

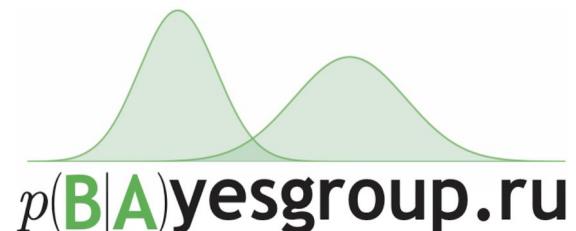
# Contrastive Self-Supervised Learning for Visual Representations (CPC and SimCLR)

Andrei Atanov



NATIONAL RESEARCH  
UNIVERSITY

**SAMSUNG**  
**Research**

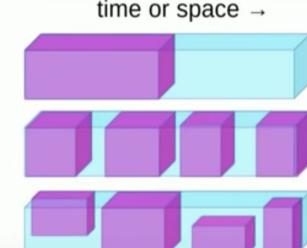


# Self-Supervised Learning. Motivation.

- Learn the model of the world prior to learning a specific task
  - Do this by predicting hidden parts of an input
  - More bits of information per training sample
- 
- Use less labeled data to train on a downstream task.

Y. LeCun

### Self-Supervised Learning = Filling in the Blanks

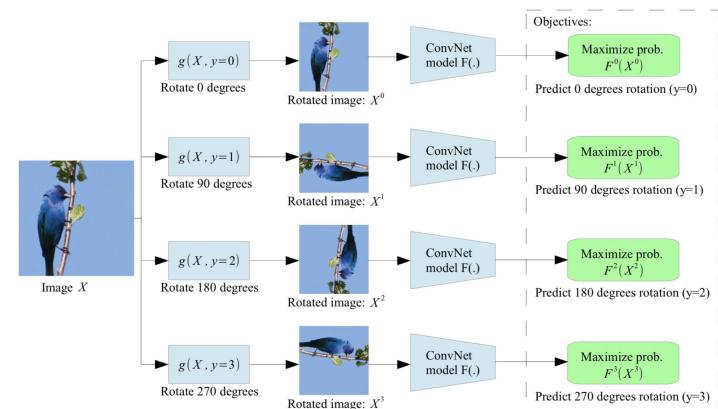


- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **masked** from the **visible**.
- ▶ Predict the **any occluded part** from **all available parts**.
- ▶ Pretend there is a part of the input you don't know and predict that.
- ▶ Reconstruction = SSL when any part could be known or unknown

Yann Le Cunn, AAAI 2020 Keynotes Turing Award  
Winners Event

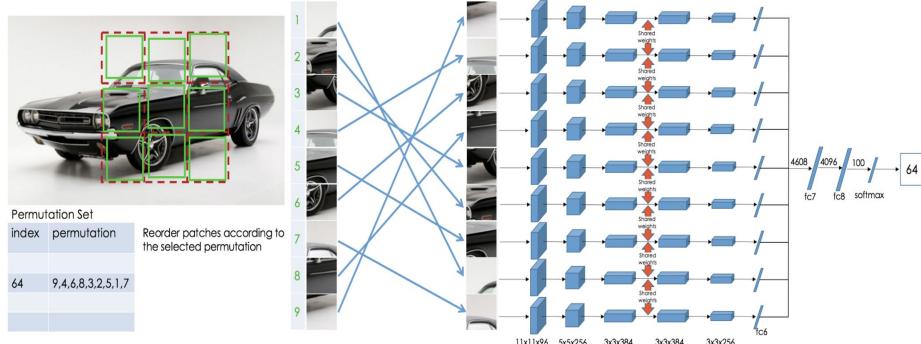
# Pretext approaches

- Transform an image  $I$  with a random transformation
- Given  $I^t$ , predict the transformation  $t$



## Rotation

[ <https://arxiv.org/abs/1803.07728> ]

