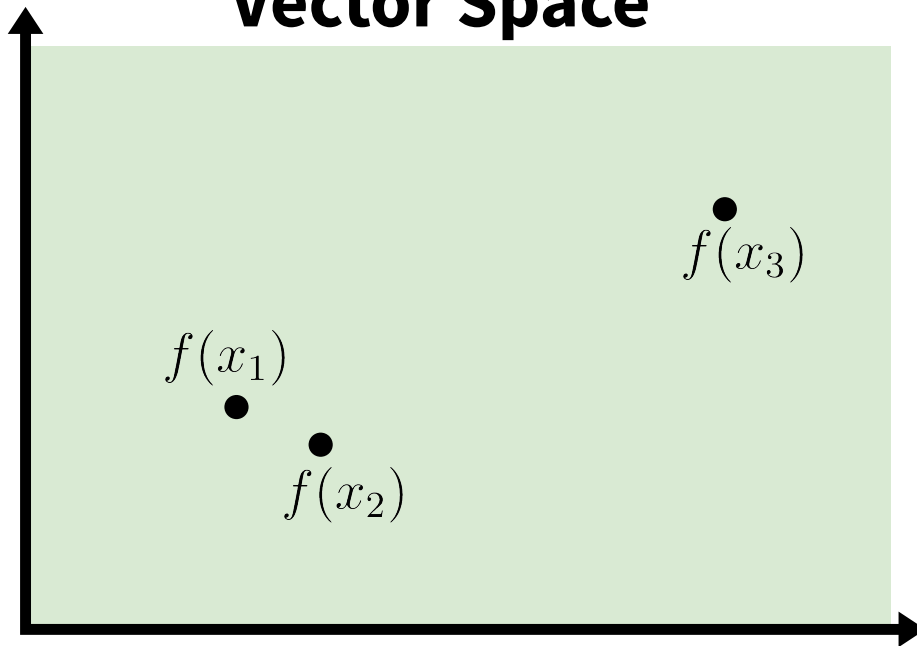
petting

baby

dogs

$f(x) = \text{bag of words}$

Vector Space



“orange”
 x_1

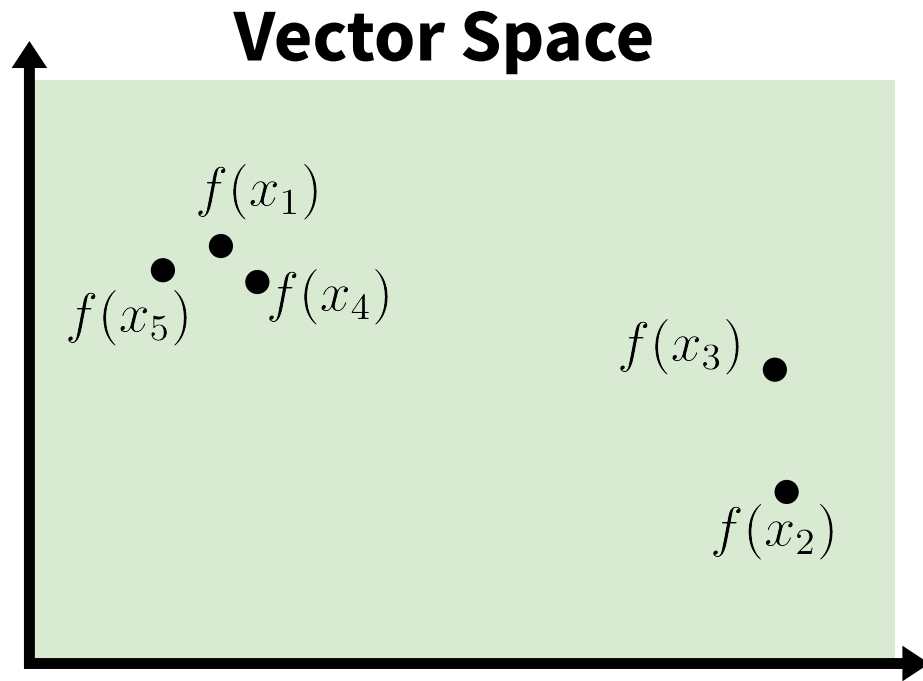
“purple”
 x_4

“puppies”
 x_2

“apple”
 x_5

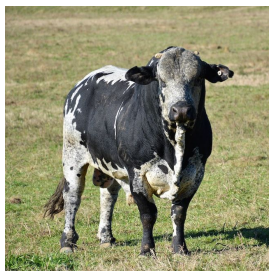
“dogs”
 x_3

$f(x)$ = word embedding



Cornell Bowers C-IS

Semantically different:
puppy vs. cow



Structurally similar:
black and white
animal, grass



Structurally different:
hands, different
backgrounds



Semantically similar:
Bernese puppies

$$f(x) = \text{raw pixels}$$

x_1

Vector Space

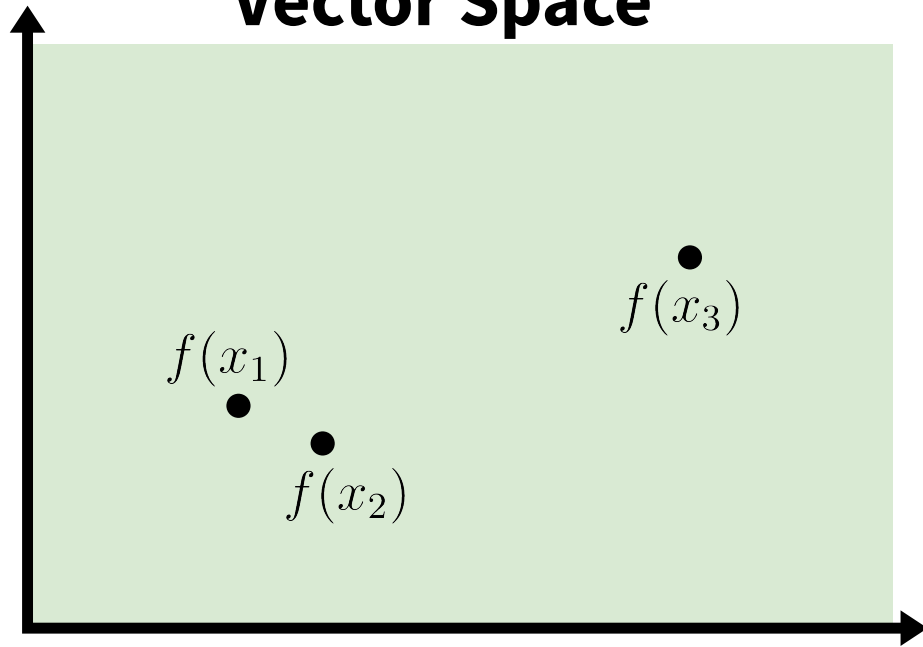
x_2

$f(x_1)$

$f(x_3)$

$f(x_2)$

x_3



Pixel-Space: Nearest Neighbors

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

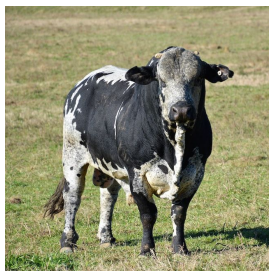


Cifar-10 Example

(http://cs231n.stanford.edu/slides/2023/lecture_13.pdf)

Cornell Bowers C-IS

Semantically different:
puppy vs. cow



Structurally similar:
black and white
animal, grass



Structurally different:
hands, different
backgrounds



Semantically similar:
Bernese puppies

$f(x)$ = classification network

x_1

Vector Space

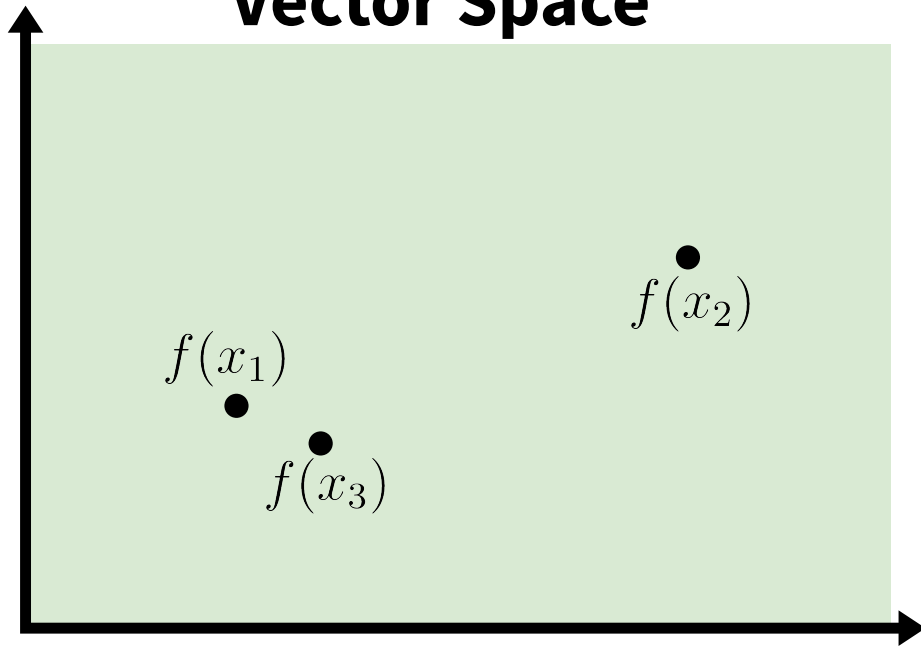
x_2

$f(x_1)$

$f(x_2)$

$f(x_3)$

x_3

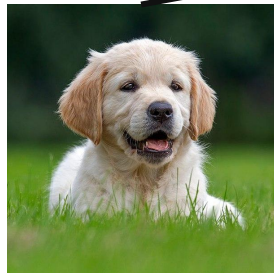


$$f(x) = \text{classification network}$$

How does the network know that these should be mapped to similar space?

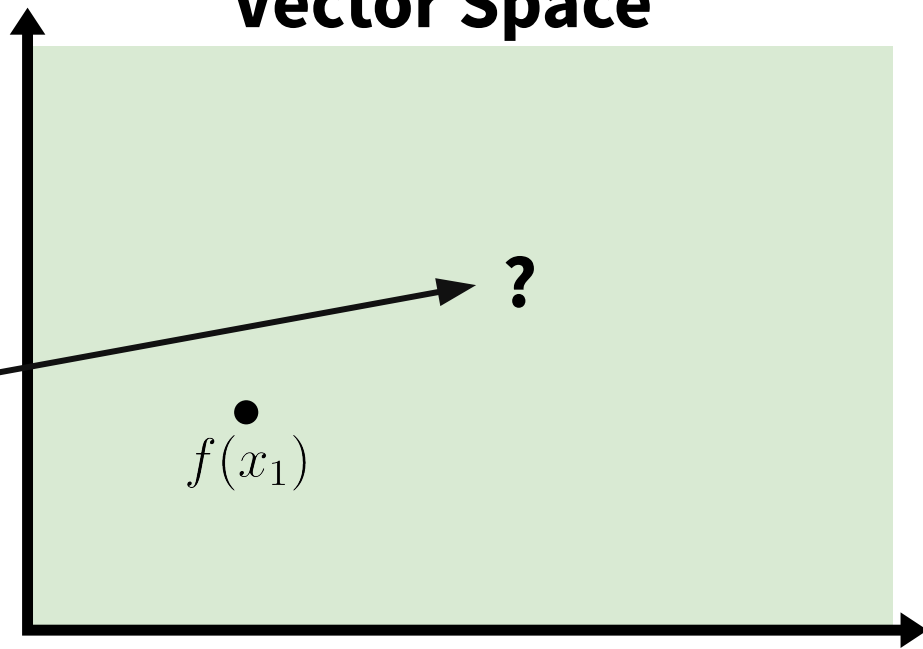


x_1



x_2

Vector Space



$$f(x) = \text{classification network}$$

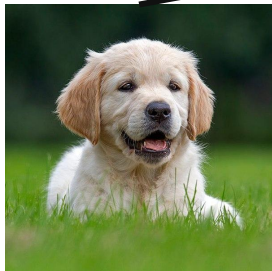
How does the network know that these should be mapped to similar space?

