

Can contrastive learning avoid shortcut solutions?

Joshua Robinson[†]

Li Sun^{*}

Ke Yu^{*}

Kayhan Batmanghelich^{*}

Stefanie Jegelka[†]

Suvrit Sra[†]

[†]*Massachusetts Institute of Technology*

^{*}*University of Pittsburgh*

joshrob@mit.edu

lis118@pitt.edu

yu.ke@pitt.edu

kayhan@pitt.edu

stefje@csail.mit.edu

suvrit@mit.edu

arXiv:2106.11230v1 [cs.LG] 21 Jun 2021

<https://github.com/joshr17/IFM>

Problem

- **Generalization** of representations learned via contrastive learning depends on what features of the data are extracted.
- **Contrast loss** does not always sufficiently guide which features are extracted, which may suppress important predictive features.

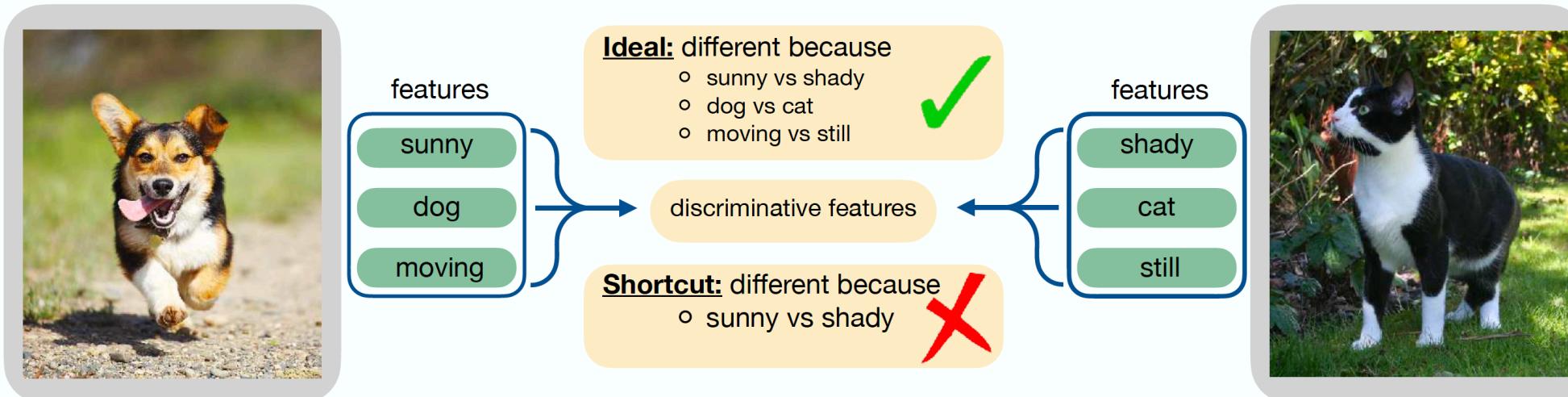


Figure 1: An ideal encoder would discriminate between instances using multiple distinguishing features instead of finding simple shortcuts that suppress features. We show that InfoNCE-trained encoders can suppress features (Sec. 2.2). However, making instance discrimination harder during training can trade off representation of different features (Sec. 2.3). To avoid the need for trade-offs we propose *implicit feature modification* (Sec. 3), which reduces suppression in general, and improves generalization (Sec. 4).

Definition of feature suppression

Formally, we assume that the data has underlying feature spaces $\mathcal{Z}^1, \dots, \mathcal{Z}^n$ with a distribution p_j on each \mathcal{Z}^j . Each $j \in [n]$, corresponding to a latent space \mathcal{Z}^j , models a distinct feature. We write the product as $\mathcal{Z}^S = \prod_{j \in S} \mathcal{Z}^j$, and simply write \mathcal{Z} instead of $\mathcal{Z}^{[n]}$ where $[n] = \{1, \dots, n\}$. A set of features $z = (z^j)_{j \in [n]} \in \mathcal{Z}$ is generated by sampling each coordinate $z^j \in \mathcal{Z}^j$ independently, and we denote the measure on \mathcal{Z} induced by z by $\lambda(\cdot | z^S)$. Further, let $\lambda(\cdot | z^S)$ denote the conditional measure on \mathcal{Z} for fixed z^S . For $S \subseteq [n]$ we use z^S to denote the projection of z onto \mathcal{Z}^S . Finally, an injective map $g : \mathcal{Z} \rightarrow \mathcal{X}$ produces observations $x = g(z)$.