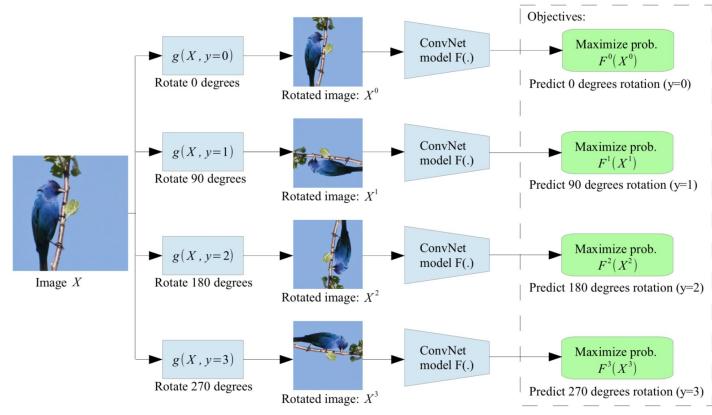


## Jigsaw puzzle

[ <https://arxiv.org/abs/1603.09246> ]

# Pretext approaches

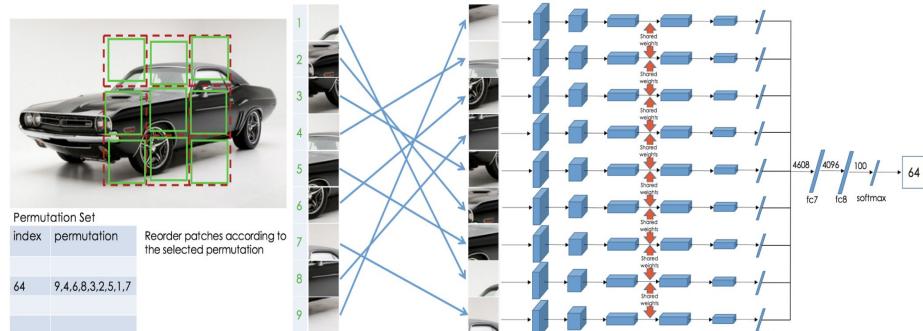
- Transform an image  $X$  with a random transformation
- Given  $X$ , predict the transformation  $t$



Rotation

[ <https://arxiv.org/abs/1803.07728> ]

- Don't learn representations directly.
- Want representations invariant to transformations



Jigsaw puzzle

[ <https://arxiv.org/abs/1603.09246> ]

# Contrastive Prediction Coding

- Predict “future” patches of an image:



Problems:

- MSE loss doesn't capture data variability.
- Learning generative model over pixels is time and computational consuming.
- Need to model local low-level information.

