# Unsupervised Visual Representation Learning via Mutual Information Regularized Assignment

Anonymous Author(s) Affiliation Address email

## Abstract

This paper proposes Mutual Information Regularized Assignment (MIRA), a 1 pseudo-labeling algorithm for unsupervised representation learning inspired by 2 information maximization. We formulate online pseudo-labeling as an optimization 3 problem to find pseudo-labels that maximize the mutual information between the la-4 bel and data while being close to a given model probability. We derive a fixed-point 5 iteration method and prove its convergence to the optimal solution. aIn contrast to 6 baselines, MIRA combined with pseudo-label prediction enables a simple yet effec-7 tive clustering-based representation learning without incorporating extra training 8 techniques or artificial constraints such as sampling strategy, equipartition con-9 straints, etc. With relatively small training epochs, representation learned by MIRA 10 achieves state-of-the-art performance on various downstream tasks, including the 11 linear/k-NN evaluation and transfer learning. Especially, with only 400 epochs, our 12 method applied to ImageNet dataset with ResNet-50 architecture achieves 75.5% 13 linear evaluation accuracy. 14

# 15 **1 Introduction**

There has been a growing interest in using a large-scale dataset to build powerful machine learning models [43]. Self-supervised learning (SSL), which aims to learn a useful representation without labels, is suitable for this trend; is actively studied in the fields of natural language processing [19, 20] and computer vision [10, 29]. In the vision domain, recent SSL methods are commonly designed to use augmented views and train visual representation to be augmentation-invariant. They have achieved state-of-the-art performance surpassing supervised representation in a variety of visual tasks, including semi-supervised learning [8, 50], transfer learning [21], and object detection [13].

Meanwhile, a line of works use clustering for un-/self-supervised representation learning. They 23 explicitly assign pseudo-labels to embedded representation via clustering, and the model is thereby 24 trained to predict such labels. These clustering-based methods can account for inter-data similarity; 25 representations are encouraged to encode the semantic structure of data. Prior works [48, 46, 4, 31] 26 have shown encouraging results in small-scaled settings; Caron et al. [6] show that it can be also 27 28 applied to the large-scaled dataset or even to a non-curated dataset [7]. Recently, several works [2, 8, 29 37] have adapted the philosophy of augmentation invariance and achieved strong empirical results. They typically assign pseudo-labels using augmented views while predicting the labels looking at 30 other differently augmented views. 31

32 Despite its conceptual simplicity, a naive application of clustering to representation learning is hard 33 to achieve, especially in large-scale dataset. This is because clustering-based methods are prone to 34 collapse, i.e., all samples are assigned to a single cluster. To address this, recent methods heavily rely

<sup>35</sup> on extra training techniques or artificial constraints, such as pre-training [47], sampling strategy [6],

equipartition constraints [2, 8], etc. However, it is unclear if these additions are appropriate or how
 such components will affect the representation quality.

In this paper, we propose Mutual Information Regularized Assignment (MIRA), a pseudo-labeling 38 algorithm that enables clustering-based SSL without any artificial constraints or extra training 39 techniques. MIRA is designed to follow the infomax principle [38] and the intuition that good 40 labels are something that can reduce most of the uncertainty about the data. Our method assigns a 41 pseudo-label in a principled way by constructing an optimization problem. For a given training model 42 that predicts pseudo-labels, the optimization problem seeks a solution that maximizes the mutual 43 information (MI) between the pseudo-labels and data while taking the model probability into account. 44 We formulate the problem as a convex optimization problem and derive the necessary and sufficient 45 condition of solution with the Karush-Kuhn-Tucker (KKT) condition. The solution can be achieved 46 by fixed-point iteration that we prove the convergence. We remark that MIRA does not require any 47 form of extra training techniques or artificial constraints, e.g., equipartition constraints. 48

We apply MIRA to clustering-based representation learning and verify the representation quality on
 several standard self-supervised learning benchmarks. We demonstrate its state-of-the-art performance
 on linear/k-NN evaluation, semi-supervised learning, and transfer learning benchmark. We further
 experiment with convergence speed, scalability, and different components of our method.

- <sup>53</sup> Our contributions are summarized as follows:
- We propose MIRA, a simple and principled pseudo-label assignment strategy based on mutual information. Our method does not require extra training techniques or artificial constraints.

We apply MIRA to clustering-based representation learning and it shows comparable performance
 against the state-of-the-art methods with half of the training epochs. Especially it achieves
 75.5% top-1 accuracy on ImageNet linear evaluation with only 400 epochs of training and best
 performance in 9 out of 11 datasets in transfer learning.