Class
“puppy”



x_1

Class
“puppy”



x_2

Vector Space

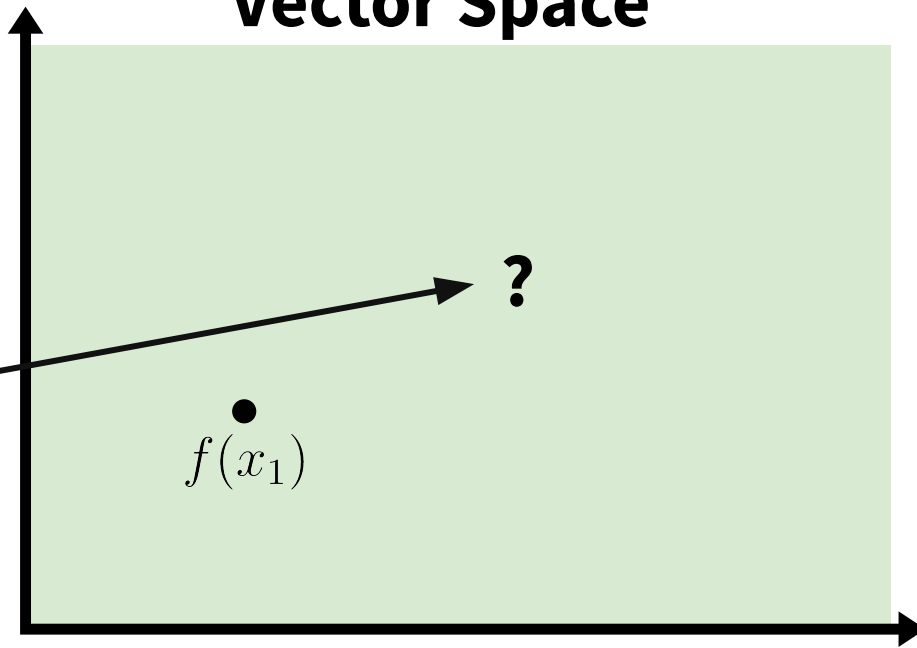
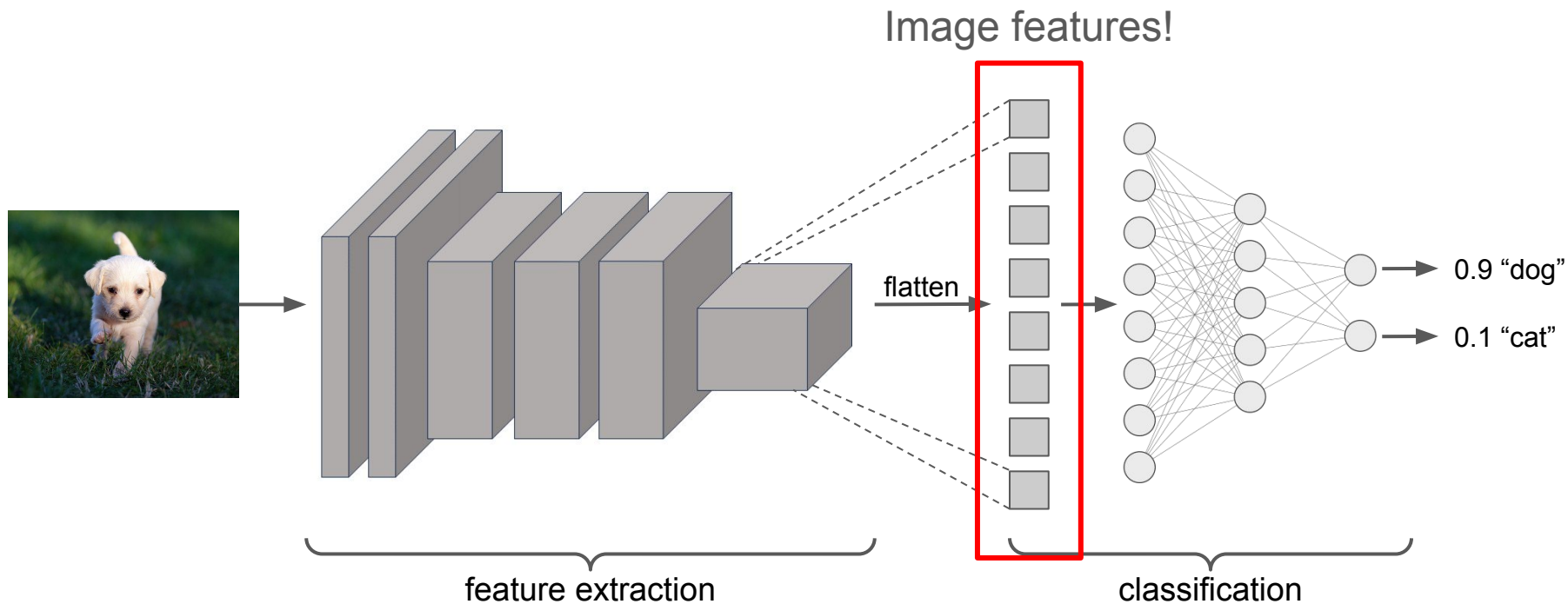


Image Classification



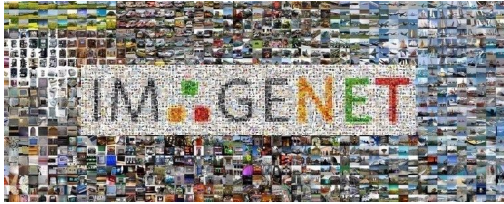
Neural Net Features: Nearest Neighbors

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



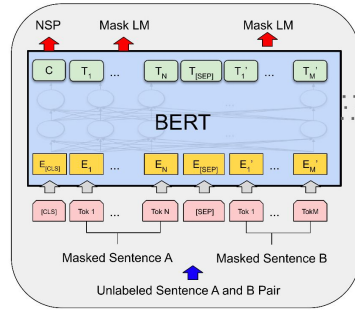
ImageNet Classification with Deep Convolutional Neural Networks by Krizhevsky et al.

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



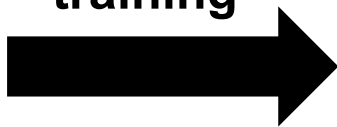
Already have a very
well-defined vector space

Image Pretraining

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



training

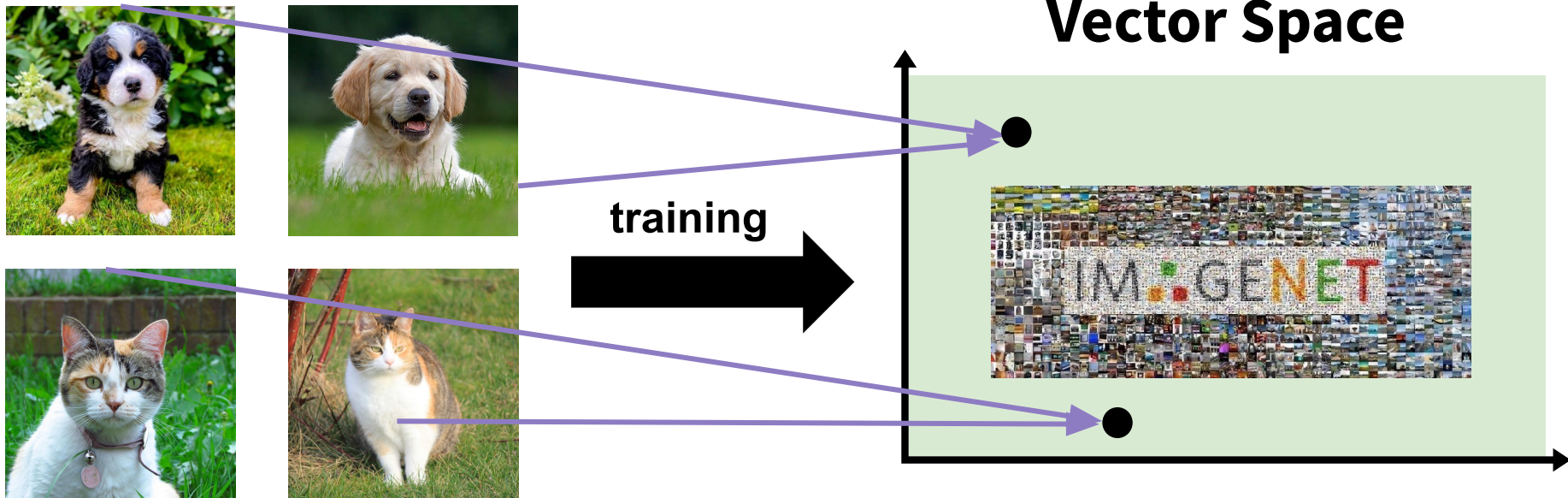


Vector Space



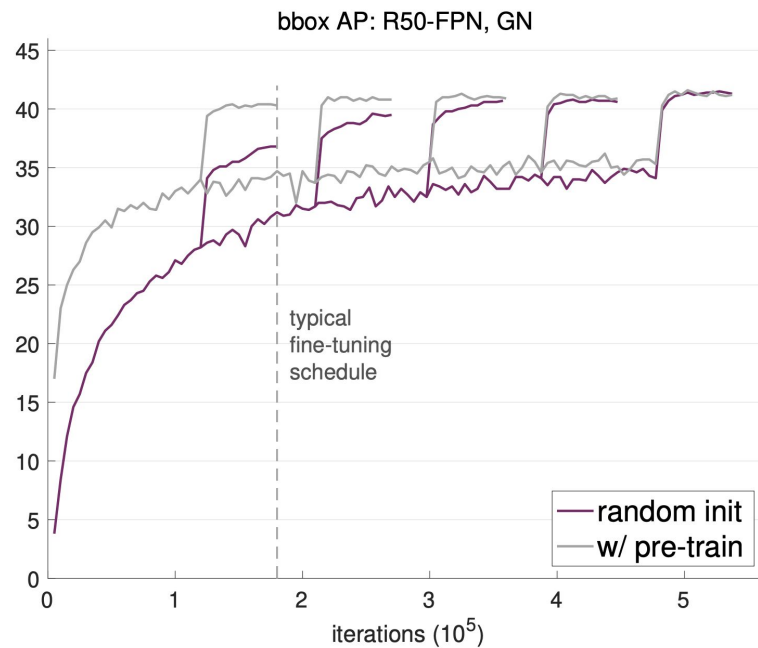
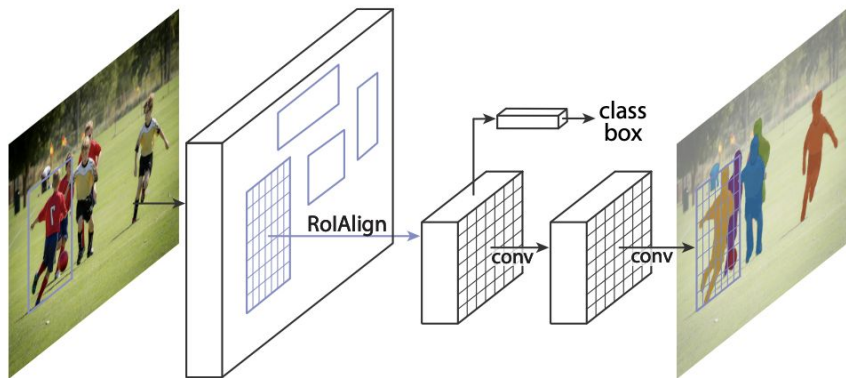
Fine-tuning

Then, finetune for a specific task



Pre-train then Fine-tune

- Use image classification backbone as a feature extractor for other vision tasks
 - E.g. Instance segmentation
- Significantly accelerates training
 - Random init requires much longer training

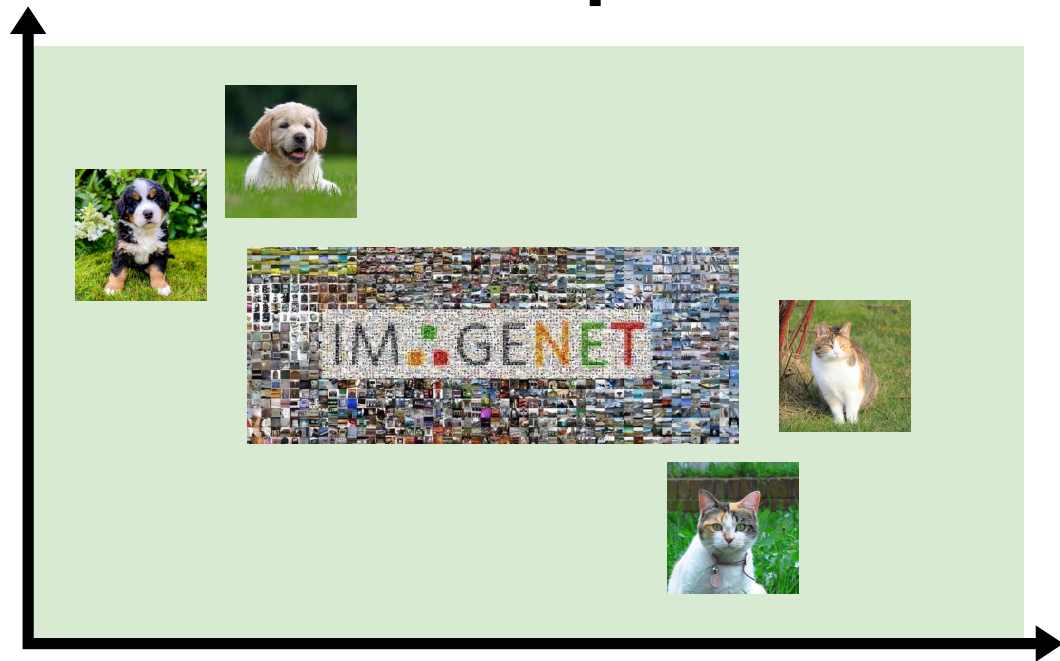


Few Shot Learning

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



Potential Problems?