# Contrastive Prediction Coding

- Encode patches into a vector space and predict next patch in this space

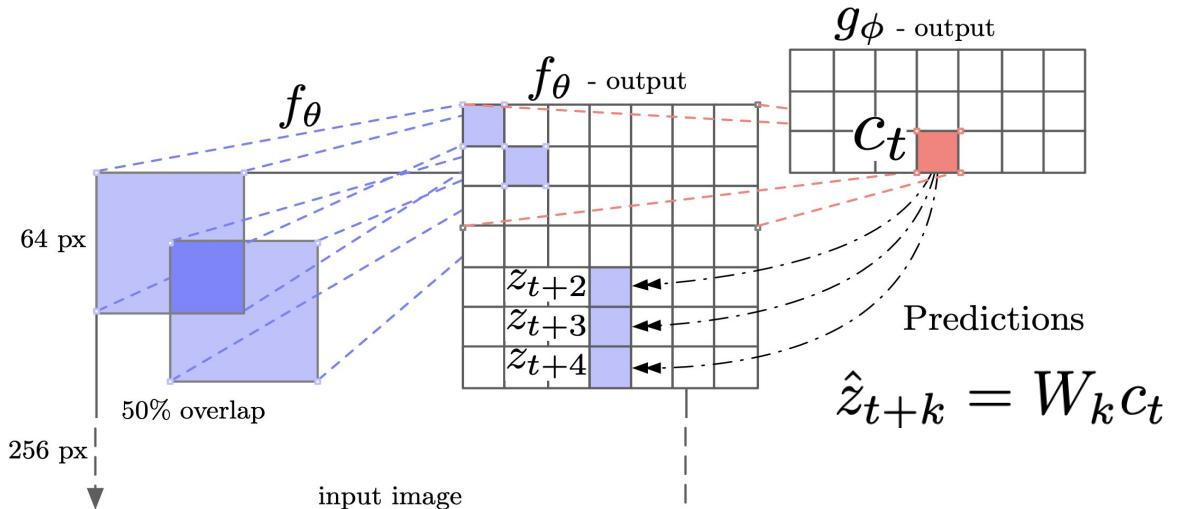
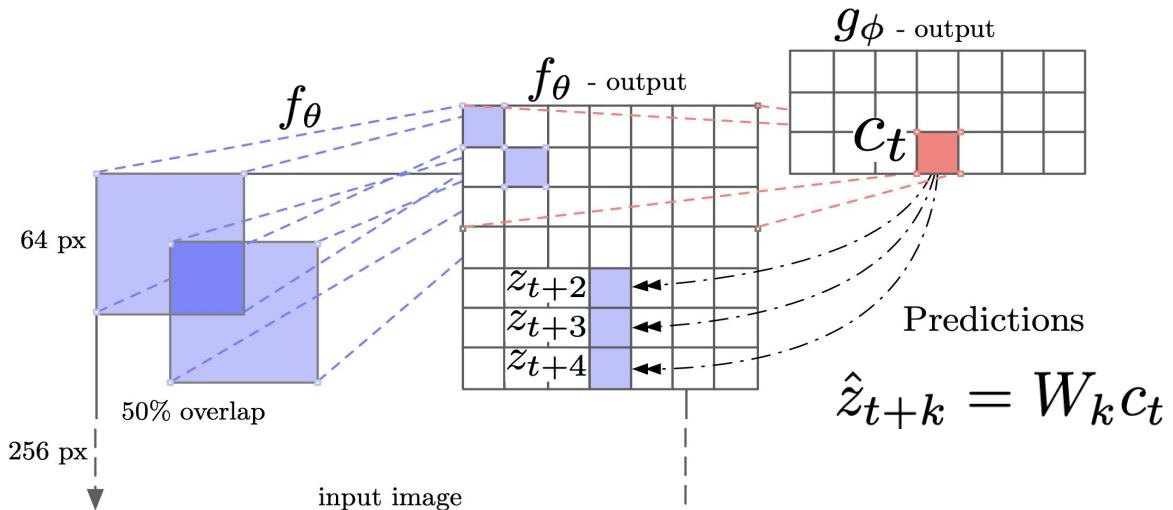


Figure 4: Visualization of Contrastive Predictive Coding for images (2D adaptation of Figure 1).

# Contrastive Prediction Coding

- Encode patches into a vector space and predict next patch in this space



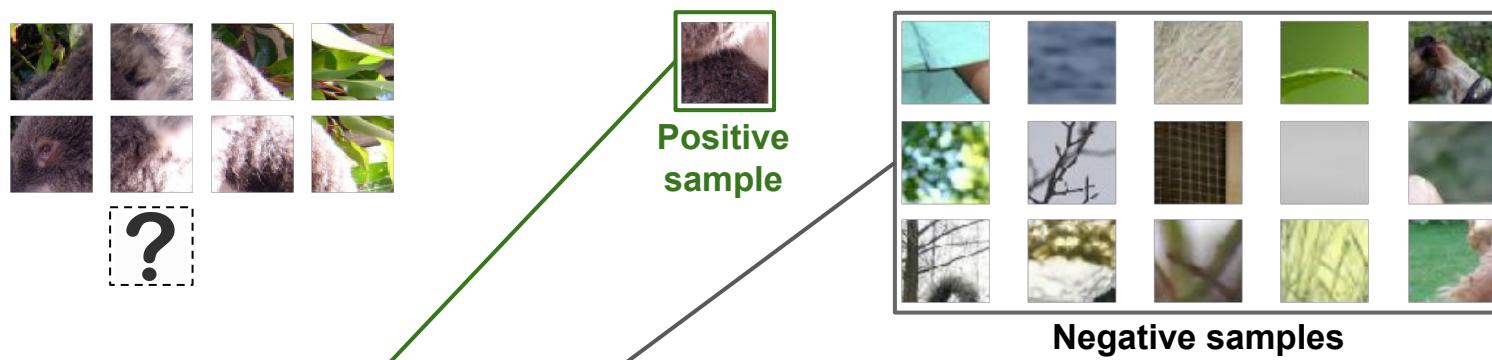
Problems with MSE?

$$\|\hat{z}_t - z_t\|^2$$

Figure 4: Visualization of Contrastive Predictive Coding for images (2D adaptation of Figure 1).