Our aim is to train an encoder $f : \mathcal{X} \rightarrow \mathbb{S}^{d-1}$ to map input data x to the surface of the unit sphere $\mathbb{S}^{d-1} = \{u \in \mathbb{R}^d : \|u\|_2 = 1\}$ in such a way that f extracts useful information. To formally define feature suppression, we need the *pushforward* $h\#\nu(V) = \nu(h^{-1}(V))$ of a measure ν on a space \mathcal{U} for a measurable map $h : \mathcal{U} \rightarrow \mathcal{V}$ and measurable $V \subseteq \mathcal{V}$, where $h^{-1}(V)$ denotes the preimage.

Definition 1. Consider an encoder $f : \mathcal{X} \rightarrow \mathbb{S}^{d-1}$ and features $S \subseteq [n]$. For each $z^S \in \mathcal{Z}^S$, let $\mu(\cdot | z^S) = (f \circ g)\#\lambda(\cdot | z^S)$ be the pushforward measure on \mathbb{S}^{d-1} by $f \circ g$ of the conditional $\lambda(\cdot | z^S)$.

1. f suppresses S if for any pair $z^S, \bar{z}^S \in \mathcal{Z}^S$, we have $\mu(\cdot | z^S) = \mu(\cdot | \bar{z}^S)$.
2. f distinguishes S if for any pair of distinct $z^S, \bar{z}^S \in \mathcal{Z}^S$, measures $\mu(\cdot | z^S), \mu(\cdot | \bar{z}^S)$ have disjoint support.

Why optimizing the InfoNCE loss can still lead to feature suppression

$$\mathcal{L}_m(f) = \mathbb{E}_{x, x^+, \{x_i^-\}_{i=1}^m} \left[-\log \frac{e^{f(x)^\top f(x^+)/\tau}}{e^{f(x)^\top f(x^+)/\tau} + \sum_{i=1}^m e^{f(x)^\top f(x_i^-)/\tau}} \right], \quad (1)$$

loss that do and solutions that do not suppress a given feature. Following previous work [40, 49, 56], we analyze the loss as the number of negatives goes to infinity,

$$\mathcal{L} = \lim_{m \rightarrow \infty} \left\{ \mathcal{L}_m(f) - \log m - \frac{2}{\tau} \right\} = \frac{1}{2\tau} \mathbb{E}_{x, x^+} \|f(x) - f(x^+)\|^2 + \mathbb{E}_{x^+} \log \left[\mathbb{E}_{x^-} e^{f(x^+)^\top f(x^-)/\tau} \right].$$

We subtract $\log m$ to ensure the limit is finite, and use x^- to denote a random sample with the same distribution as x_i^- . Prop. 1 (proved in App. A) shows that, assuming the marginals p_j are uniform, the InfoNCE loss is optimized both by encoders that suppress feature j , and by encoders that distinguish j .

Proposition 1. *Suppose that p_j is uniform on $\mathcal{Z}^j = \mathbb{S}^{d-1}$ for all $j \in [n]$. Then for any feature $j \in [n]$ there exists an encoder f_{supp} that suppresses feature j and encoder f_{disc} that discriminates j but both attain $\min_{f: \text{measurable}} \mathcal{L}(f)$.*

Controlling feature learning via the difficulty of instance discrimination

$$\mathcal{L}_m(f) = \mathbb{E}_{x, x^+, \{x_i^-\}_{i=1}^m} \left[-\log \frac{e^{f(x)^\top f(x^+)/\tau}}{e^{f(x)^\top f(x^+)/\tau} + \sum_{i=1}^m e^{f(x)^\top f(x_i^-)/\tau}} \right], \quad (1)$$

1. Temperature τ in the InfoNCE loss (Eqn. 1). Smaller τ places higher importance on positive and negative pairs with high similarity [47].
 2. Hard negative sampling method of Robinson et al. [40], which uses importance sampling to sample harder negatives. The method introduces a hardness concentration parameter β , with larger β corresponding to harder negatives (see [40] for full details).
- 40.** Joshua Robinson, Ching-Yao Chuang, Suvrit Sra, and Stefanie Jegelka. Contrastive learning with hard negative samples. In *Int. Conf. on Learning Representations (ICLR)*, 2021.

