

# Weakly- & Self-Supervised Learning

Lantian Li

2020.02.24

# Weakly Supervised Learning

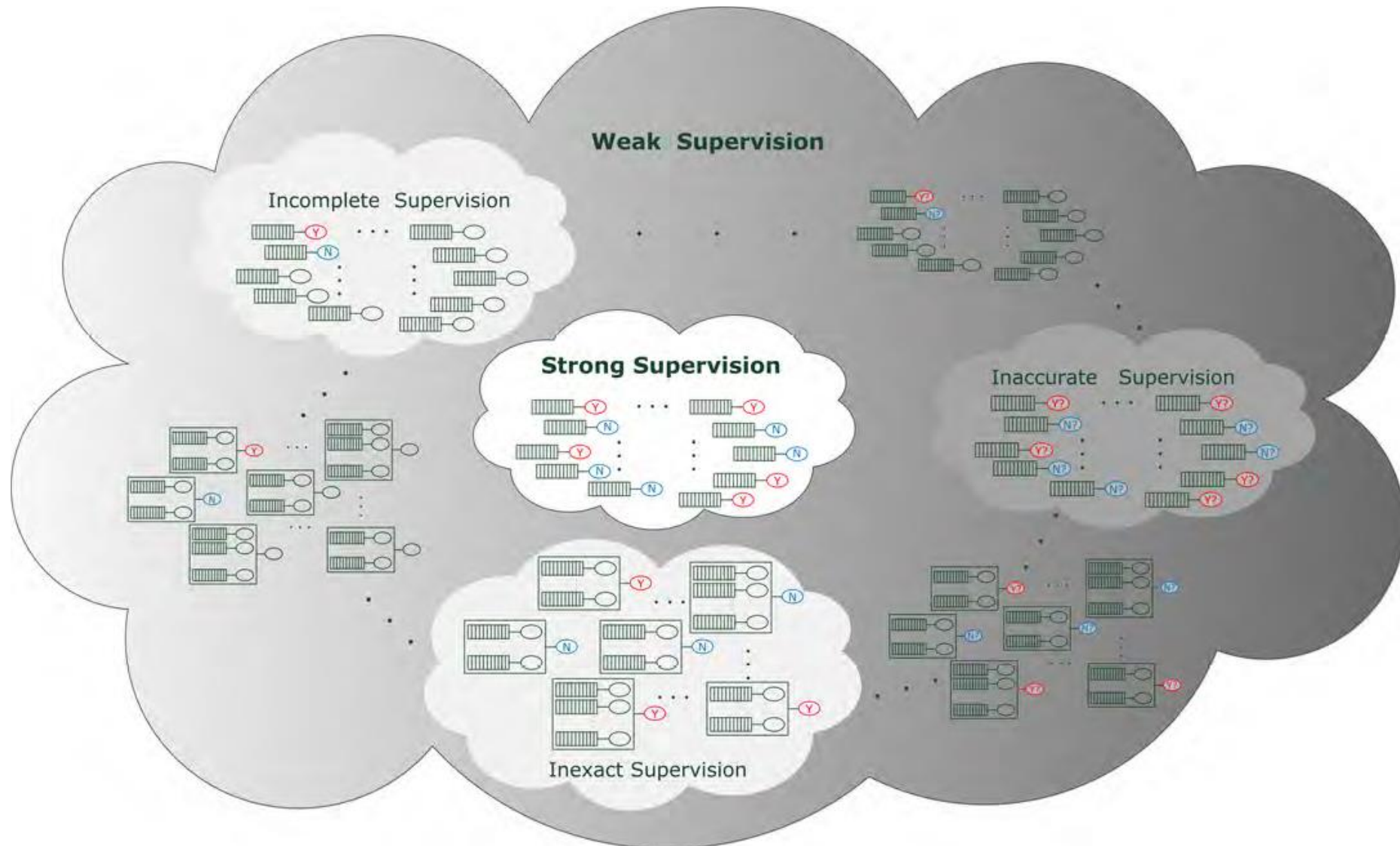
# Supervised learning

- Concepts
  - learning from a large number of examples
  - each example has its individual label.
- Pros and Cons
  - task-related, good performance (deep neural networks)
  - high cost of data labeling.

# Weakly supervised learning

- Concepts
  - learning with weak supervision.
  - noisy, limited, or imprecise sources
- Three types of weak supervision
  - *incomplete*: speaker / image categorization
  - *inexact*: object in a video / image / doc.
  - *inaccurate*: crowdsourcing