# Contrastive loss

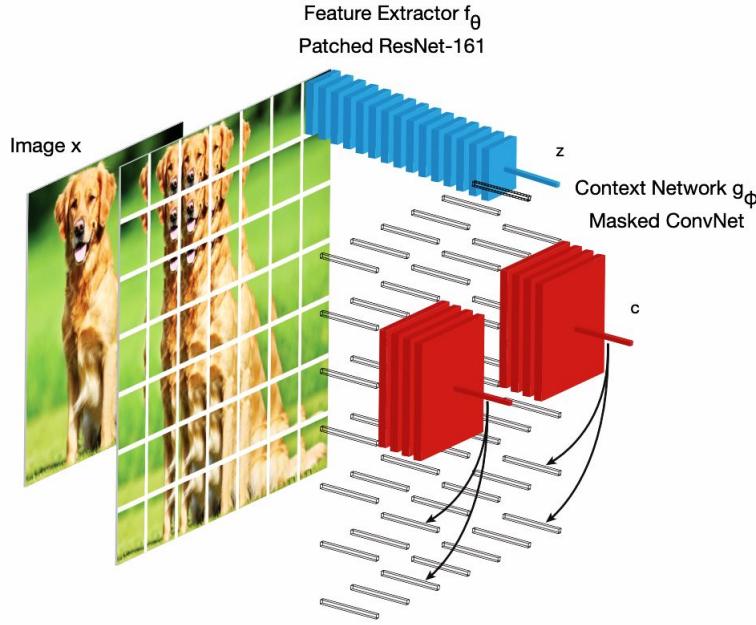
- MSE loss  $\|\hat{z}_t - z_t\|^2$  would lead to a constant prediction.
- Generative models are hard to train.
- Instead, use Noise Contrastive Estimation loss:



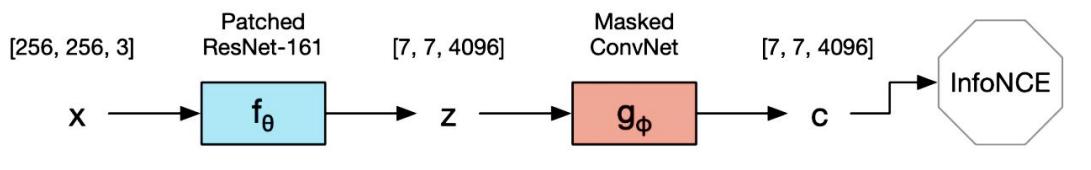
$$\log p(z_t | \hat{z}_t, \{z_l\}) = \log \left( \frac{\exp(\hat{z}_t^T z_t)}{\exp(\hat{z}_t^T z_t) + \sum_l \exp(\hat{z}_t^T z_l)} \right) \rightarrow \max$$

Predicted vector

# Contrastive Prediction Coding



**Self-supervised  
pre-training**  
100% images; 0% labels



Pre-training

# Linear Evaluation

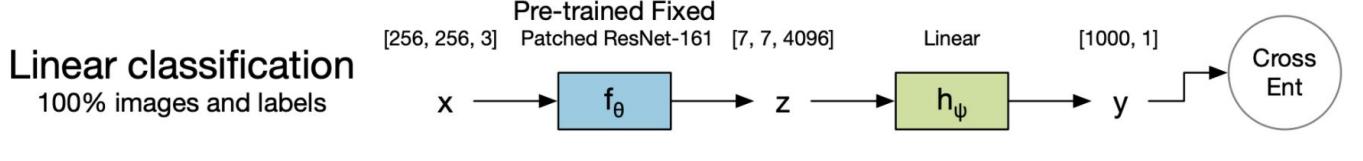


Table 1: Linear classification accuracy, and comparison to other self-supervised methods. In all cases the feature extractor is optimized in an unsupervised manner (using one of the methods listed below), and a linear classifier is trained on top using all labels in the ImageNet dataset.