Proposition 2 (Informal). Suppose that p_j is uniform on $\mathcal{Z}^j = \mathbb{S}^{d-1}$ for all $j \in [n]$. Further, for $S \subseteq [n]$ suppose that $x, x^+, \{x_i^-\}_i$ are conditioned on the event that they have the same features S . Then any f that minimizes the (limiting) InfoNCE loss suppresses features S .

Implicit feature modification for reducing feature suppression

To avoid this effect, we develop a technique that *adaptively* modifies samples to **remove** whichever features are used to discriminate a particular positive pair from negatives, then trains an encoder to discriminate instances using *both* the original features, and the features left over after modification. While a natural method for modifying features is to directly transform raw input data, it is very **challenging** to modify the semantics of an input in this way. So instead we propose **modifying features** by applying transformations to encoded samples $v = f(x)$. Since we modify the encoded samples, instead of raw inputs x , we describe our method as *implicit*.

We set up our notation. Given batch $x, x^+, \{x_i^-\}_{i=1}^m$ we write $v = f(x)$, $v^+ = f(x^+)$, and $v_i^- = f(x_i^-)$ to denote the corresponding embeddings. As in Eqn. 1, the point-wise InfoNCE loss is,

$$\ell(v, v^+, \{v_i^-\}_{i=1}^m) = -\log \frac{e^{v^\top v^+ / \tau}}{e^{v^\top v^+ / \tau} + \sum_{i=1}^m e^{v^\top v_i^- / \tau}}.$$

Definition 2 (Implicit feature modification). Given budget $\epsilon \in \mathbb{R}_+^m$, and encoder $f : \mathcal{X} \rightarrow \mathbb{S}^d$, an adversary removes features from f that discriminates batch $x, x^+, \{x_i^-\}_{i=1}^m$ by maximizing the point-wise InfoNCE loss, $\ell_\epsilon(v, v^+, \{v_i^-\}_{i=1}^m) = \max_{\delta^+ \in \mathcal{B}_{\epsilon^+}, \{\delta_i^- \in \mathcal{B}_{\epsilon_i}\}_{i=1}^m} \ell(v, v^+ + \delta^+, \{v_i^- + \delta_i^-\}_{i=1}^m)$.

Implicit feature modification for reducing feature suppression

Lemma 1. For any $v, v^+, \{v_i^-\}_{i=1}^m \in \mathbb{R}^d$ we have,

$$\nabla_{v_j^-} \ell = \frac{e^{v^\top v_j^- / \tau}}{e^{v^\top v^+ / \tau} + \sum_{i=1}^m e^{v^\top v_i^- / \tau}} \cdot \frac{v}{\tau} \quad \text{and} \quad \nabla_{v^+} \ell = \left(\frac{e^{v^\top v^+ / \tau}}{e^{v^\top v^+ / \tau} + \sum_{i=1}^m e^{v^\top v_i^- / \tau}} - 1 \right) \cdot \frac{v}{\tau}.$$

In particular, $\nabla_{v_j^-} \ell \propto v$ and $\nabla_{v^+} \ell \propto -v$.

This expression shows that the adversary perturbs v_j^- (resp. v^+) in the direction of the anchor v (resp $-v$). Since the derivative directions are *independent* of $\{v_i^-\}_{i=1}^m$ and v^+ , we can analytically compute optimal perturbations in \mathcal{B}_ϵ . Indeed, following the constant ascent direction shows the optimal updates are simply $v_i^- \leftarrow v_i^- + \epsilon_i v$ and $v^+ \leftarrow v^+ - \epsilon^+ v$. The positive (resp. negative) perturbations increase (resp. decrease) cosine similarity to the anchor $\text{sim}(v, v_i^- + \epsilon_i v) \rightarrow 1$ as $\epsilon_i \rightarrow \infty$ (resp. $\text{sim}(v, v^+ - \epsilon^+ v) \rightarrow -1$ as $\epsilon^+ \rightarrow \infty$). In Fig. 4 we visualize the newly synthesized v_i^-, v^+ and find meaningful interpolation of semantics. Plugging the update rules for v^+ and v_i^- into the point-wise InfoNCE loss yields,