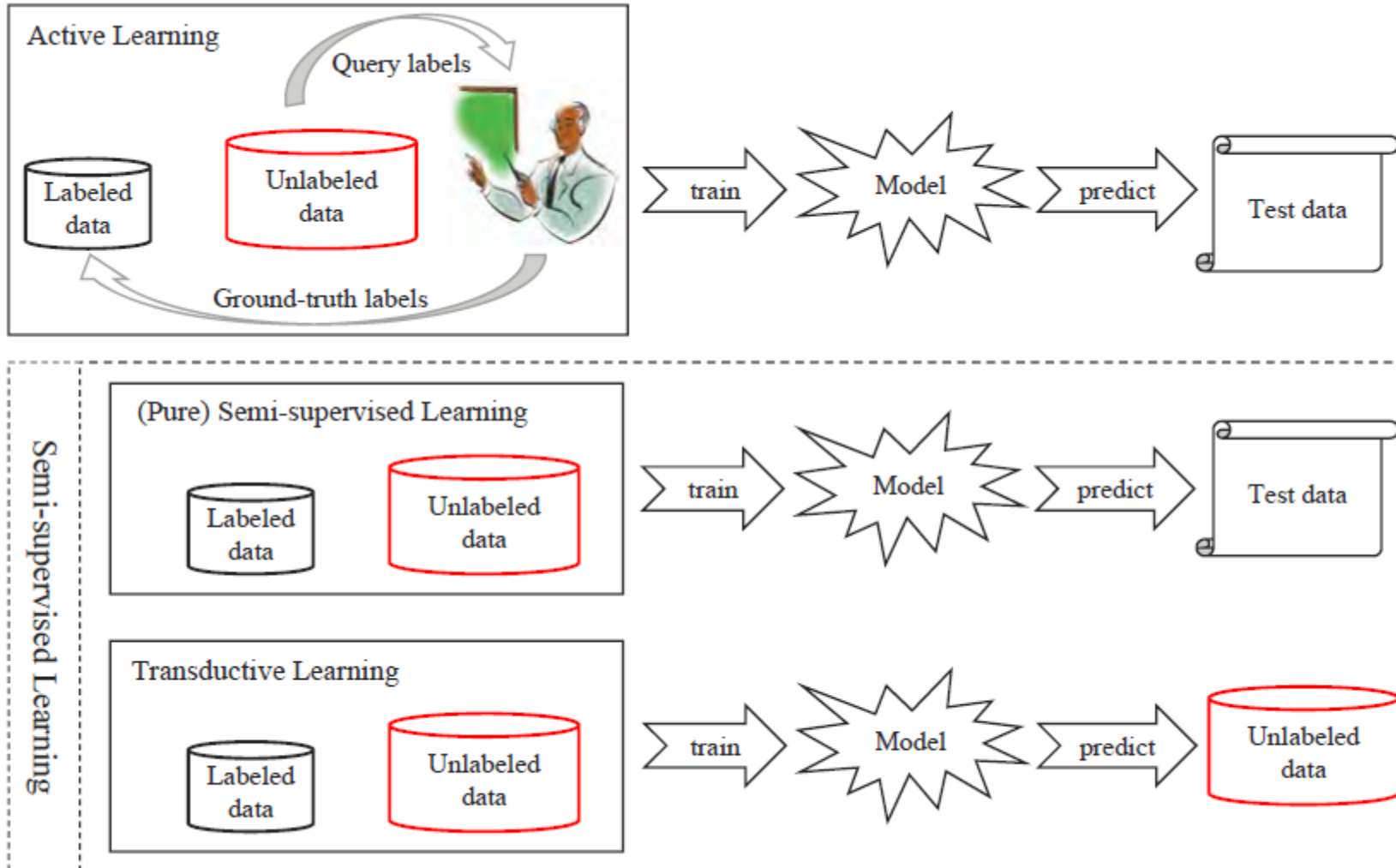
# Weak supervision



# Incomplete Supervision

- Active learning
  - with human intervention
  - labels can be queried from an oracle.
- Semi-supervised learning
  - without human intervention
  - automatically exploit unlabeled data to improve performance

# Incomplete Supervision



# Active learning

- Goal
  - minimize the number of queries to minimize labeling cost
  - to select valuable unlabeled data
- Selection criteria
  - *informativeness*: uncertainty and entropy (0.55 vs. 0.99; 4:3 vs 6:1)
  - *representativeness*: sampling distribution (clusters)



# Semi-supervised learning

- Goal
  - data without labels to help construct models

# Explanation by GMM

$$f(\mathbf{x}|\Theta) = \sum_{j=1}^n \alpha_j f(\mathbf{x}|\theta_j), \quad (1)$$

where  $\alpha_i$  is the mixture coefficient,  $\sum_{i=1}^n \alpha_i = 1$ , and  $\Theta = \{\theta_i\}$  are the model parameters. In this case, label  $y_i$  can be considered as a random variable whose distribution  $P(y_i|\mathbf{x}_i, g_i)$  is determined by the mixture component  $g_i$  and the feature vector  $\mathbf{x}_i$ . According to the **maximum a posteriori** criterion, we have the model

$$h(\mathbf{x}) = \arg \max_{c \in \{Y, N\}} \sum_{j=1}^n P(y_i = c | g_i = j, \mathbf{x}_i) P(g_i = j | \mathbf{x}_i), \quad (2)$$

where

$$P(g_i = j | \mathbf{x}_i) = \frac{\alpha_j f(\mathbf{x}_i | \theta_j)}{\sum_{k=1}^n \alpha_k f(\mathbf{x}_i | \theta_k)}.$$