Method	Architecture	Parameters (M)	Top-1	Top-5
<b>Methods using ResNet-50:</b>				
Local Aggregation [66]	ResNet-50	24	60.2	-
Momentum Contrast [25]	ResNet-50	24	60.6	-
<b>CPC v2</b>	ResNet-50	24	<b>63.8</b>	<b>85.3</b>
<b>Methods using different architectures:</b>				
Multi-task [13]	ResNet-101	28	-	69.3
Rotation [32]	RevNet-50 $\times 4$	86	55.4	-
CPC v1 [58]	ResNet-101	28	48.7	73.6
BigBiGAN [15]	RevNet-50 $\times 4$	86	61.3	81.9
AMDIM [5]	Custom-103	626	68.1	-
CMC [57]	ResNet-50 $\times 2$	188	68.4	88.2
Momentum Contrast [25]	ResNet-50 $\times 4$	375	68.6	-
<b>CPC v2</b>	ResNet-161	305	<b>71.5</b>	<b>90.1</b>

# Linear Evaluation

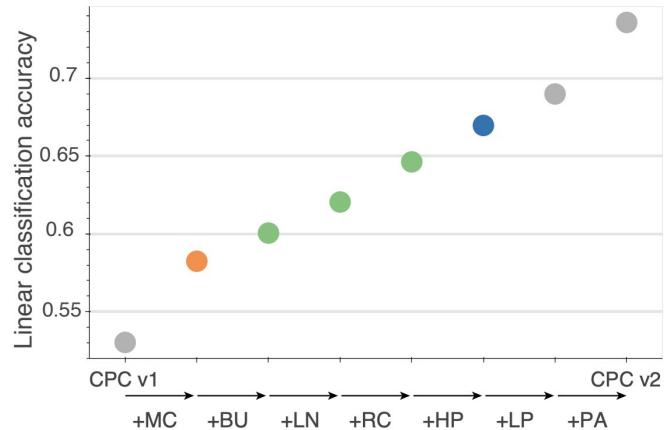
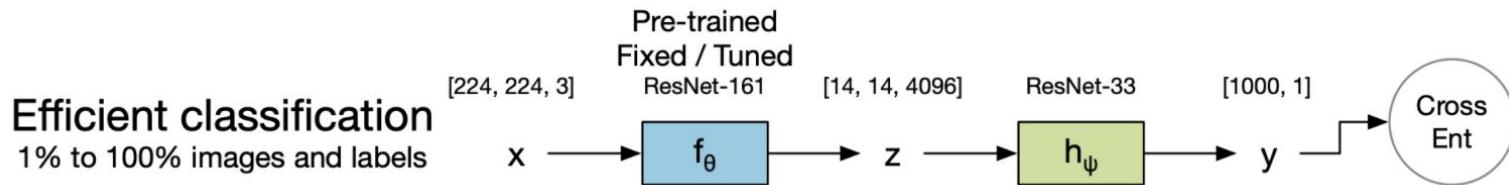
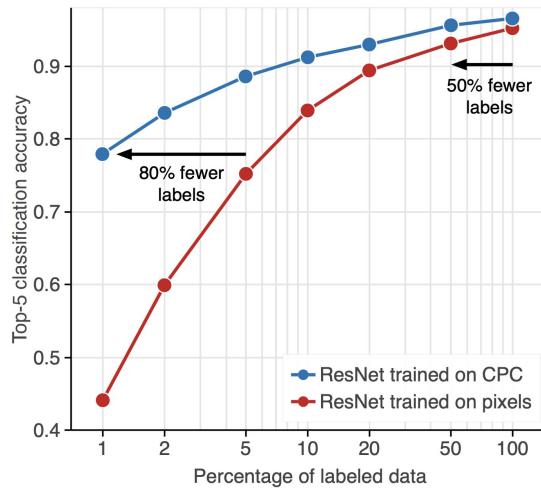


Figure 3: Linear classification performance of new variants of CPC, which incrementally add a series of modifications. MC: model capacity. BU: bottom-up spatial predictions. LN: layer normalization. RC: random color-dropping. HP: horizontal spatial predictions. LP: larger patches. PA: further patch-based augmentation. We use color to indicate the number of spatial predictions used (orange, green, blue for 1, 2 and 4 directions). Note that these accuracies are evaluated on a custom validation set and are therefore not directly comparable to the results we report and compare to.