$$\ell_\epsilon(v, v^+, \{v_i^-\}_{i=1}^m) = -\log \frac{e^{(v^\top v^+ - \epsilon^+)/\tau}}{e^{(v^\top v^+ - \epsilon^+)/\tau} + \sum_{i=1}^m e^{(v^\top v_i^- + \epsilon_i)/\tau}}. \quad (2)$$

$$\min_f \{\mathcal{L}(f) + \alpha \mathcal{L}_\epsilon(f)\}/2$$

Implicit feature modification for reducing feature suppression

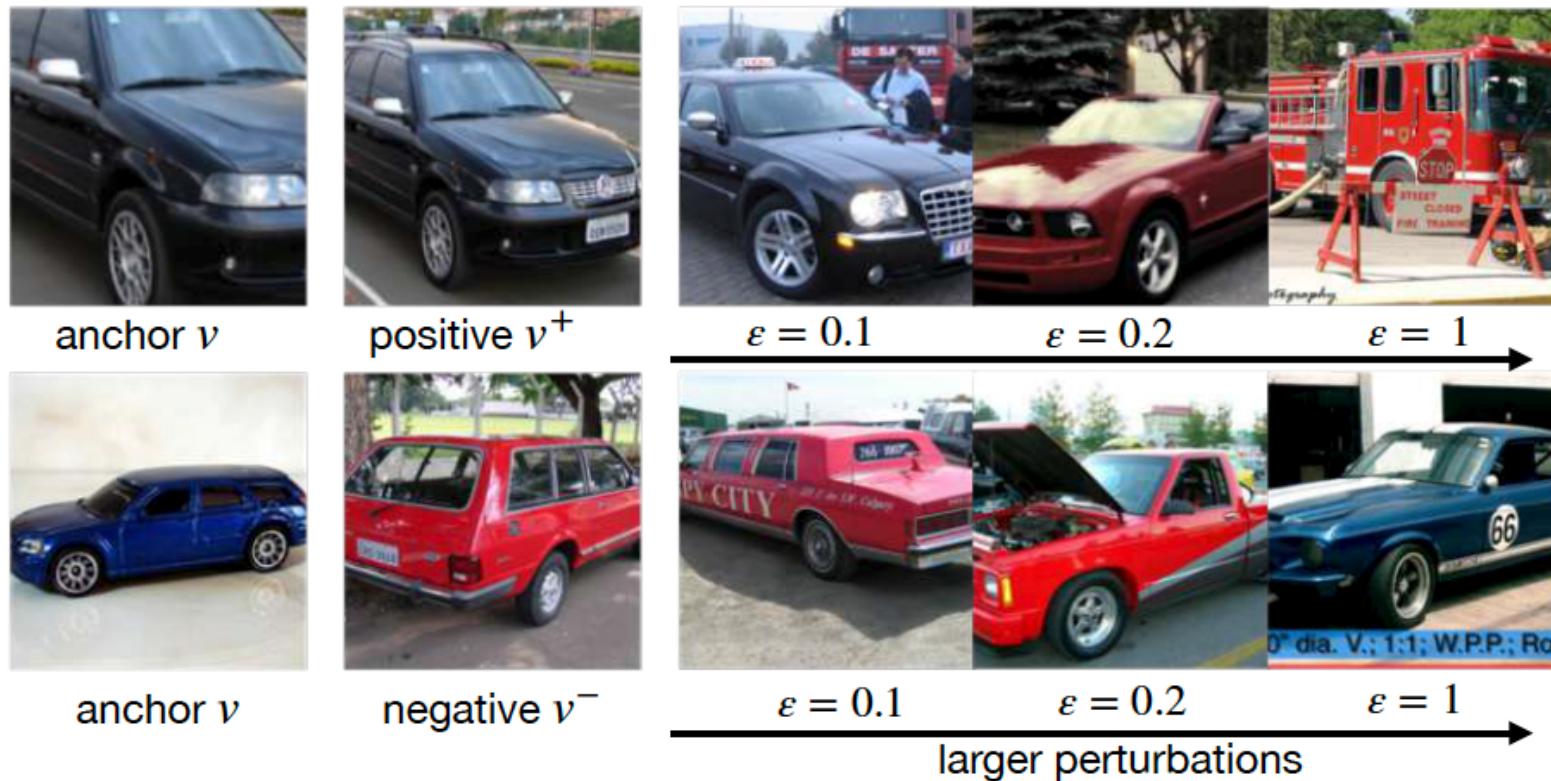


Figure 4: Visualizing implicit feature modification. **Top row:** progressively moving positive sample away from anchor. **Bottom row:** progressively moving negative sample away from anchor. In both cases, semantics such as color, orientation, and vehicle type are modified, showing the suitability of implicit feature modification for altering instance discrimination tasks.

Experimental Results

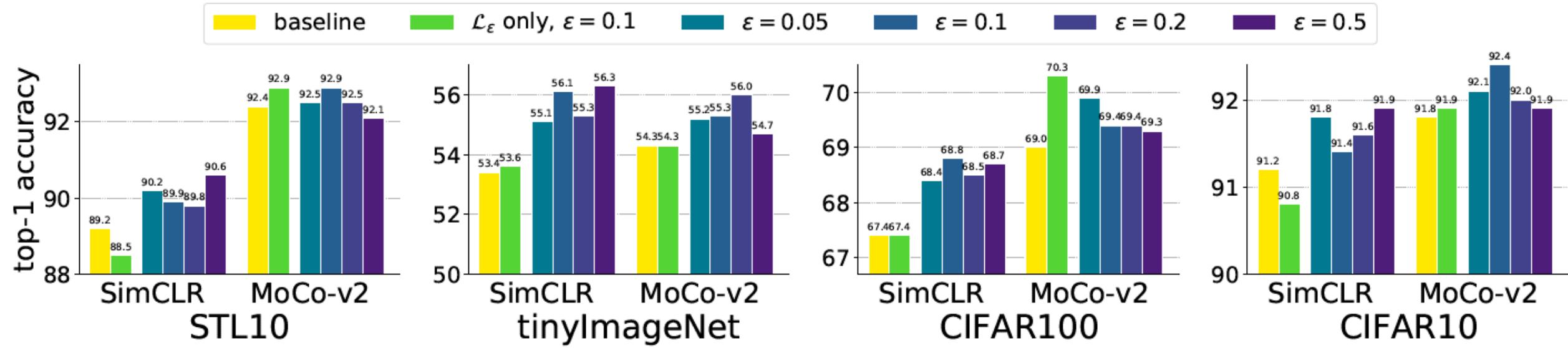


Figure 7: IFM improves linear readout performance on all datasets for all $\epsilon \in \{0.05, 0.1, 0.2\}$ compared to SimCLR and MoCo-v2 baselines. Protocol uses 400 epochs of training with ResNet-50 backbone.