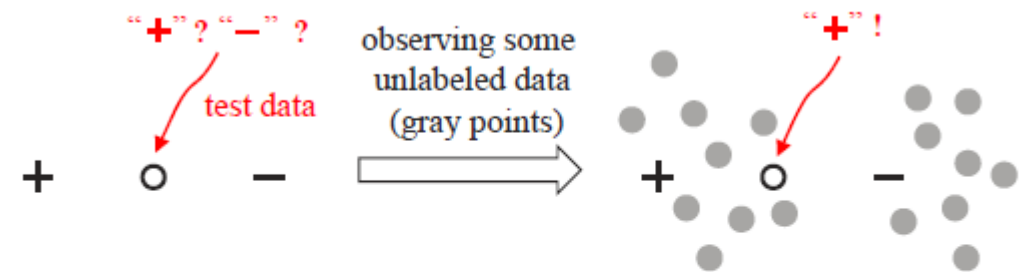


Figure 3. Illustration of the usefulness of unlabeled data

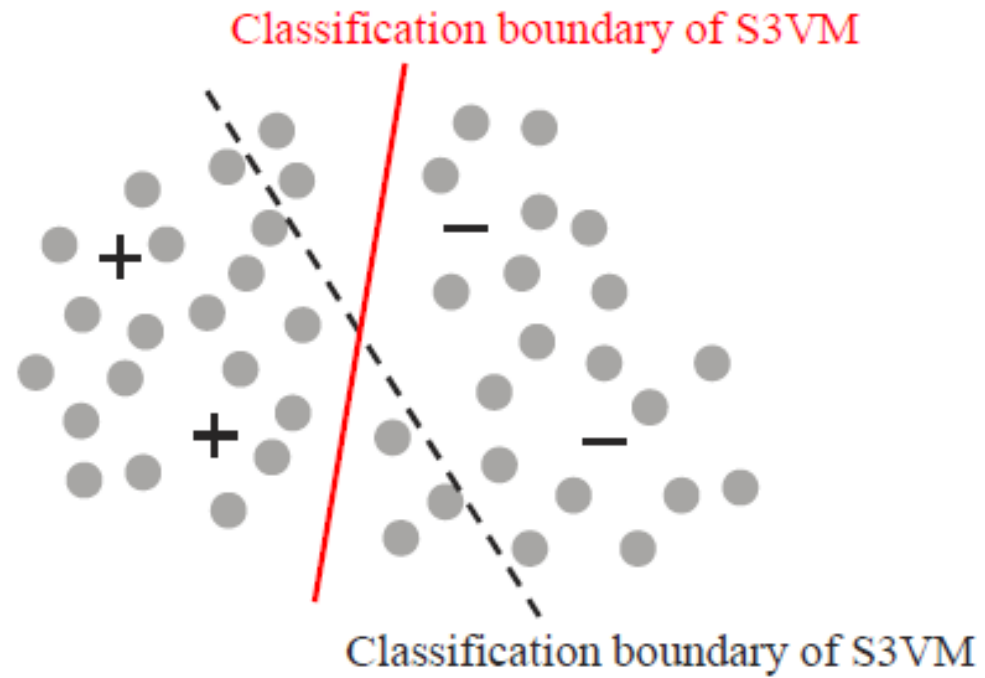
# Semi-supervised learning

- Goal
  - data without labels to help construct models
- Data assumptions
  - cluster assumption (the same cluster has the same class)
  - manifold assumption (nearby instances have similar predictions)

# Semi-supervised learning

- Categories
  - generative methods (GMMs)
  - graph-based methods (knowledge graph: relation completion)
  - low-density separation methods (S3VMs)

# S3VMs vs. SVM



**Figure 4. Illustration of the usefulness of unlabeled data**

# Semi-supervised learning

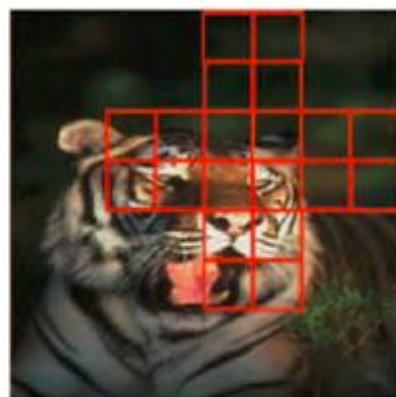
- Categories
  - generative methods (GMMs)
  - graph-based methods (knowledge graph: relation completion)
  - low-density separation methods (S3VMs)
  - disagreement-based methods (co-training)