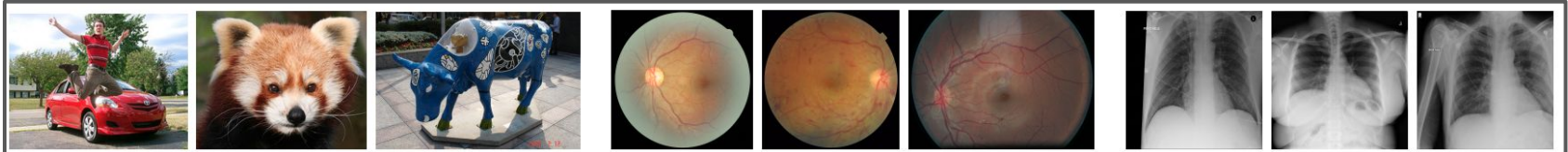


Figure 1: Example images from the IMAGENET, the *retinal fundus photographs*, and the CHEXPART datasets, respectively. The fundus photographs and chest x-rays have much higher resolution than the IMAGENET images, and are classified by looking for small local variations in tissue.

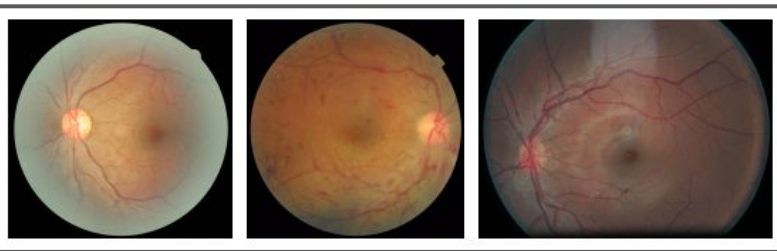
Transfer Learning

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

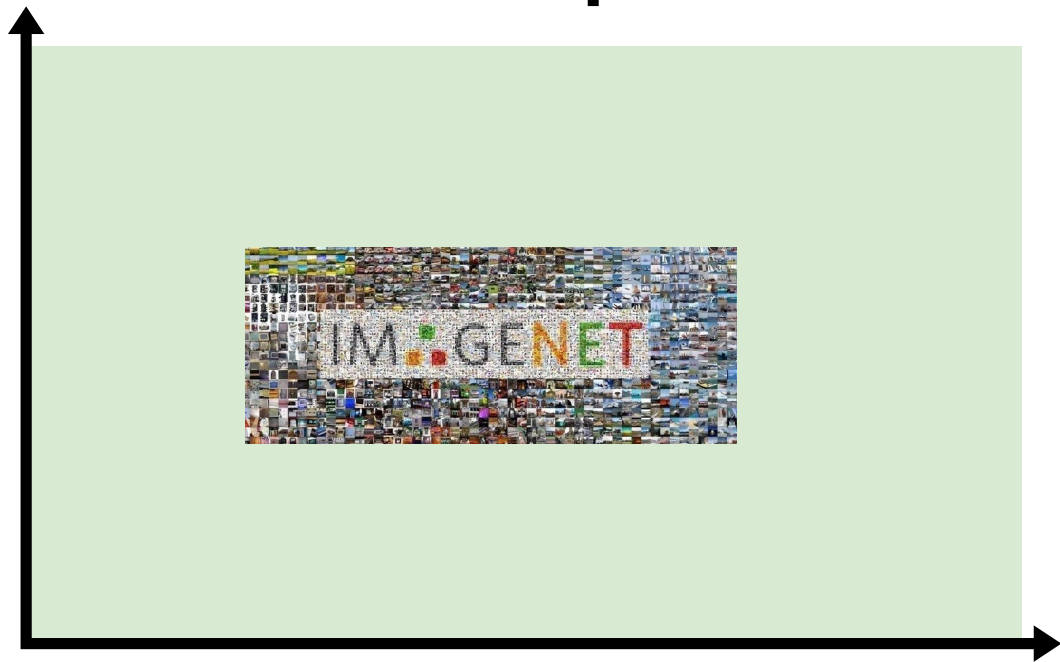
x_1

x_2

x_3



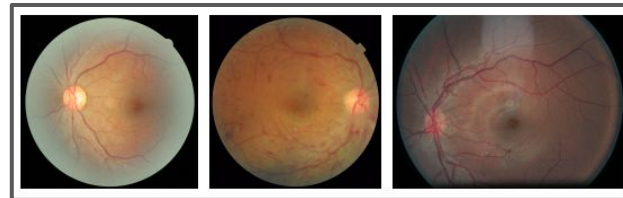
Vector Space



$f(x_3)$ ●
 $f(x_1)$ ● $f(x_2)$ ●

Potential Problems?

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



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

Potential Problems?

- Classify pathologies in chest x-rays

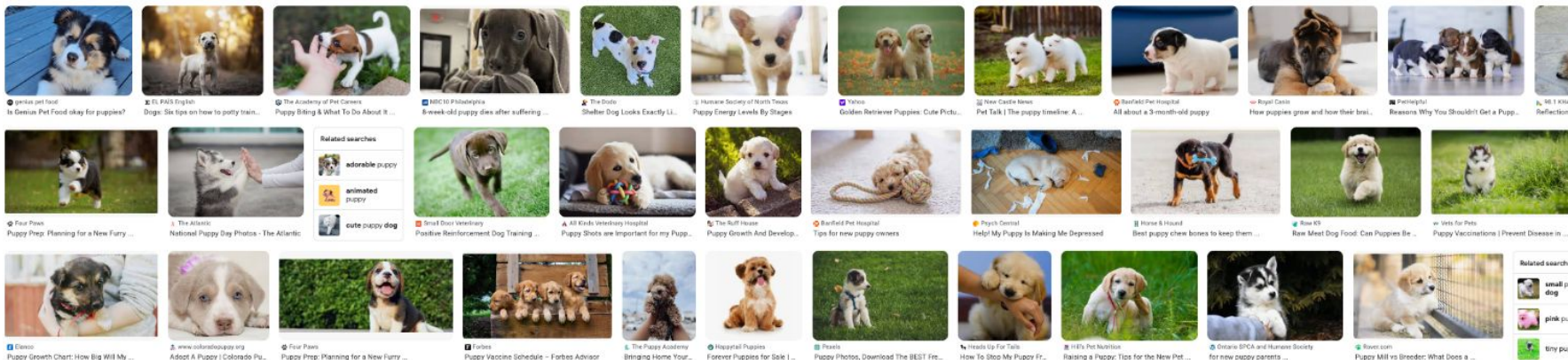


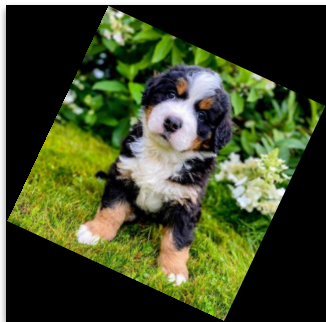
Model Architecture	Atelectasis	Cardiomegaly	Consolidation	Edema	Pleural Effusion
Resnet-50	79.52±0.31	75.23±0.35	85.49±1.32	88.34±1.17	88.70±0.13
Resnet-50 (trans)	79.76±0.47	74.93±1.41	84.42±0.65	88.89±1.66	88.07±1.23
CBR-LargeT	81.52±0.25	74.83±1.66	88.12±0.25	87.97±1.40	88.37±0.01
CBR-LargeT (trans)	80.89±1.68	76.84±0.87	86.15±0.71	89.03±0.74	88.44±0.84
CBR-LargeW	79.79±0.79	74.63±0.69	86.71±1.45	84.80±0.77	86.53±0.54
CBR-LargeW (trans)	80.70±0.31	77.23±0.84	86.87±0.33	89.57±0.34	87.29±0.69
CBR-Small	80.43±0.72	74.36±1.06	88.07±0.60	86.20±1.35	86.14±1.78
CBR-Small (trans)	80.18±0.85	75.24±1.43	86.48±1.13	89.09±1.04	87.88±1.01
CBR-Tiny	80.81±0.55	75.17±0.73	85.31±0.82	84.87±1.13	85.56±0.89
CBR-Tiny (trans)	80.02±1.06	75.74±0.71	84.28±0.82	89.81±1.08	87.69±0.75

Table 2: **Transfer learning provides mixed performance gains on chest x-rays.** Performances (AUC%) of diagnosing different pathologies on the CHEXPART dataset. Again we see that transfer learning does not help significantly, and much smaller models performing comparably.

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?





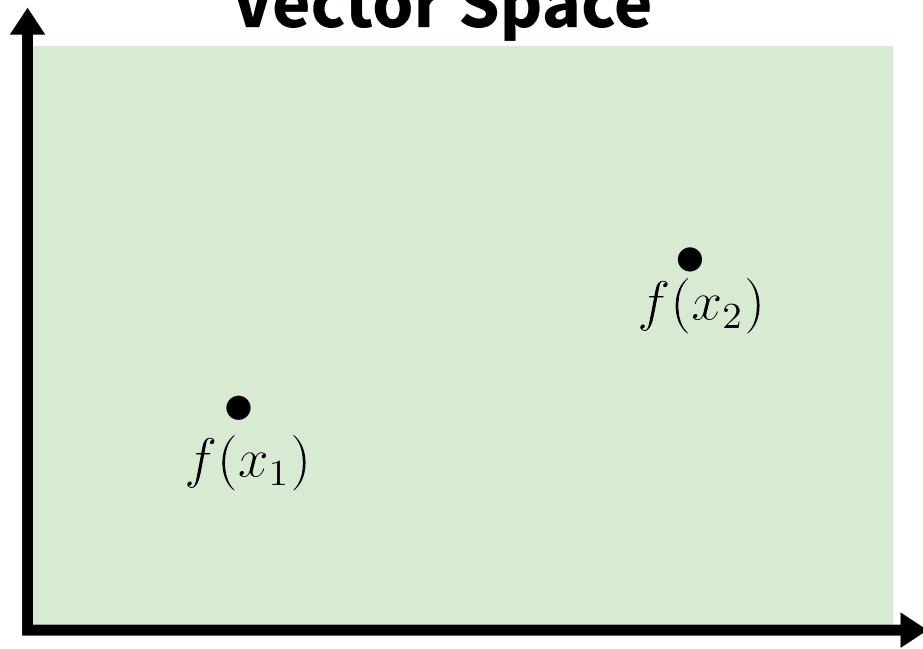
The exact same image, rotated, maps to a completely different location in vector space

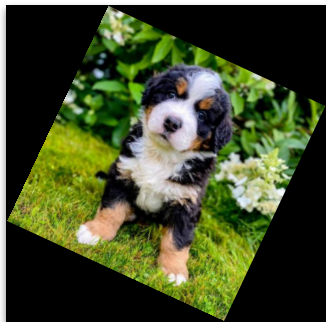
$$f(x) = \text{raw pixel distance}$$

x_1

Vector Space

x_2





How do we learn structure so that these map to similar points in vector space?