# Semi-Supervised setting



Method Labeled data	Architecture	Top-5 accuracy				
		1%	5%	10%	50%	100%
† Supervised baseline	ResNet-200	44.1	75.2*	83.9	93.1	95.2#
<b>Methods using label-propagation:</b>						
Pseudolabeling [63]	ResNet-50	51.6	-	82.4	-	-
VAT + Entropy Minimization [63]	ResNet-50	47.0	-	83.4	-	-
Unsup. Data Augmentation [61]	ResNet-50	-	-	88.5	-	-
Rotation + VAT + Ent. Min. [63]	ResNet-50 $\times 4$	-	-	<b>91.2</b>	-	95.0
<b>Methods using representation learning only:</b>						
Instance Discrimination [60]	ResNet-50	39.2	-	77.4	-	-
Rotation [63]	ResNet-152 $\times 2$	57.5	-	86.4	-	-
ResNet on BigBiGAN (fixed)	RevNet-50 $\times 4$	55.2	73.7	78.8	85.5	87.0
ResNet on AMDIM (fixed)	Custom-103	67.4	81.8	85.8	91.0	92.2
ResNet on CPC v2 (fixed)	ResNet-161	77.1	87.5	90.5	95.0	96.2
ResNet on CPC v2 (fine-tuned)	ResNet-161	<b>77.9*</b>	<b>88.6</b>	<b>91.2</b>	<b>95.6#</b>	<b>96.5</b>



# Transfer Learning

**Transfer learning**  
100% images and labels



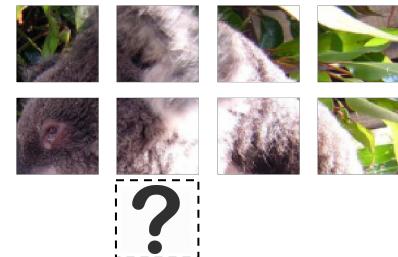
Method	Architecture	mAP
<i>Transfer from labeled data:</i>		
Supervised baseline	ResNet-152	74.7
<i>Transfer from unlabeled data:</i>		
Exemplar [17] by [13]	ResNet-101	60.9
Motion Segmentation [47] by [13]	ResNet-101	61.1
Colorization [64] by [13]	ResNet-101	65.5
Relative Position [14] by [13]	ResNet-101	66.8
Multi-task [13]	ResNet-101	70.5
Instance Discrimination [60]	ResNet-50	65.4
Deep Cluster [7]	VGG-16	65.9
Deeper Cluster [8]	VGG-16	67.8
Local Aggregation [66]	ResNet-50	69.1
Momentum Contrast [25]	ResNet-50	74.9
Faster-RCNN trained on CPC v2	ResNet-161	<b>76.6</b>

# CPC in a nutshell

$$L_t = -\log \left( \frac{\exp(\hat{z}_t^T z_t)}{\exp(\hat{z}_t^T z_t) + \sum_l \exp(\hat{z}_t^T z_l)} \right)$$