Experimental Results

Method	logFEV1pp	logFEV ₁ FVC	CLE	CLE 1-off	Para-septal	Para-septal 1-off	mMRC	mMRC 1-off
Loss	R-Square		Accuracy (%)					
\mathcal{L} (baseline)	$0.566 \pm .005$	$0.661 \pm .005$	49.6 ± 0.4	81.8 ± 0.5	55.7 ± 0.3	84.4 ± 0.2	50.4 ± 0.5	72.5 ± 0.3
$\mathcal{L}_\epsilon, \epsilon = 0.1$	$0.591 \pm .008$	$0.681 \pm .008$	49.4 ± 0.4	81.9 ± 0.3	55.6 ± 0.3	85.1 ± 0.2	50.3 ± 0.8	72.7 ± 0.4
IFM, $\epsilon = 0.1$	$0.615 \pm .005$	$0.691 \pm .006$	48.2 ± 0.8	80.6 ± 0.4	55.3 ± 0.4	84.7 ± 0.3	50.4 ± 0.5	72.8 ± 0.2
IFM, $\epsilon = 0.2$	$0.595 \pm .006$	$0.683 \pm .006$	48.5 ± 0.6	80.5 ± 0.6	55.3 ± 0.3	85.1 ± 0.1	49.8 ± 0.8	72.0 ± 0.3
IFM, $\epsilon = 0.5$	$0.607 \pm .006$	$0.683 \pm .005$	49.6 ± 0.4	82.0 ± 0.3	54.9 ± 0.2	84.7 ± 0.2	50.6 ± 0.4	73.1 ± 0.2
IFM, $\epsilon = 1.0$	$0.583 \pm .005$	$0.675 \pm .006$	50.0 ± 0.5	82.9 ± 0.4	56.3 ± 0.6	85.7 ± 0.2	50.3 ± 0.6	71.9 ± 0.3