# Inexact Supervision

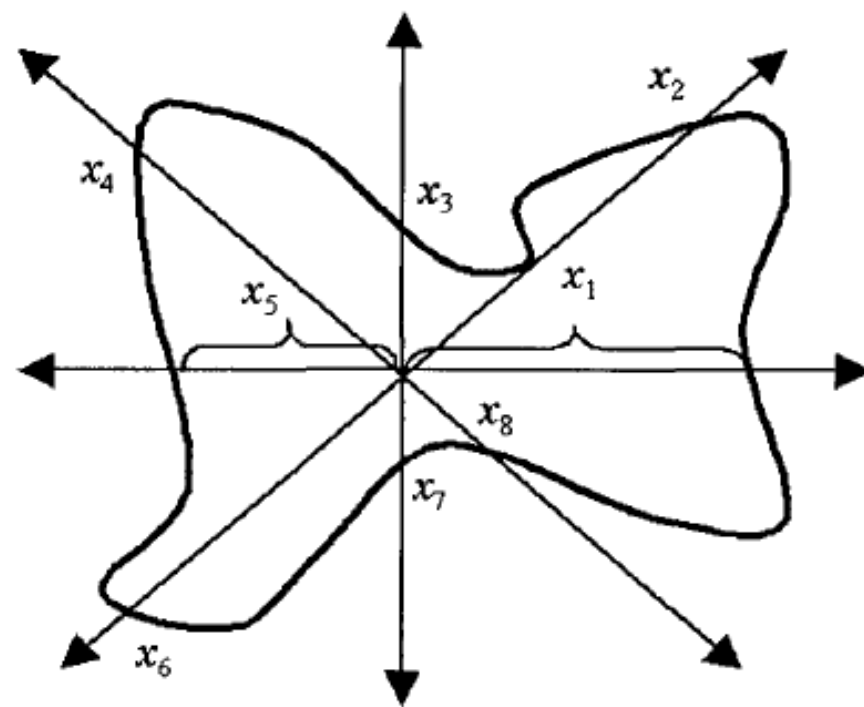
- Multi-instance learning
  - object in a video / image / doc.
  - Bag generates instances based on concept.



(a) SB



(b) SBN



# Inaccurate Supervision

- Learning with label noise
  - add error rate in the cost function
  - data editing

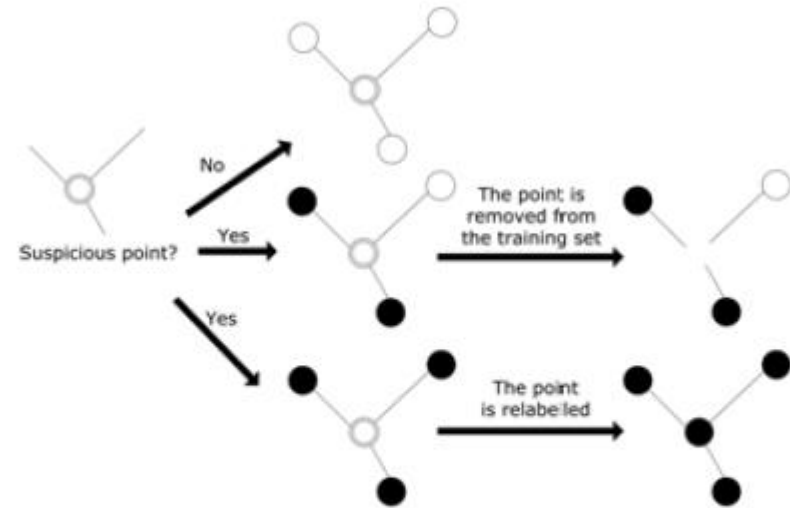


Figure 6. Identify and remove/relabel suspicious points



# Inaccurate Supervision

- Learning with label noise
  - add error rate in the cost function
  - data editing
  - crowdsourcing
    - ensemble methods with voting
    - spammer elimination
    - combine with economics (Nash equilibrium)