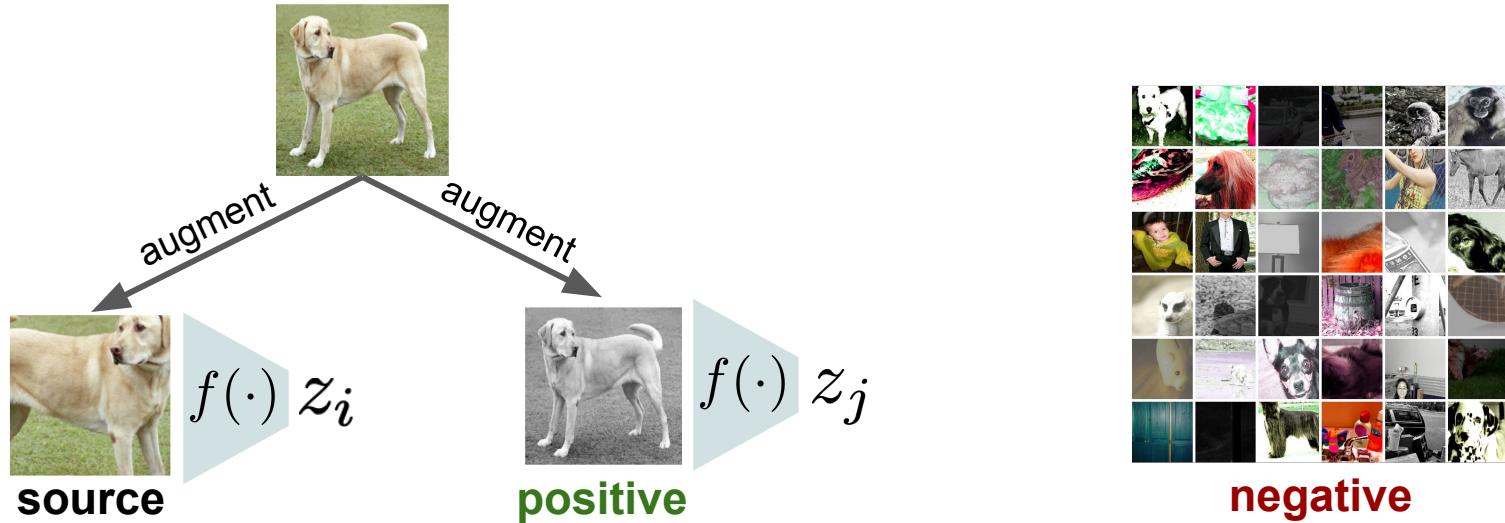
**source**      **positive**      **negative**

- **source** --- predicted embedding for a patch
- **positive** --- true embedding
- **negative** --- embeddings of random patches



Other (simple) ways to sample source, positive and negatives?

# A Simple Framework for Contrastive Learning (SimCLR)



$$L_{i,j} = -\log \left( \frac{\exp(z_i^T z_j)}{\exp(z_i^T z_j) + \sum_l \exp(z_i^T z_l)} \right)$$

# SimCLR

- Given an image find its augmentation among other images.
- Use the rest of batch as negative examples instead of a memory bank  
(need large batches, N=2-4K)

Batch of size N:



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

# Data augmentation strategy



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



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



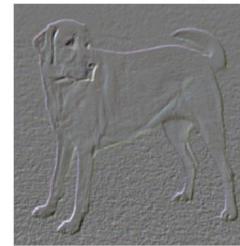
(g) Cutout



(h) Gaussian noise



(i) Gaussian blur



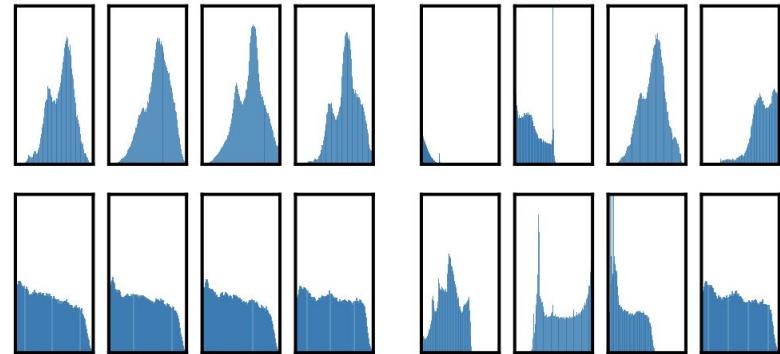
(j) Sobel filtering

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)

# Data augmentation strategy



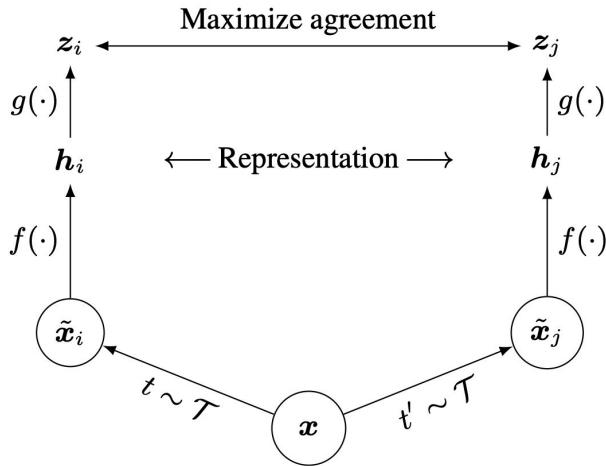
Figure 5. Linear evaluation (ImageNet top-1 accuracy) under individual or composition of data augmentations, applied only to one branch. For all columns but the last, diagonal entries correspond to single transformation, and off-diagonals correspond to composition of two transformations (applied sequentially). The last column reflects the average over the row.