Table 2: Linear readout performance on COPDGene dataset. The values are the average of 5-fold cross validation with standard deviations. IFM yields improvements on all phenotype predictions.

Experimental Results

Constructing non-robust features. Given encoder f we finetune a linear probe (classifier) h on-top of f using training data (to avoid smoothing effects we do not use data augmentation). Once h is trained, we consider each labeled example (x, y) from training data $\mathcal{D}_{\text{train}} \in \{\text{tinyImageNet}, \text{STL10}, \text{CIFAR10}, \text{CIFAR100}\}$. A hallucinated target label t is sampled uniformly at random, and we perturb $x = x_0$ until $h \circ f$ predicts t using repeated FGSM attacks [12] $x_k \leftarrow x_{k-1} - \varepsilon \text{sign}(\nabla_x \ell(h \circ f(x_{k-1}), t))$. At each step we check if $\arg \max_i h \circ f(x_k)_i = t$ (we use the maximum of logits for inference) and stop iterating and set $x_{\text{adv}} = x_k$ for the first k for which the prediction is t . This usually takes no more than a few FGSM steps with $\varepsilon = 0.01$. We form a dataset of “robust” features by adding (x_{adv}, y) to \mathcal{D}_R , and a dataset of “non-robust” features by adding (x_{adv}, t) to \mathcal{D}_{NR} . To a human the pair (x_{adv}, t) will look mislabeled, but for the encoder x_{adv} contains features predictive of t . Finally, we re-finetune (i.e. re-train) linear classifier g using \mathcal{D}_R (resp. \mathcal{D}_{NR}) as training data.

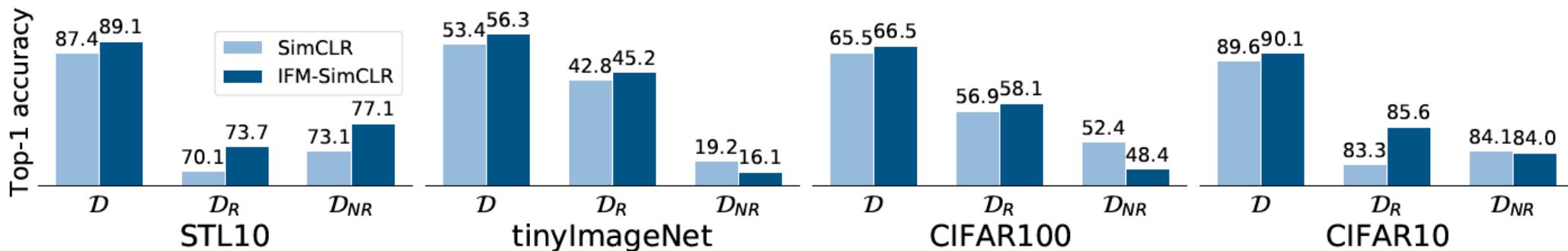


Figure 8: Label $\{\mathcal{D}, \mathcal{D}_R, \mathcal{D}_{NR}\}$ indicates which dataset was used to train the linear readout function. Improved performance of IFM on standard data \mathcal{D} can be attributed to improved representation of *robust* features \mathcal{D}_R . See Sec. 4.3 for construction of robust (\mathcal{D}_R) and non-robust (\mathcal{D}_{NR}) feature datasets.

THANK YOU!

? Q & A !