x_1

x_2

$$f(x) = ???$$

Vector Space

$f(x_2)$
 $f(x_1)$

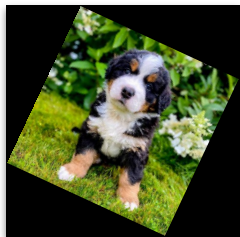
$$f(x) = \text{classification network}$$

Class
“puppy”



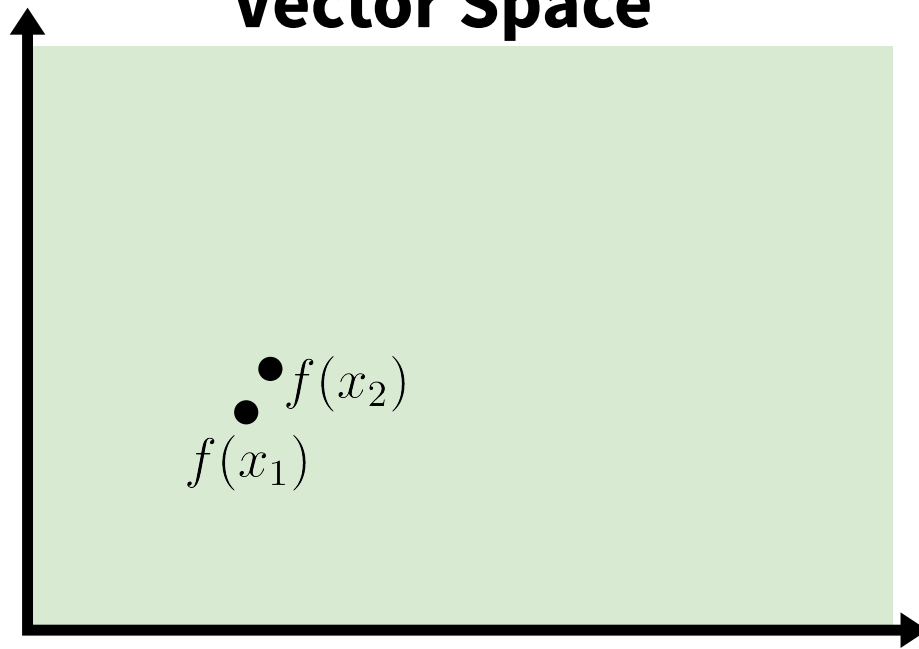
x_1

Class
“puppy”

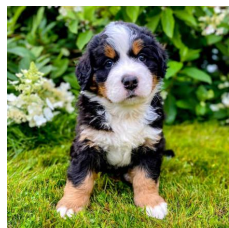


x_2

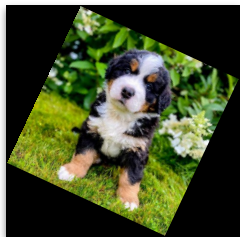
Vector Space



And what if they are unlabeled?