(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.

# A Simple Framework for Contrastive Learning



*Figure 2.* A simple framework for contrastive learning of visual representations. Two separate data augmentation operators are sampled from the same family of augmentations ( $t \sim \mathcal{T}$  and  $t' \sim \mathcal{T}$ ) and applied to each data example to obtain two correlated views. A base encoder network  $f(\cdot)$  and a projection head  $g(\cdot)$  are trained to maximize agreement using a contrastive loss. After training is completed, we throw away the projection head  $g(\cdot)$  and use encoder  $f(\cdot)$  and representation  $h$  for downstream tasks.

# Projection head

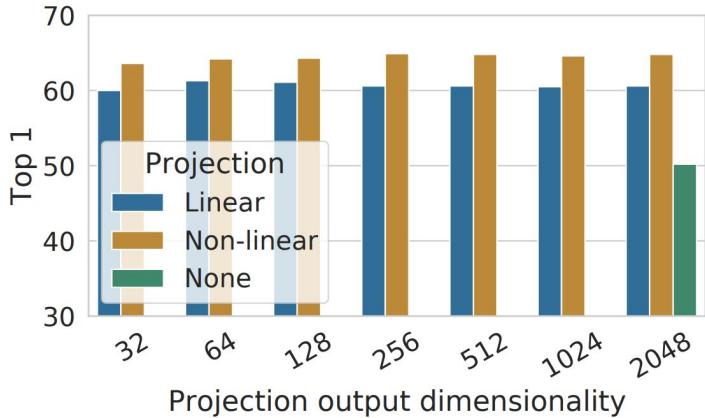


Figure 8. Linear evaluation of representations with different projection heads  $g(\cdot)$  and various dimensions of  $z = g(\mathbf{h})$ . The representation  $\mathbf{h}$  (before projection) is 2048-dimensional here.

What to predict?	Random guess	Representation $\mathbf{h}$	Representation $g(\mathbf{h})$
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  $\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}$ ), Gaussian noise, and Sobel filtering transformation during the pretraining for the last three rows. Both  $\mathbf{h}$  and  $g(\mathbf{h})$  are of the same dimensionality, i.e. 2048.

# Linear evaluation

Method	Architecture	Param.	Top 1	Top 5
<i>Methods using ResNet-50:</i>				
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	<b>69.3</b>	<b>89.0</b>
<i>Methods using other architectures:</i>				
Rotation	RevNet-50 (4×)	86	55.4	-
BigBiGAN	RevNet-50 (4×)	86	61.3	81.9
AMDIM	Custom-ResNet	626	68.1	-
CMC	ResNet-50 (2×)	188	68.4	88.2
MoCo	ResNet-50 (4×)	375	68.6	-
CPC v2	ResNet-161 (*)	305	71.5	90.1
SimCLR (ours)	ResNet-50 (2×)	94	74.2	92.0
SimCLR (ours)	ResNet-50 (4×)	375	<b>76.5</b>	<b>93.2</b>

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

# Semi-supervised evaluation

Method	Architecture	Label fraction		
		1%	10%	Top 5
Supervised baseline	ResNet-50	48.4	80.4	
<i>Methods using other label-propagation:</i>				
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×)	-	91.2	
<i>Methods using representation learning only:</i>				
InstDisc	ResNet-50	39.2	77.4	
BigBiGAN	RevNet-50 (4×)	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×)	83.0	91.2	
SimCLR (ours)	ResNet-50 (4×)	<b>85.8</b>	<b>92.6</b>	

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

# Conclusion

- Contrastive losses reduce the “object prediction” problem to classification
- The main difference is how to construct positive and negative samples
- Require large amount of negative samples and long training  
(best SimCLR takes 1.5 hours of 128 TPUs for 100 epochs)
- Contrastive self-supervised pretraining reaches (and outperforms) the level of supervised pretraining
- Performance saturated with increasing number of training examples