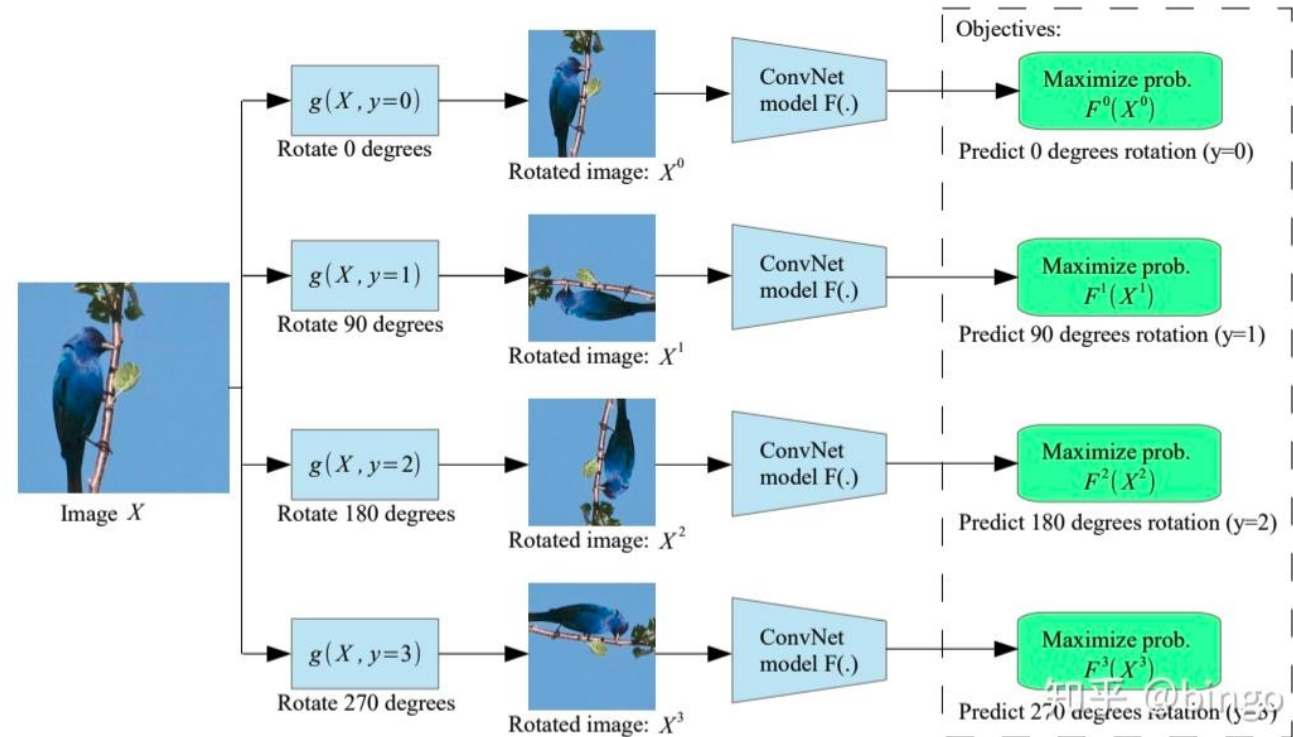
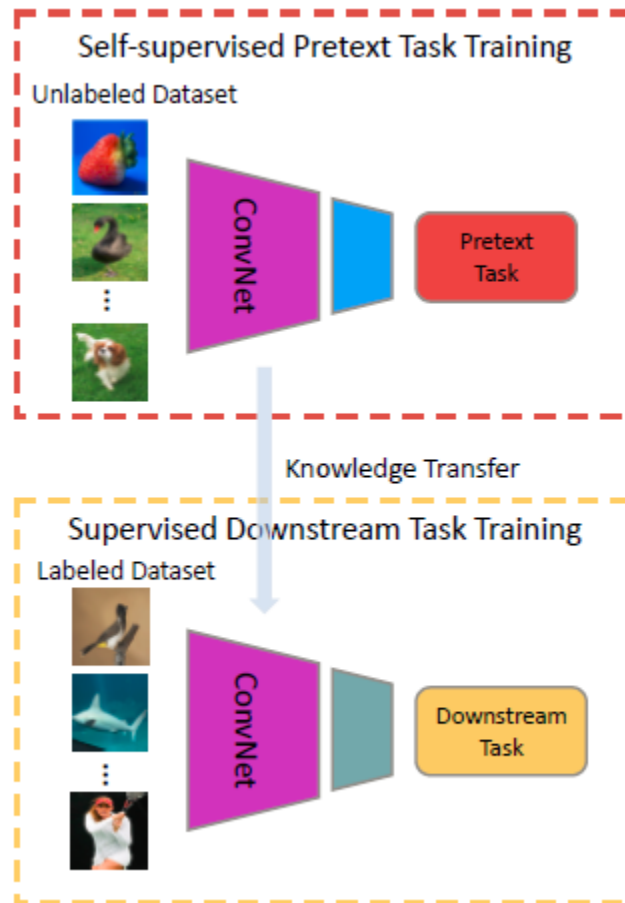
# Quick summary

- Weakly supervised learning
- Weak supervision
  - Semi-supervised learning

# Self-Supervised Learning

# Self-supervised learning

- Self-supervised learning is supervised learning without human-annotated labels.