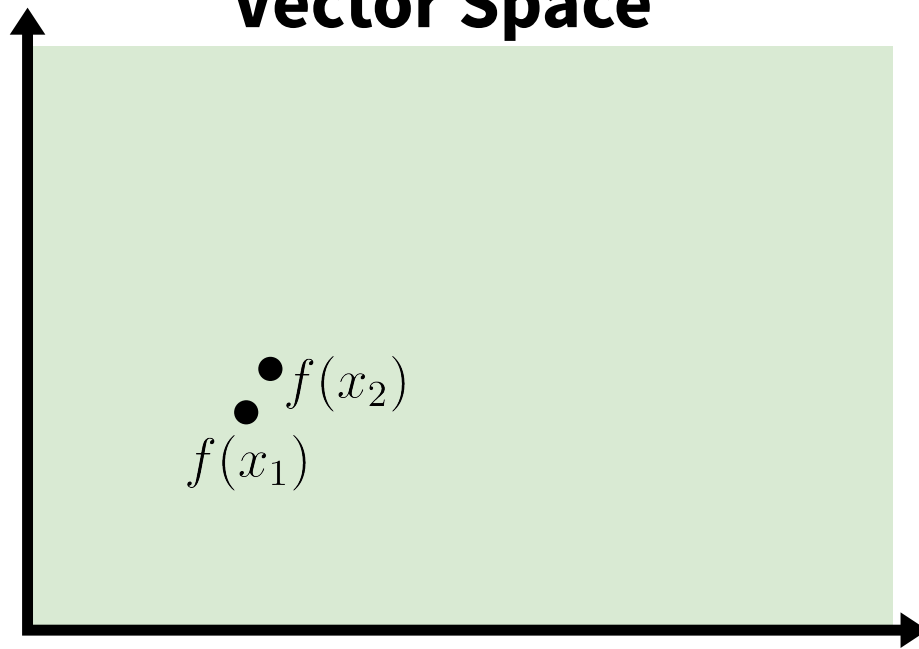
x_1



x_2

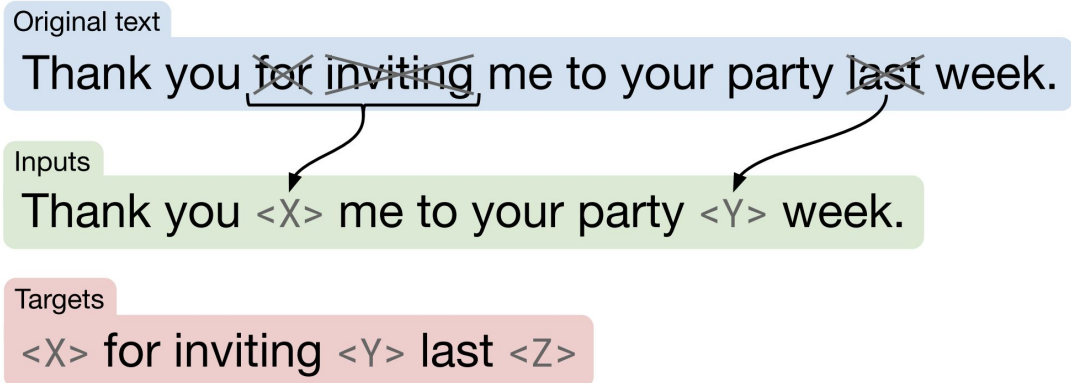
$f(x)$ = classification network

Vector Space

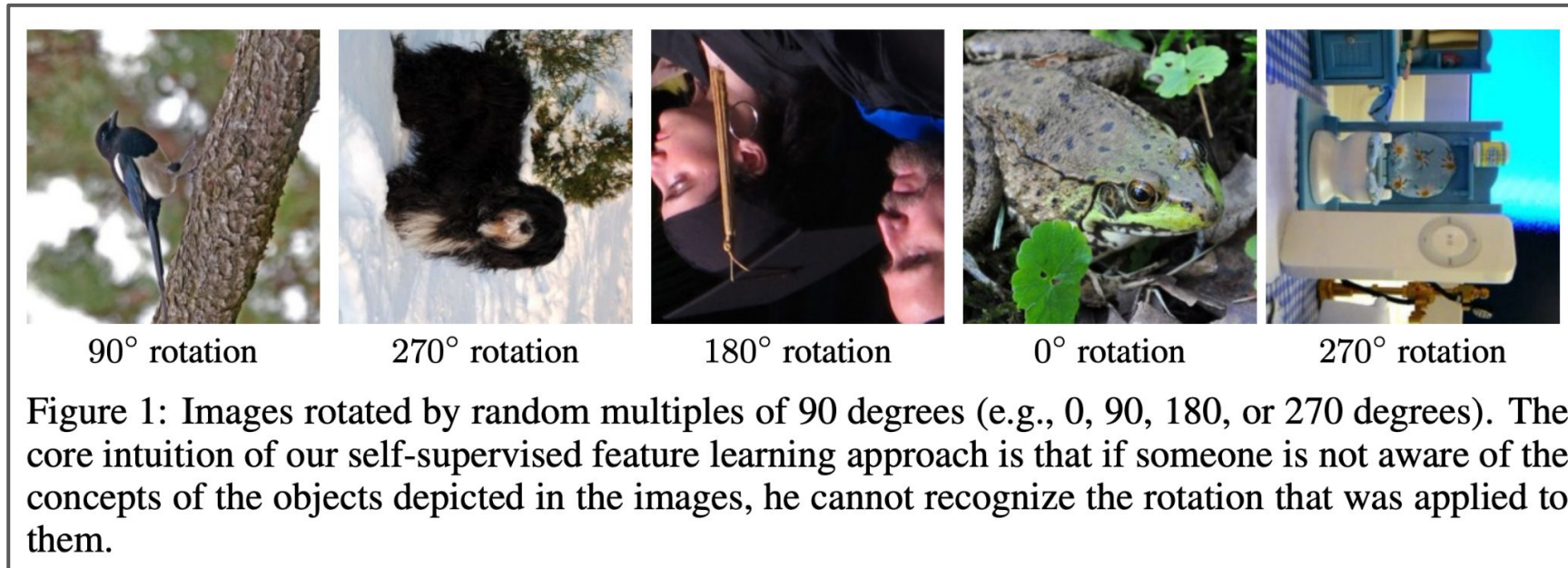


Self-Supervised Learning

- 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

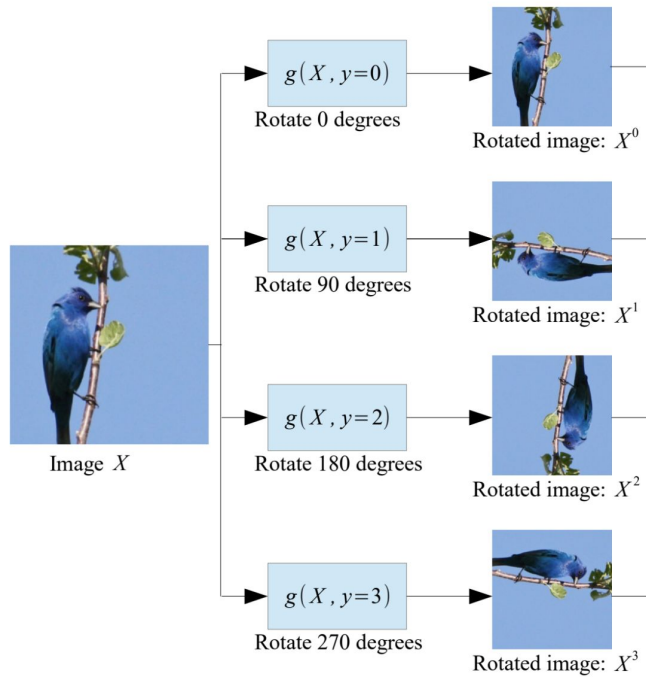


Self-Supervised Learning: Rotation Prediction



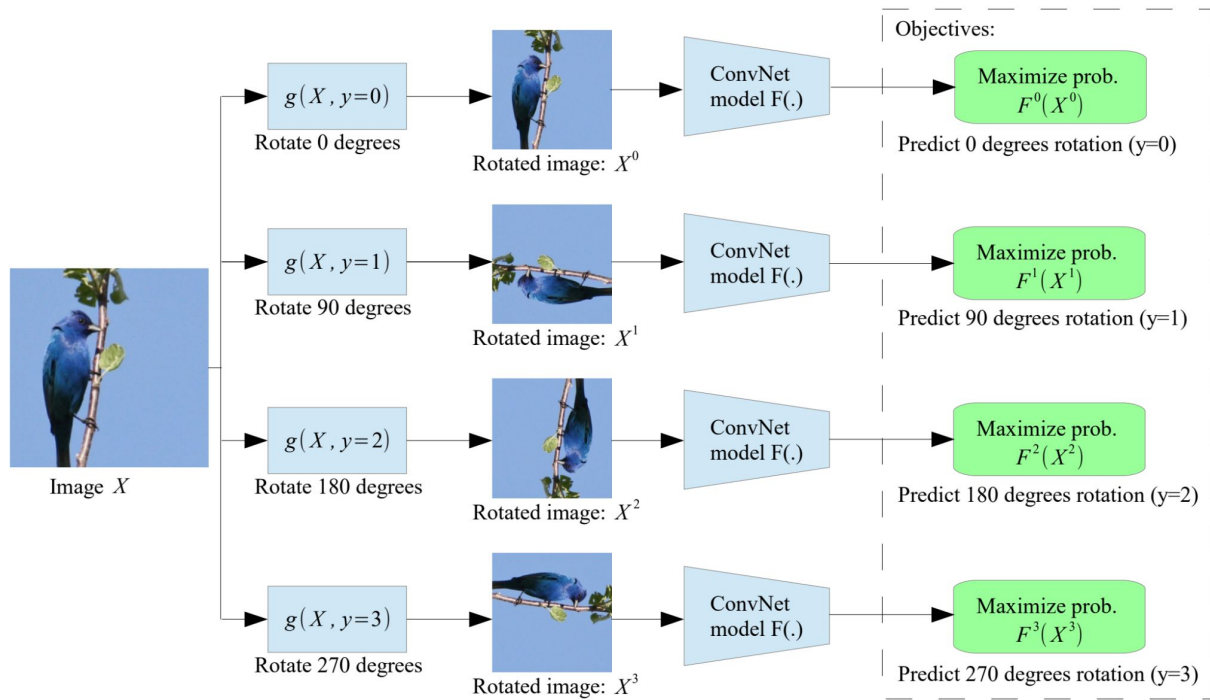
Rotation Prediction

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



Rotation Prediction

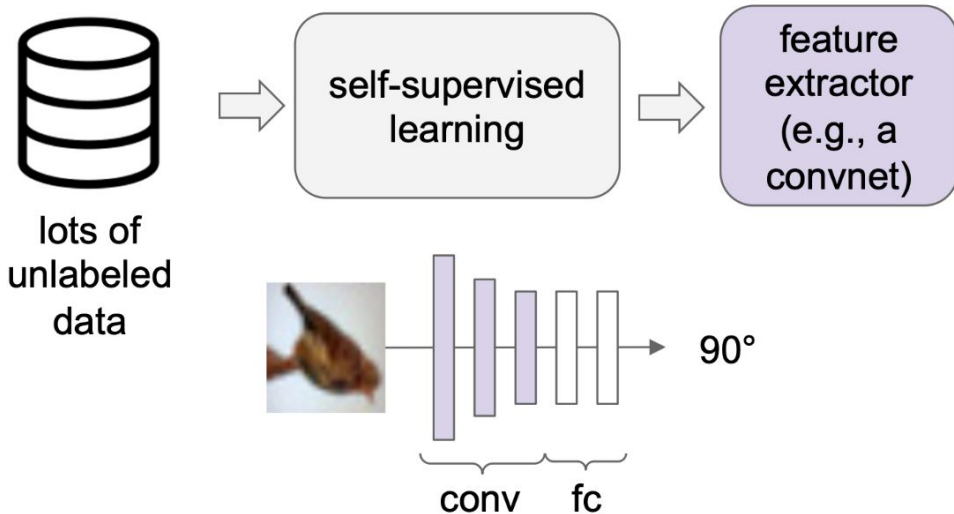
- Self-supervised learning by rotating the input image
- Predict which rotation is applied
 - 4-way classification



How to evaluate a self-supervised learning method?

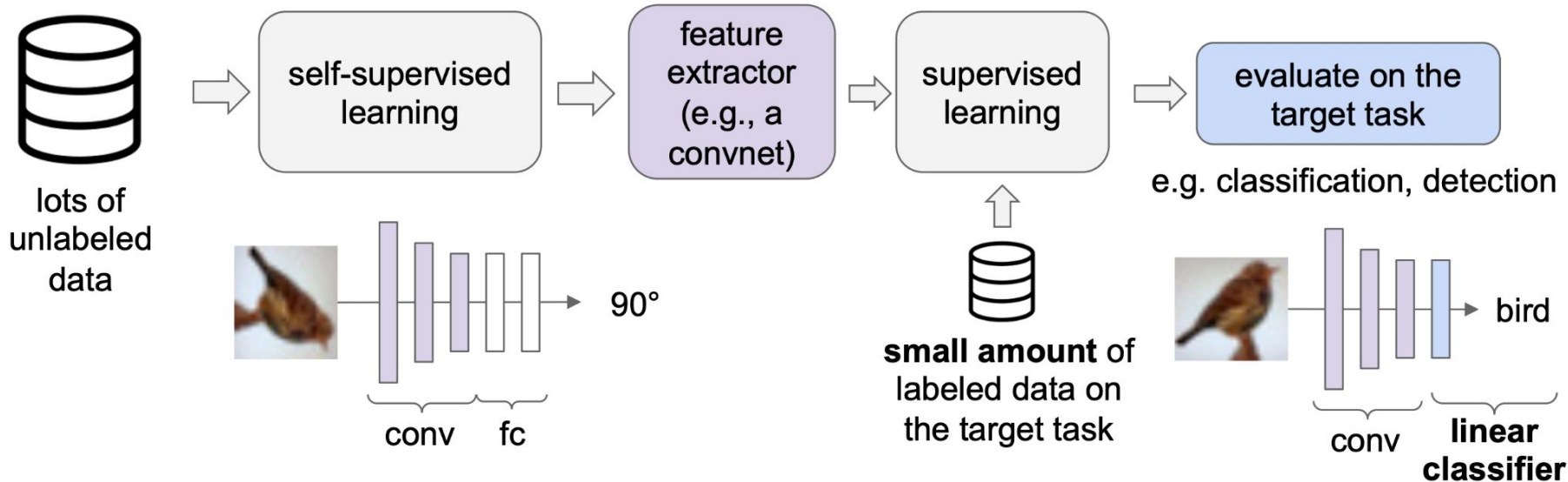
- Don't care about the performance of the self-supervised learning task
 - E.g. Image rotation prediction
- Evaluate the learned feature encoder on downstream target tasks

How to evaluate a self-supervised learning method?



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

How to evaluate a self-supervised learning method?

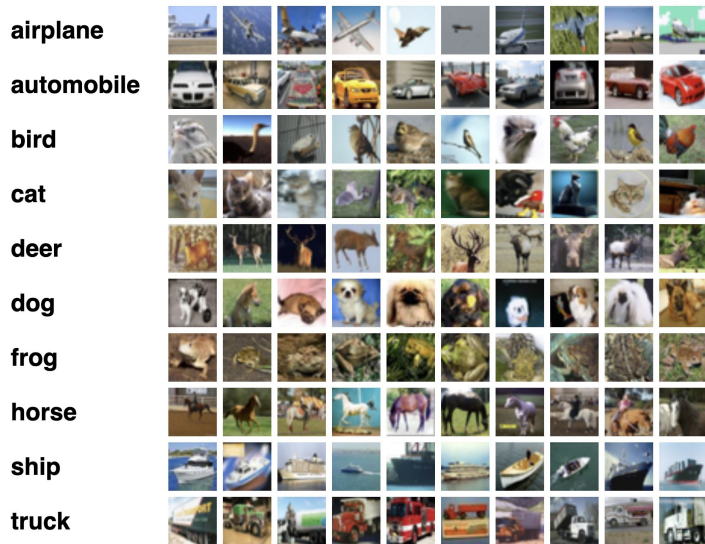


1. Learn good feature extractors from self-supervised pretext tasks, e.g., predicting image rotations

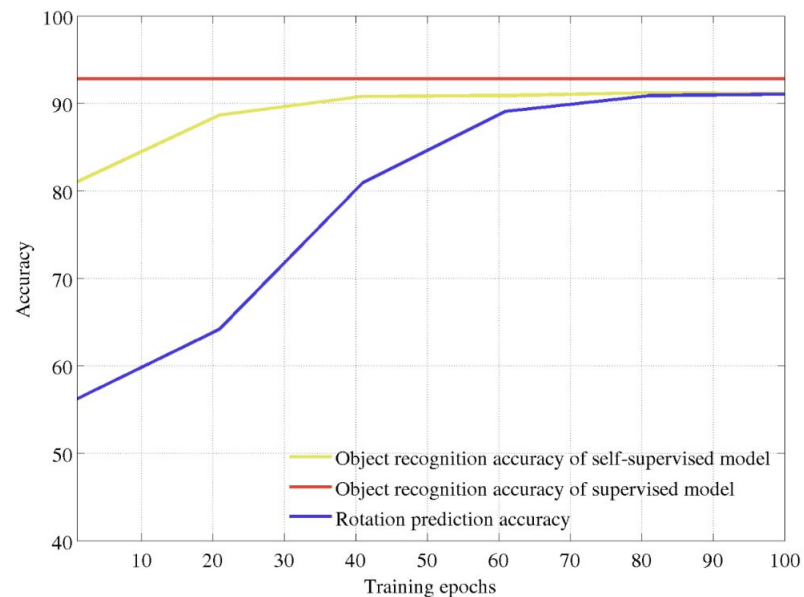
2. Attach a shallow network on the feature extractor; train the shallow network on the target task with small amount of labeled data

Self-Supervised Evaluation

- Downstream performance correlates with prefix task: rotation prediction

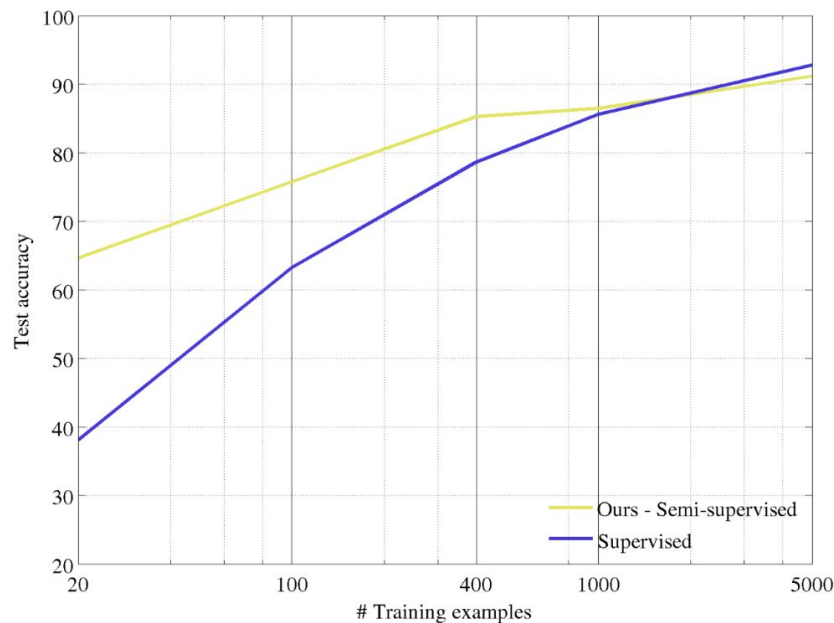


Cifar-10 Image Classification



Self-Supervised Evaluation

- Self-supervised learning outperforms supervised learning with limited data
 - Can use large volumes of unlabeled data!



Discuss

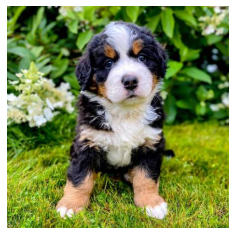
We are provided this image without labels: what are some other tasks we can do with it?

How can we perform self-supervised learning with images?

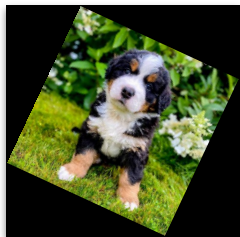


$$f(x) = ???$$

Can we learn this directly?

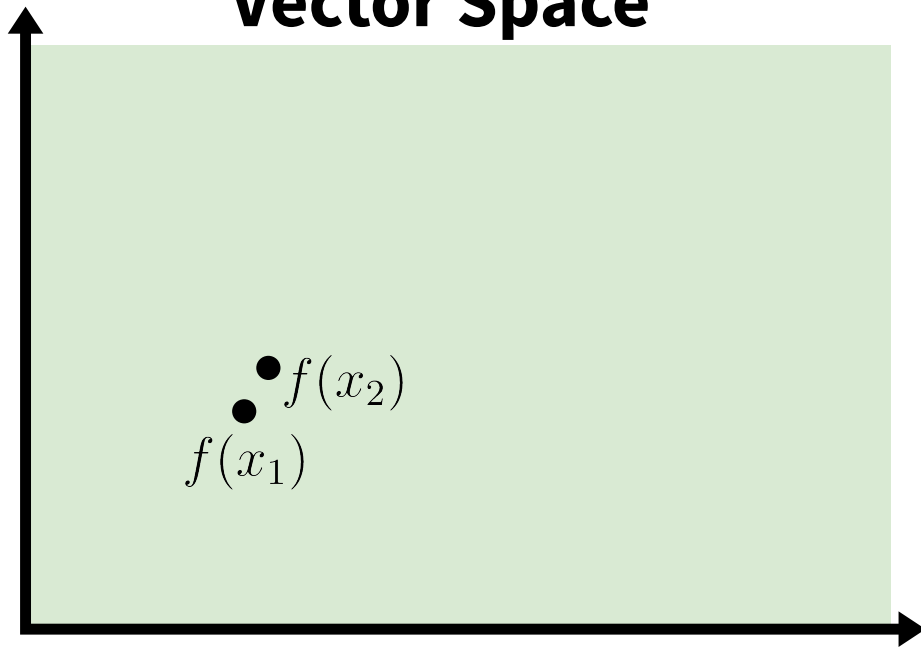


x_1



x_2

Vector Space

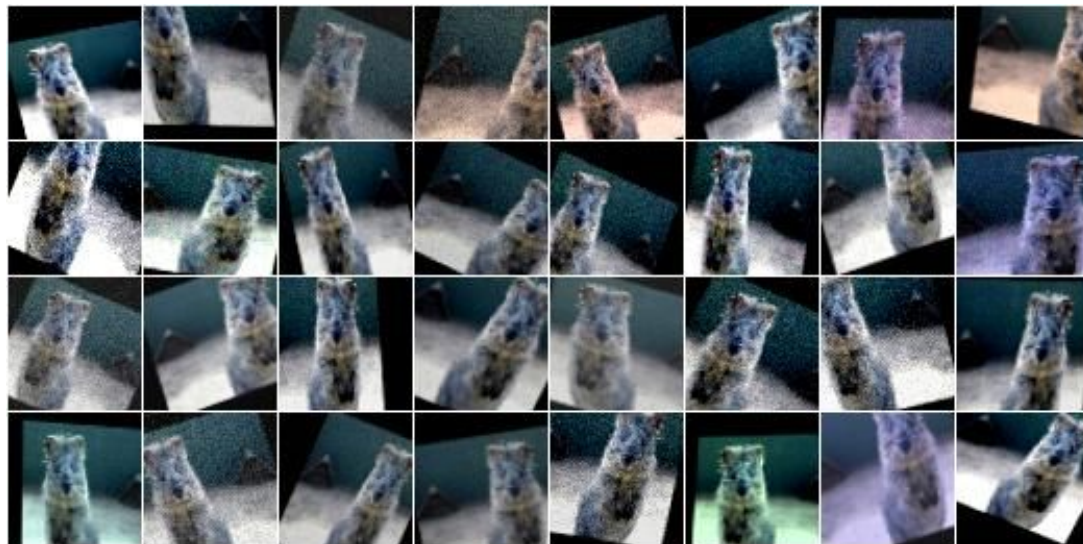


Review: Image Augmentation

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



Augmentation





x_1

x_2

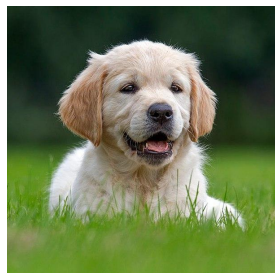
x_3

$f(x) =$ contrastive learning

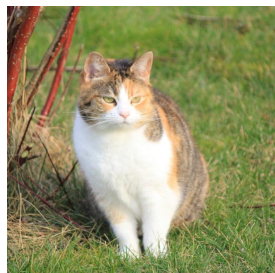
Vector Space

$f(x_1)$
● $f(x_2)$
 $f(x_3)$

All positive pairs are augmentations of the original image



x_4



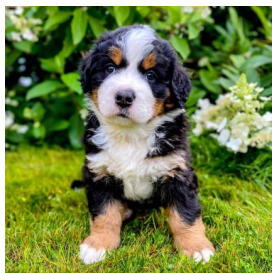
x_5



x_1



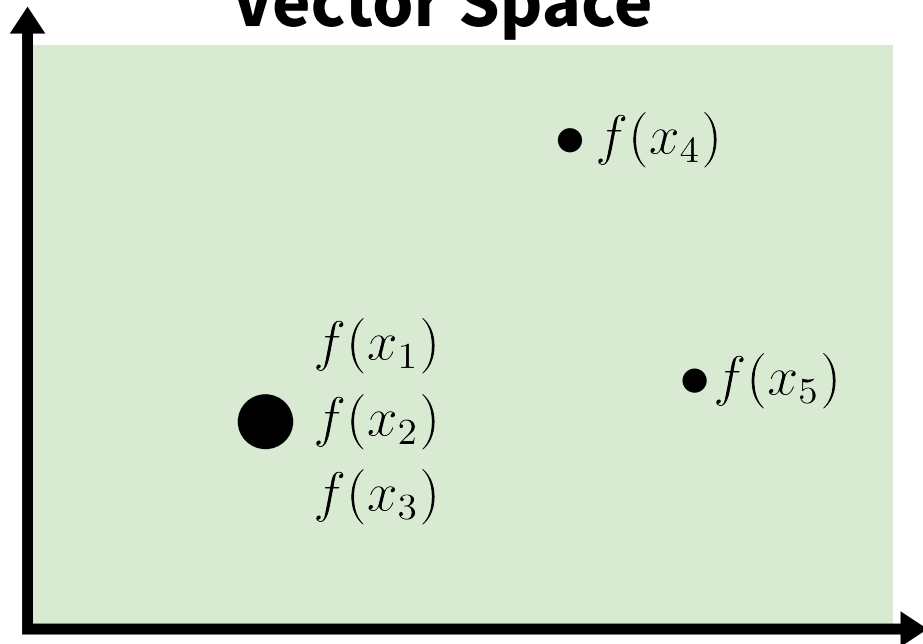
x_2



x_3

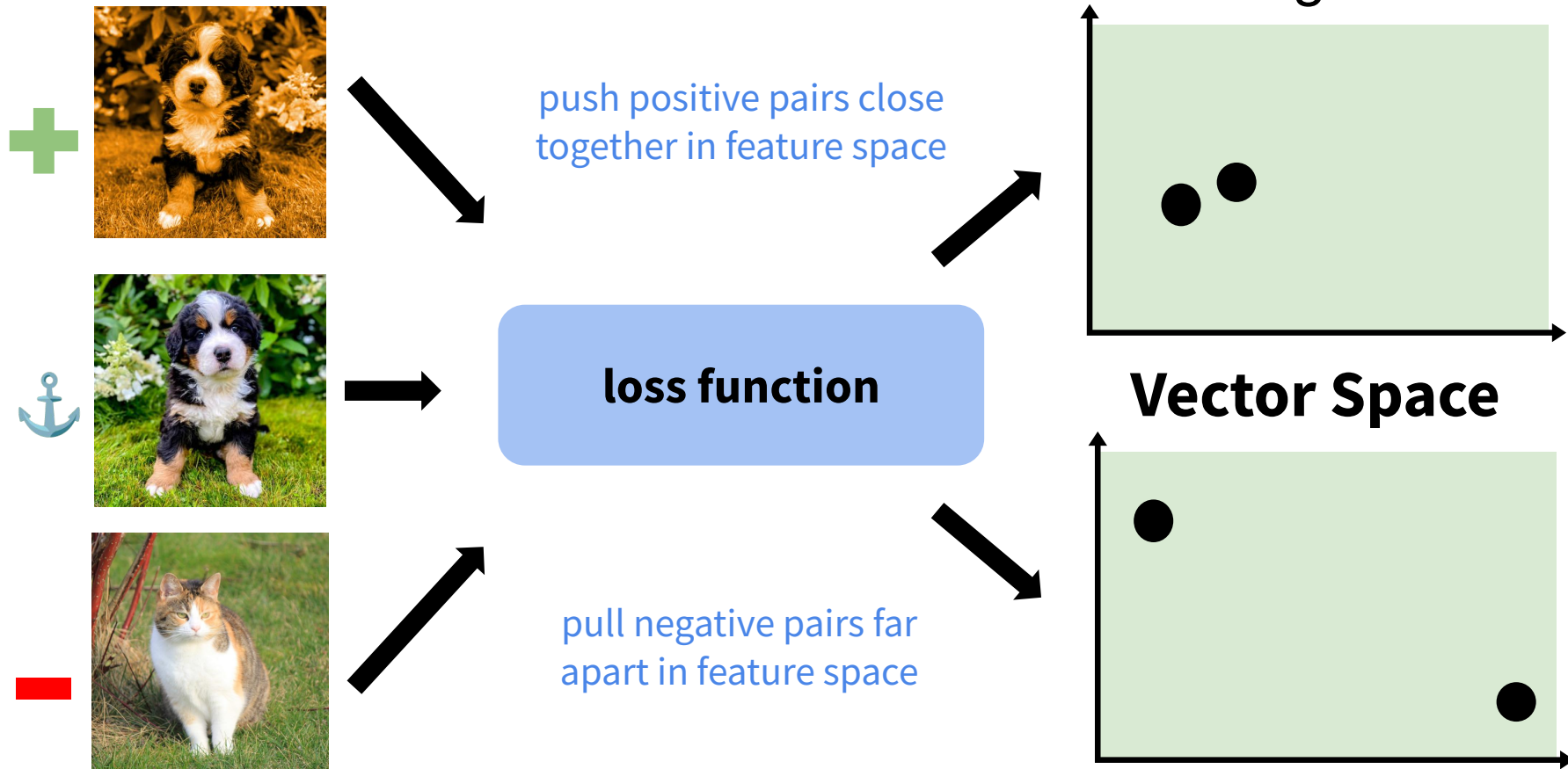
$f(x) =$ contrastive learning

Vector Space



Any other image is a negative pair

Basic Model for Contrastive Learning



Triplet loss function



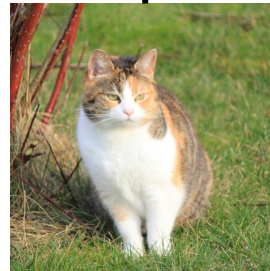
Anchor
example

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

Model



Positive pair



Negative pair

Triplet loss function

$$\ell = \max(0, \underbrace{\|f(\mathbf{x}_i) - f(\mathbf{x}^+)\|^2}_{\text{Model should map positive examples close together}} - \underbrace{\|f(\mathbf{x}_i) - f(\mathbf{x}^-)\|^2}_{\text{Model should map negative examples far apart}} + c)$$

Ensures
loss is not
negative

Model should map
positive examples
close together

Model should map
negative examples far
apart

Margin



Discuss

- Any potential problems with the triplet loss?
- Any ideas to remedy those problems

$$\ell = \max(0, \underbrace{\|f(\mathbf{x}_i) - f(\mathbf{x}^+)\|^2}_{\text{Model should map positive examples close together}} - \underbrace{\|f(\mathbf{x}_i) - f(\mathbf{x}^-)\|^2}_{\text{Model should map negative examples far apart}} + c)$$

Ensures
loss is not
negative

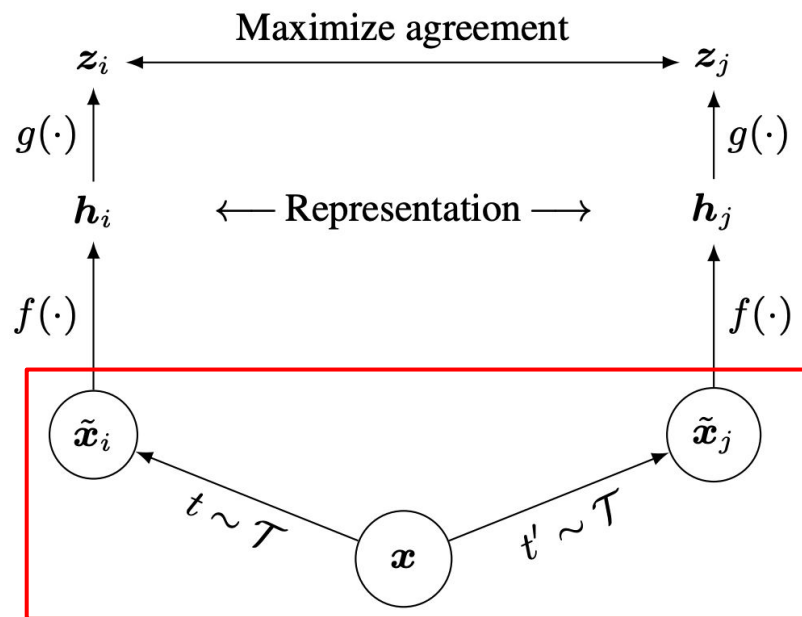
Model should map
positive examples
close together