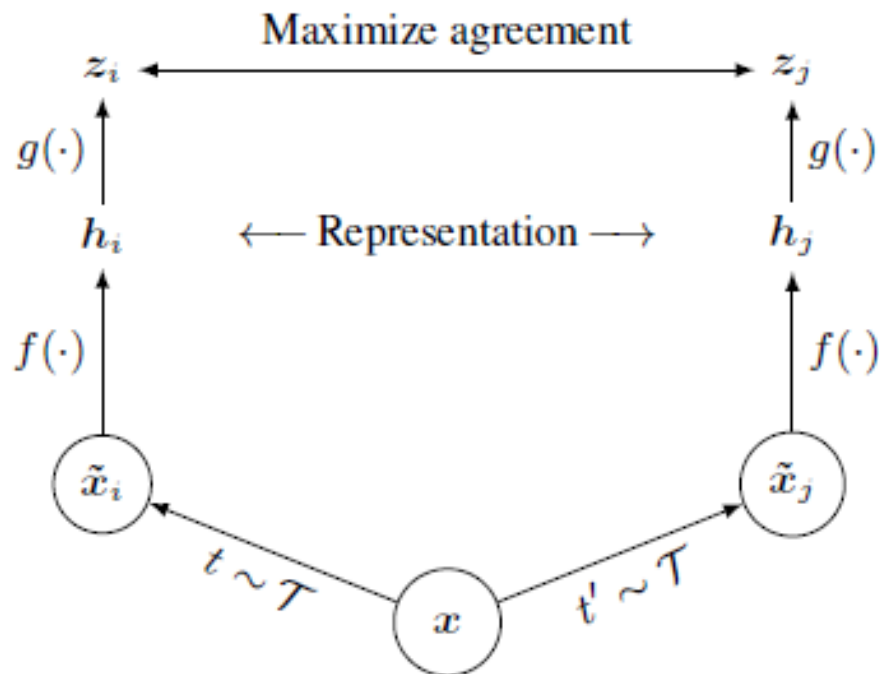


# Learning representation without supervision

- Generative model
  - AE, VAE, PixRNNs ...
- Discriminative model
  - objective function is the same as supervised learning.
  - perform on pretext task.
  - but inputs and labels are derived from an unlabeled dataset.

# SimCLR

- data augmentation operations
- projection head
- contrastive loss function




---

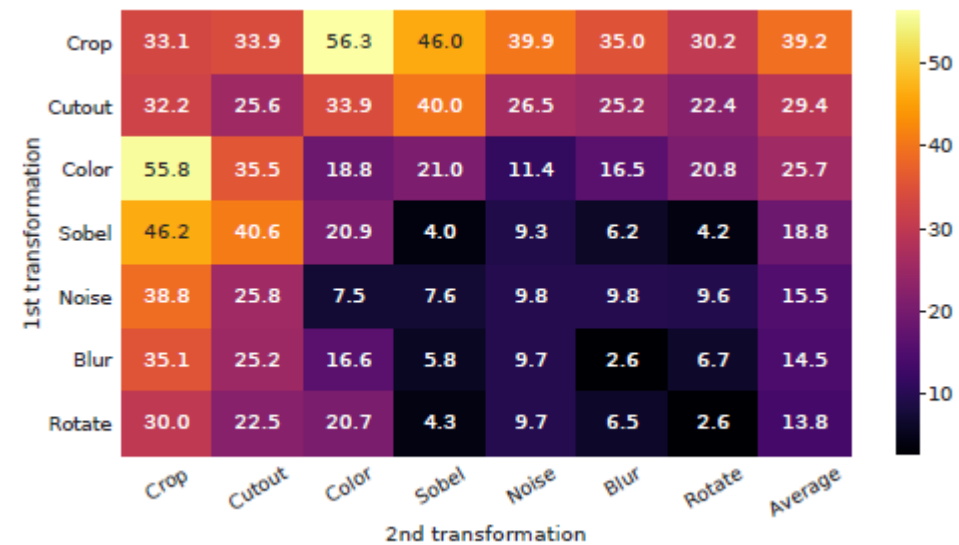
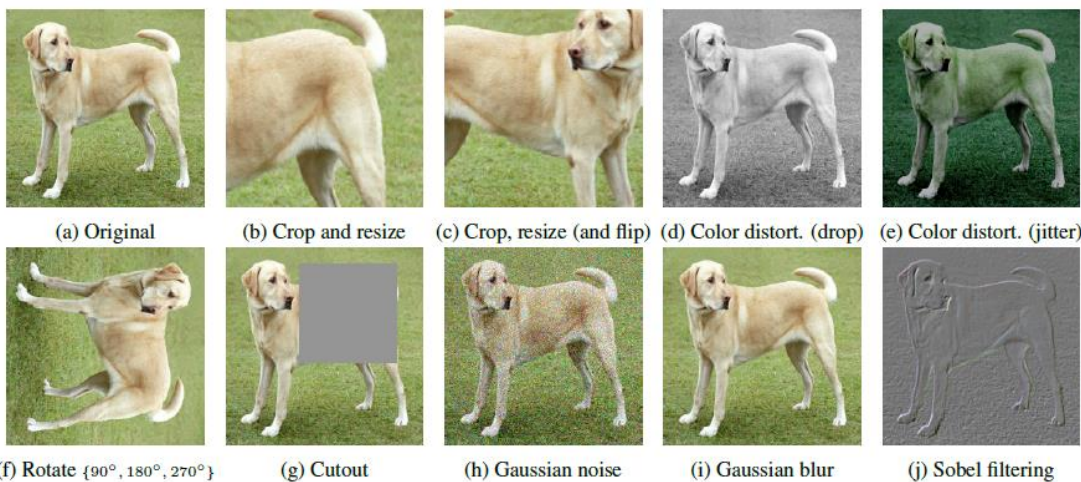
**Algorithm 1** SimCLR's main learning algorithm.