Model should map
negative examples far
apart

Margin

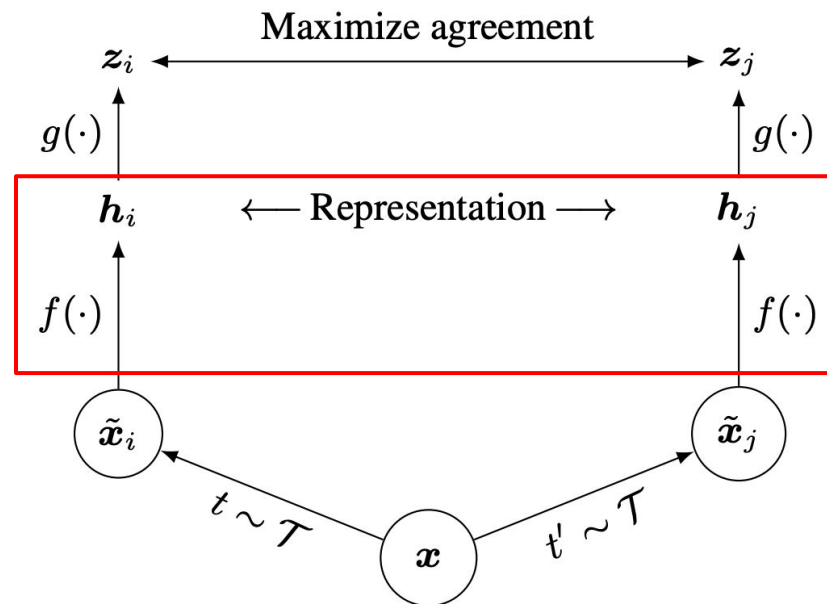
SimCLR: A Simple Contrastive Learning Framework for Images

- Sample two different augmentations of an image



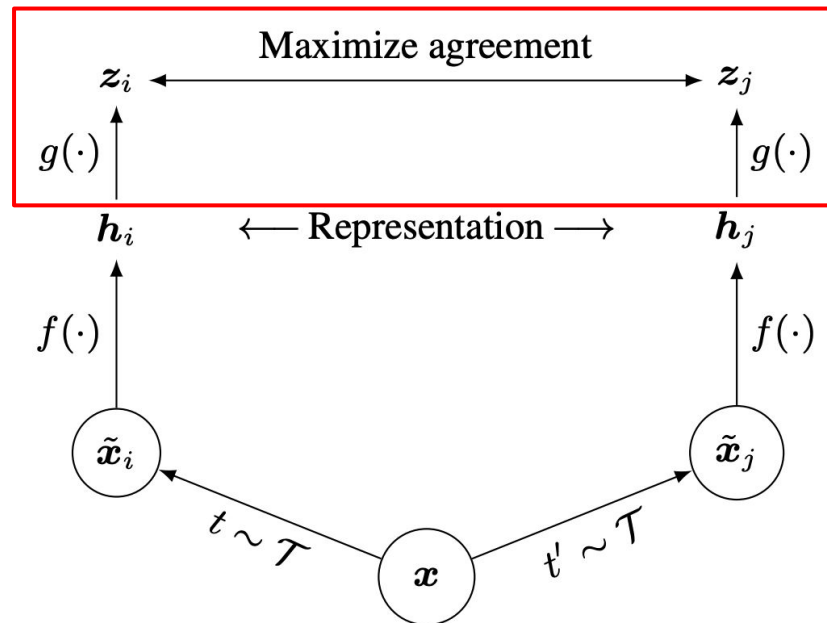
SimCLR: A Simple Contrastive Learning Framework for Images

- 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



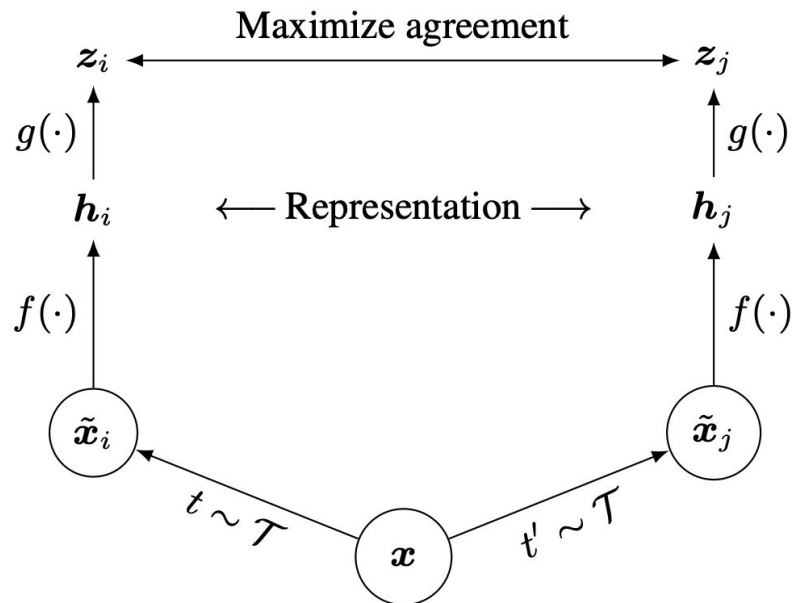
SimCLR: A Simple Contrastive Learning Framework for Images

- 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



SimCLR: A Simple Contrastive Learning Framework for Images

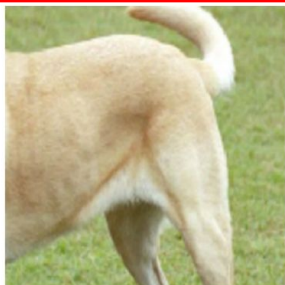
- Sample two different augmentations of an image
- Apply a base encoder to each view of the image to extract an image feature
 - e.g. ResNet
- Apply an MLP projection head to generate final representations
 - Throw away projection head after training



SimCLR Augmentations



(a) Original



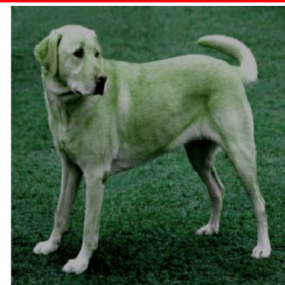
(b) Crop and resize



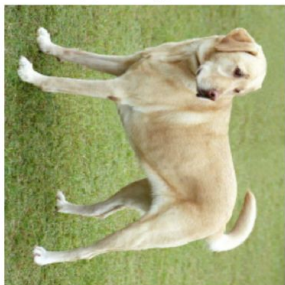
(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate $\{90^\circ, 180^\circ, 270^\circ\}$



(g) Cutout



(h) Gaussian noise



(i) Gaussian blur



(j) Sobel filtering

SimCLR Loss

- Temperature-scaled cross-entropy loss

Model should map positive
examples close together

$$\mathcal{L}_{\text{SimCLR}} = -\log \left(\frac{\overbrace{\exp(d(\mathbf{x}_i, \mathbf{x}_i^+)/\tau)}^{\text{positive}}}{\underbrace{\exp(d(\mathbf{x}_i, \mathbf{x}_i^+)/\tau) + \exp(d(\mathbf{x}_i, \mathbf{x}_i^-)/\tau)}_{\text{negative}}}} \right)$$

Model should map negative
examples far apart

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 τ , structure of f, g, \mathcal{T} .
for sampled minibatch $\{\mathbf{x}_k\}_{k=1}^N$ **do**
 for all $k \in \{1, \dots, N\}$ **do**
 draw two augmentation functions $t \sim \mathcal{T}, t' \sim \mathcal{T}$
 # the first augmentation
 $\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$
 $\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$ # representation
 $\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$ # projection
 # the second augmentation
 $\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$
 $\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$ # representation
 $\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$ # projection
 end for
 for all $i \in \{1, \dots, 2N\}$ and $j \in \{1, \dots, 2N\}$ **do**
 $s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$ # pairwise similarity
 end for
 define $\ell(i, j)$ **as** $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$
 $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$
 update networks f and g to minimize \mathcal{L}
end for
return encoder network $f(\cdot)$, and throw away $g(\cdot)$

Comparison of Loss Functions

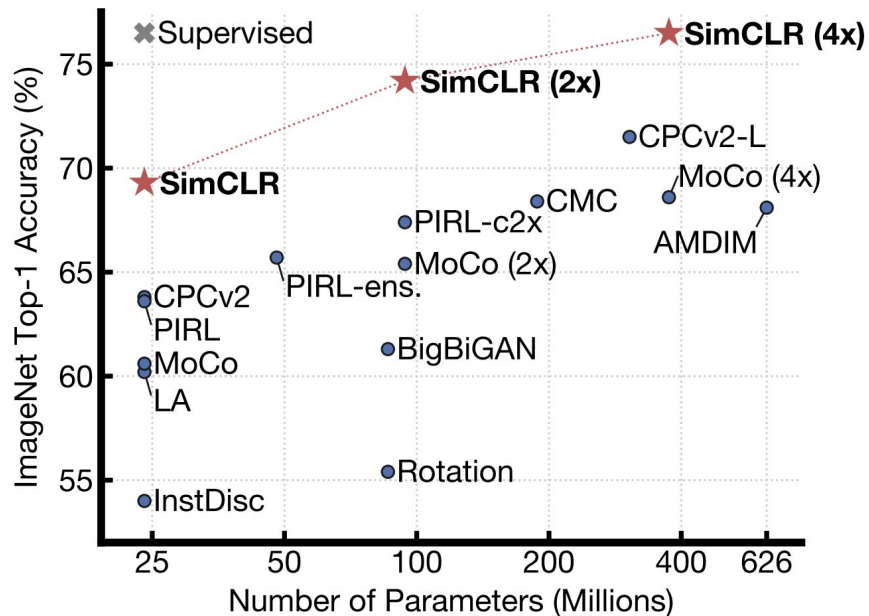
- Temperature-scaled cross entropy places more weight on hard negatives
 - Don't need to mine hard negatives

Name	Negative loss function	Gradient w.r.t. \mathbf{u}
NT-Xent	$\mathbf{u}^T \mathbf{v}^+ / \tau - \log \sum_{v \in \{v^+, v^-\}} \exp(\mathbf{u}^T \mathbf{v} / \tau)$	$(1 - \frac{\exp(\mathbf{u}^T \mathbf{v}^+ / \tau)}{Z(\mathbf{u})}) / \tau \mathbf{v}^+ - \sum_{v^-} \frac{\exp(\mathbf{u}^T \mathbf{v}^- / \tau)}{Z(\mathbf{u})} / \tau \mathbf{v}^-$
NT-Logistic	$\log \sigma(\mathbf{u}^T \mathbf{v}^+ / \tau) + \log \sigma(-\mathbf{u}^T \mathbf{v}^- / \tau)$	$(\sigma(-\mathbf{u}^T \mathbf{v}^+ / \tau)) / \tau \mathbf{v}^+ - \sigma(\mathbf{u}^T \mathbf{v}^- / \tau) / \tau \mathbf{v}^-$
Margin Triplet	$-\max(\mathbf{u}^T \mathbf{v}^- - \mathbf{u}^T \mathbf{v}^+ + m, 0)$	$\mathbf{v}^+ - \mathbf{v}^-$ if $\mathbf{u}^T \mathbf{v}^+ - \mathbf{u}^T \mathbf{v}^- < m$ else $\mathbf{0}$

Table 2. Negative loss functions and their gradients. All input vectors, i.e. \mathbf{u} , \mathbf{v}^+ , \mathbf{v}^- , are ℓ_2 normalized. NT-Xent is an abbreviation for “Normalized Temperature-scaled Cross Entropy”. Different loss functions impose different weightings of positive and negative examples.