---

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

---

# Data augmentation



Methods	Color distortion strength					AutoAug
	1/8	1/4	1/2	1	1 (+Blur)	
SimCLR	59.6	61.0	62.6	63.2	64.5	61.1
Supervised	77.0	76.7	76.5	75.7	75.4	77.1

# Projection head

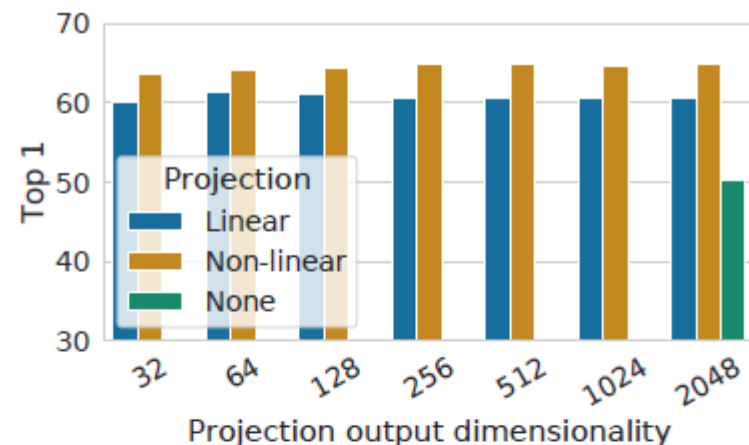
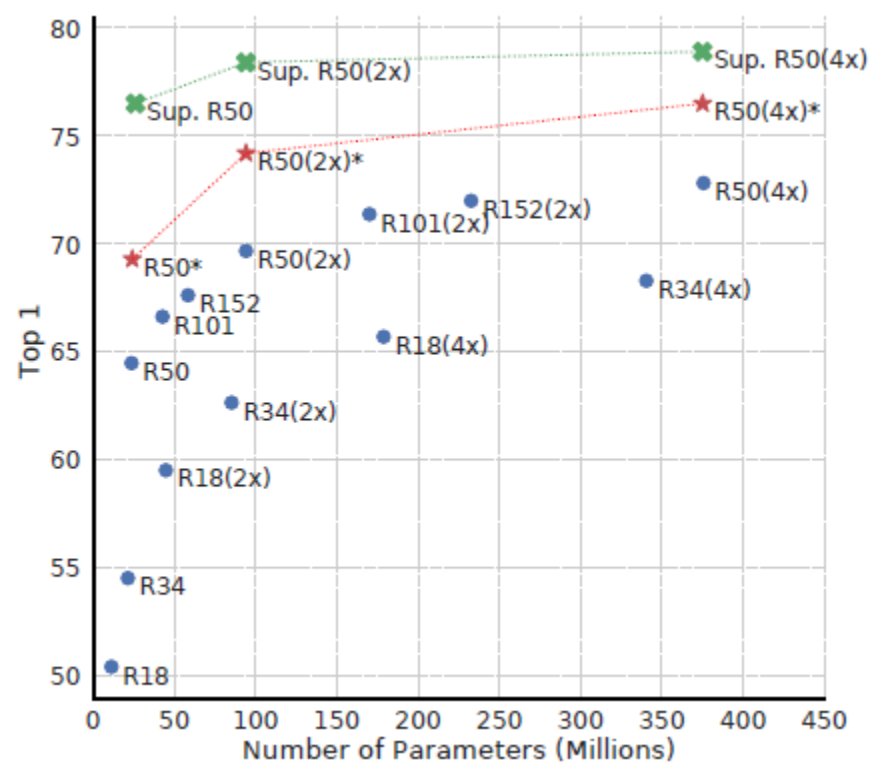


Figure 8. Linear evaluation of representations with different projection heads  $g(\cdot)$  and various dimensions of  $z = g(h)$ . The representation  $h$  (before projection) is 2048-dimensional here.

What to predict?	Random guess	Representation $h$	Representation $g(h)$
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  $\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}$ ), 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.



# Contrastive loss

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

Table 2. Negative loss functions and their gradients. All input vectors, i.e.  $u, v^+, v^-$ , are  $\ell_2$  normalized. NT-Xent is an abbreviation for “Normalized Temperature-scaled Cross Entropy”. Different loss functions impose different weightings of positive and negative examples.