SimCLR Results

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



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
<i>Linear evaluation:</i>												
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
<i>Fine-tuned:</i>												
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

Effect Of Projection Head

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

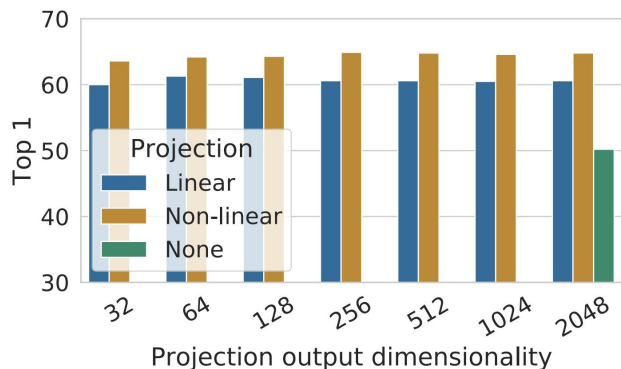
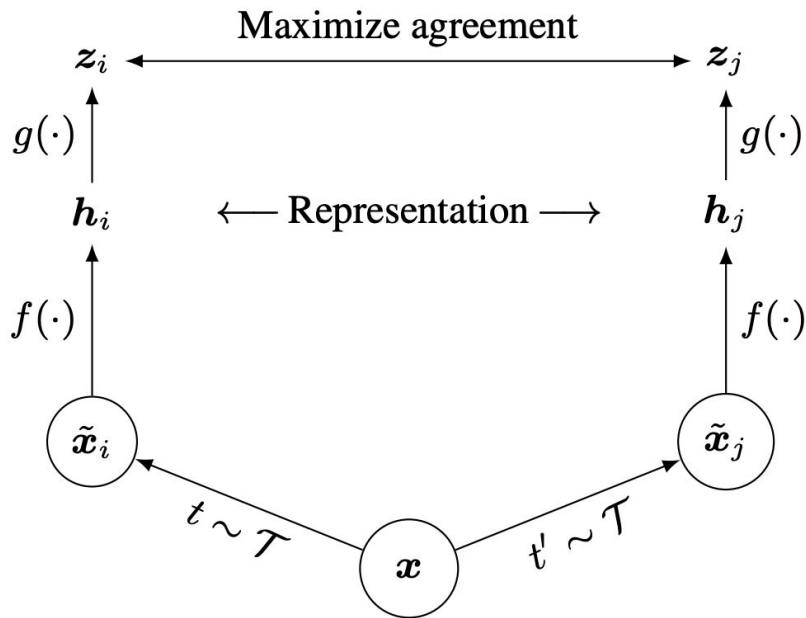
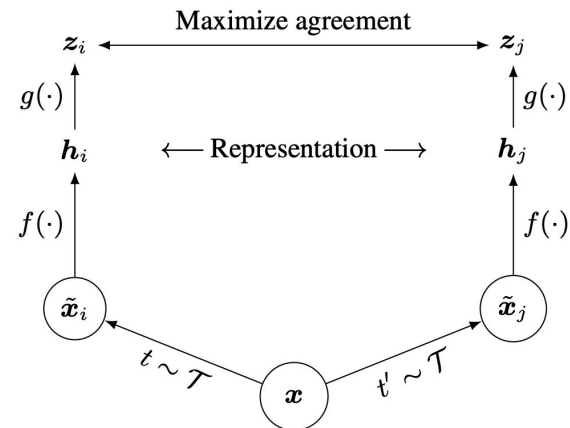
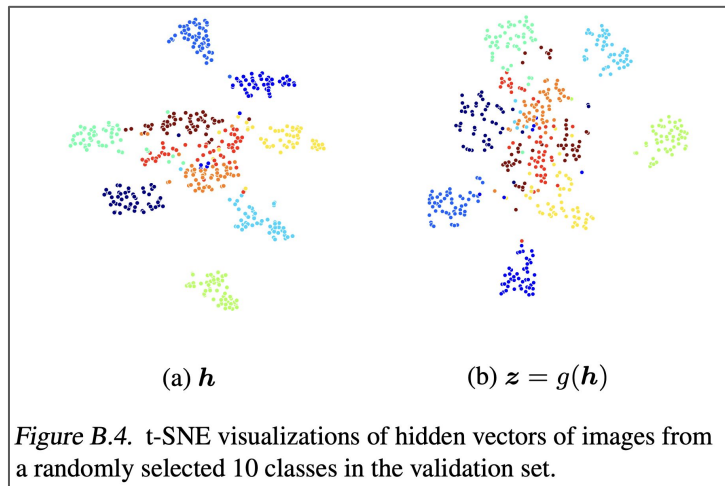


Figure 8. Linear evaluation of representations with different projection heads $g(\cdot)$ and various dimensions of $z = g(\mathbf{h})$. The representation \mathbf{h} (before projection) is 2048-dimensional here.



Effect Of Projection Head

- Projects data to “augmentation-invariant” representation
 - Features less useful for downstream tasks



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

Impact of Loss Function

- Proposed loss outperforms the margin loss
 - Even with negative mining
- L2 normalization is useful
- Sensitive to cross-entropy temperature

Margin	NT-Logi.	Margin (sh)	NT-Logi.(sh)	NT-Xent
50.9	51.6	57.5	57.9	63.9

Table 4. Linear evaluation (top-1) for models trained with different loss functions. “sh” means using semi-hard negative mining.

ℓ_2 norm?	τ	Entropy	Contrastive acc.	Top 1
Yes	0.05	1.0	90.5	59.7
	0.1	4.5	87.8	64.4
	0.5	8.2	68.2	60.7
	1	8.3	59.1	58.0
No	10	0.5	91.7	57.2
	100	0.5	92.1	57.0

Table 5. Linear evaluation for models trained with different choices of ℓ_2 norm and temperature τ for NT-Xent loss. The contrastive distribution is over 4096 examples.

Impact of Batch Size

- Requires large batches
 - Harder negatives!

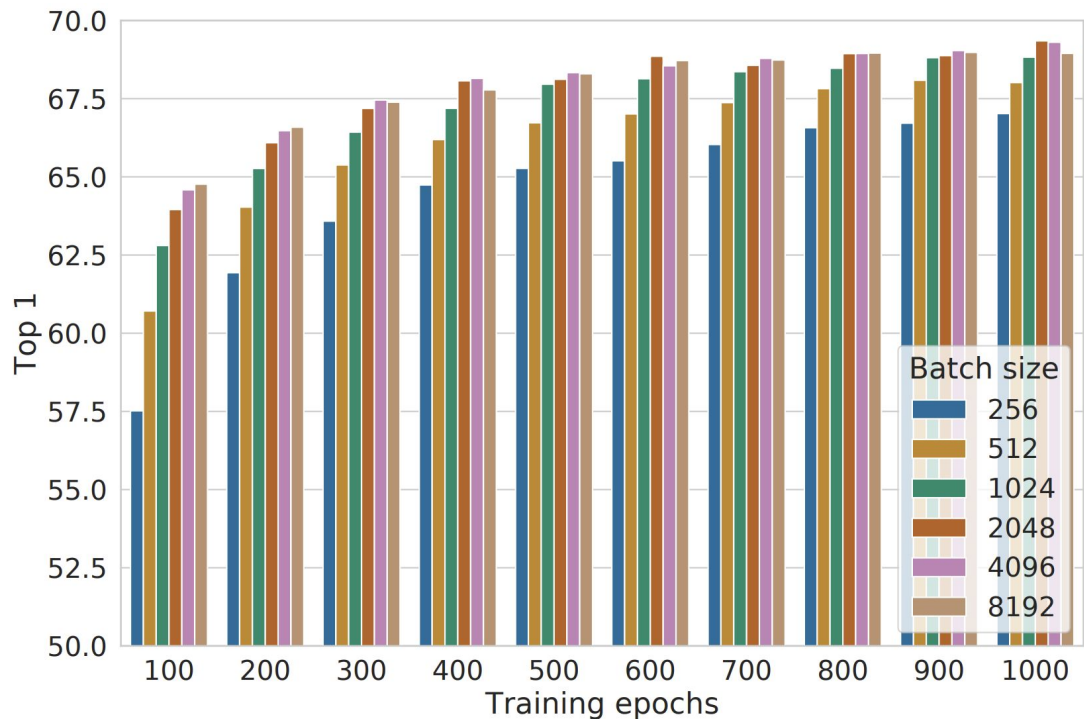
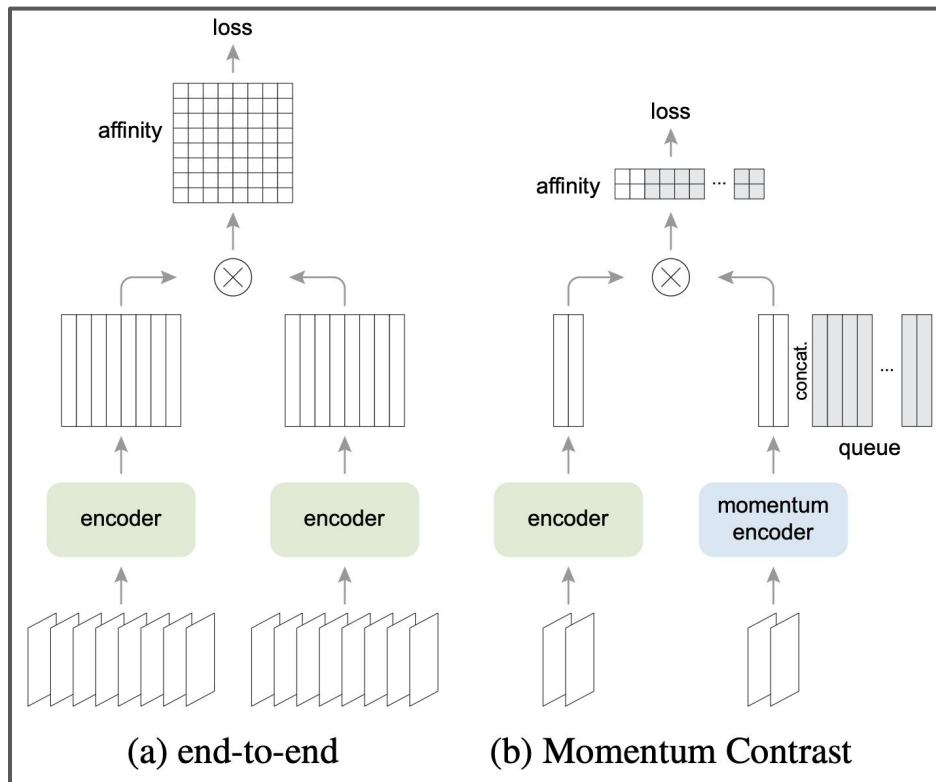


Figure 9. Linear evaluation models (ResNet-50) trained with different batch size and epochs. Each bar is a single run from scratch.¹⁰

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

$$\theta_k \leftarrow m\theta_k + (1 - m)\theta_q$$



MoCo v2

- MoCo v2: MoCo with some tricks from SimCLR
 - Stronger augmentations
 - MLP projection head
- Outperform SimCLR with modest batch sizes
 - Large numbers of negatives available from the queue

case	MLP	unsup. pre-train			batch	ImageNet acc.
		aug+	cos	epochs		
MoCo v1 [6]				200	256	60.6
SimCLR [2]	✓	✓	✓	200	256	61.9
SimCLR [2]	✓	✓	✓	200	8192	66.6
MoCo v2	✓	✓	✓	200	256	67.5
<i>results of longer unsupervised training follow:</i>						
SimCLR [2]	✓	✓	✓	1000	4096	69.3
MoCo v2	✓	✓	✓	800	256	71.1

Table 2. **MoCo vs. SimCLR**: ImageNet linear classifier accuracy (**ResNet-50, 1-crop 224×224**), trained on features from unsupervised pre-training. “aug+” in SimCLR includes blur and stronger color distortion. SimCLR ablations are from Fig. 9 in [2] (we thank the authors for providing the numerical results).

mechanism	batch	memory / GPU	time / 200-ep.
MoCo	256	5.0G	53 hrs
end-to-end	256	7.4G	65 hrs
end-to-end	4096	93.0G [†]	n/a

Table 3. **Memory and time cost** in 8 V100 16G GPUs, implemented in PyTorch. [†]: based on our estimation.

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