NT-Xent -> 1:1 + 1:2N-1

NT-Logistic -> 1:1

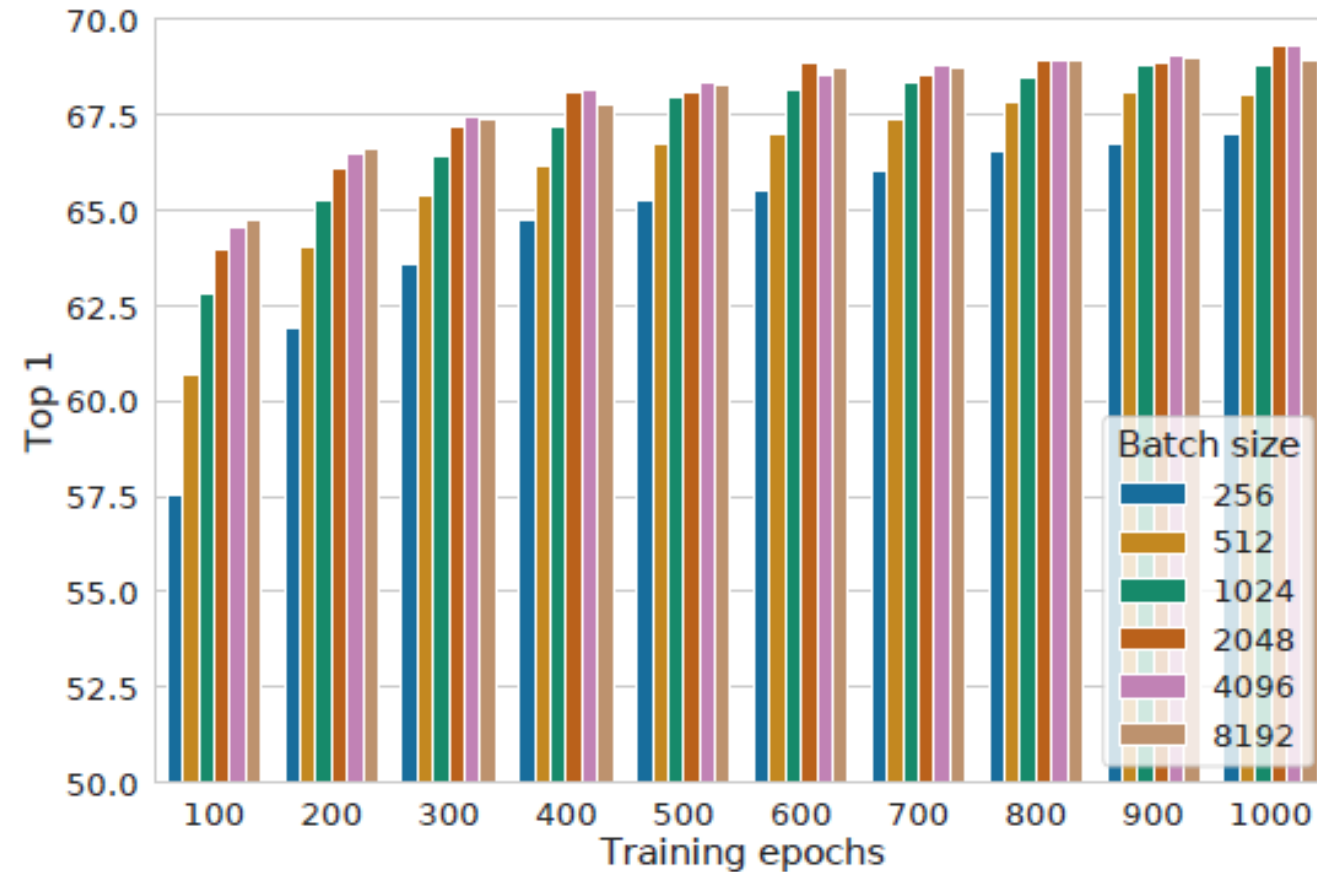
Margin Triple -> 1:1 + 1:1

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.

# Contrastive loss

Larger batch sizes and longer training compared with supervised learning



# Self-supervised methods

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 methods

Method	Architecture	Label fraction	
		1%	10%
<i>Methods using other label-propagation:</i>			
Pseudo-label	ResNet50	51.6	82.4
VAT+Entropy Min.	ResNet50	47.0	83.4
UDA (w. RandAug)	ResNet50	-	88.5
FixMatch (w. RandAug)	ResNet50	-	89.1
S4L (Rot+VAT+En. M.)	ResNet50 (4×)	-	91.2
<i>Methods using representation learning only:</i>			
InstDisc	ResNet50	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.

# Transfer learning

	Food	CIFAR10	CIFAR100	Birdsnap	SUN397	Cars	Aircraft	VOC2007	DTD	Pets	Caltech-101	Flowers
<i>Linear evaluation:</i>												
SimCLR (ours)	<b>76.9</b>	<b>95.3</b>	80.2	48.4	<b>65.9</b>	60.0	61.2	<b>84.2</b>	<b>78.9</b>	89.2	<b>93.9</b>	<b>95.0</b>
Supervised	75.2	<b>95.7</b>	<b>81.2</b>	<b>56.4</b>	64.9	<b>68.8</b>	<b>63.8</b>	83.8	<b>78.7</b>	<b>92.3</b>	<b>94.1</b>	94.2
<i>Fine-tuned:</i>												
SimCLR (ours)	<b>89.4</b>	<b>98.6</b>	<b>89.0</b>	<b>78.2</b>	<b>68.1</b>	<b>92.1</b>	<b>87.0</b>	<b>86.6</b>	<b>77.8</b>	92.1	<b>94.1</b>	97.6
Supervised	88.7	98.3	<b>88.7</b>	<b>77.8</b>	67.0	91.4	<b>88.0</b>	86.5	<b>78.8</b>	<b>93.2</b>	<b>94.2</b>	<b>98.0</b>
Random init	88.3	96.0	81.9	<b>77.0</b>	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×) models pretrained on ImageNet. Results not significantly worse than the best ( $p > 0.05$ , permutation test) are shown in bold. See Appendix B.6 for experimental details and results with standard ResNet-50.

# Quick summary

- Speaker representation