• Representation by MIRA also consistently improves over other information-based SSL meth-

ods [22, 50]. Especially, our method without a multi-crop augmentation strategy achieves 73.8% top-1 accuracy and outperforms BarlowTwins [50], an information maximization-based self-

63 supervised method.

## 64 2 Related works

**Self-supervised learning** SSL methods are designed to learn the representation by solving pretext 65 tasks. Recent state-of-the-art SSL methods train their representation to be augmentation invariant. 66 They are based on various pretext tasks: instance discrimination [10, 11, 13, 14], metric learning [27, 67 12], self-training [51, 9], and clustering [2, 6, 8]; our method belongs to the clustering-based SSL 68 method. Meanwhile, these methods are prone to collapsing into a trivial solution where every 69 representation is map into a constant vector. To address this, a variety of schemes and mechanisms are 70 suggested, e.g., the asymmetric structure, redundancy reduction, etc. We will review more relevant 71 works in detail below. 72

**Collapse preventing** Many SSL approaches rely on extra training techniques and artificial assump-73 tions to prevent collapsing. In clustering-based methods, DeepCluster [6] adapts a sample strategy to 74 sample elements uniformly across pseudo-labels to deal with empty clusters; SeLa [2] and SwAV [8] 75 impose equipartition constraints to balance the cluster distribution. Similarly, SelfClassifier [1] uses a 76 uniform pseudo-label prior, and PCL [37] employs concentration scaling. DINO [9] and ReSSL [51] 77 address collapsing by specific combinations of implementation details, i.e., centering and scaling with 78 an exponential moving average network; while their mechanism for preventing collapse is unclear. 79 In this work, we show our method can naturally avoid collapsing without any of these assumptions 80 or training techniques. We achieve results better than baselines with a simple but novel information 81 regularization algorithm. We take a more detailed comparison with SeLa and SwAV in Sec. 3.3. 82

Information maximization Information maximization is a principal approach to learn representation and to avoid collapse. DeepInfoMax [30] propose the MI maximization between the local and
 global views for representation learning; the existence of negative pairs prevents training toward the
 trivial solution. BarlowTwins [50] and W-MSE [22] addresses collapsing with redundancy reduction,



Figure 1: Overview of representation learning via MIRA. In our representation learning, MIRA provides pseudo-labels with model probabilities, and the model is learned by predicting the pseudo-labels. Our main contribution is in the I pseudo-labeling process that accounts for mutual information between the pseudo-label and data. In MIRA, optimal pseudo-labels are computed through the fixed-point iteration (Eq. 6). Given such pseudo-labels, <sup>(2)</sup> model updates its parameters by gradient update on swapped prediction loss.

indirectly maximizing the information content of the embedding vectors [3]. Among clustering-based 87

approaches, IIC [33] maximizes the MI between the embedding codes for representation learning; 88

most similarly to ours, TWIST [25] proposes to combine the mutual information between the data 89

and class prediction as a negative loss term with a consistency loss. Both IIC and TWIST use the 90

MI as a loss function and directly optimize their model parameters with gradient descent of the loss. 91

However, direct optimization of MI terms by updating model parameters often leads to a sub-optimal 92 solution [25]; TWIST copes with this issue by appending the normalization layer before softmax

93 and introducing an additional self-labeling stage. In contrast, MIRA addresses the difficulty of MI

94

maximization in a principled way via explicit optimization. 95

#### 3 Method 96

In this section, we explain our pseudo-labeling algorithm-MIRA. When applying MIRA to repre-97 sentation learning, we follow the basic framework of clustering-based representation learning that 98 alternates between *pseudo-labeling*, i.e., cluster assignments, and *model training* to predict such labels. 99 Figure 1 illustrates our representation training cycle. We will first explain our main contribution, 100

MIRA (pseudo-labeling) and then explain how it applies to model training. 101

Our idea is to employ the information maximization principle into pseudo-labeling. We formulate 102 an optimization problem for online clustering that assigns soft pseudo-labels to mini-batch samples 103 (Sec. 3.1). The problem accounts for the model probability and mutual information between the 104 pseudo-labels and data. We propose an iterative method to solve the optimization problem (Sec. 3.2). 105 For the model training, we use the swapped prediction loss as in [8] (Sec. 3.3). 106

### 3.1 MI regularized cluster assignment 107

We have a model<sup>1</sup>  $f_{\theta}$  parametrized by  $\theta$  that outputs K-dimensional logit  $f_{\theta}(x) \in \mathbb{R}^{K}$  for an image 108 x, where K is a predefined number of clusters. The model probability p of an image x is then given 109 by the temperature  $\tau_t$  scaled output of the model— $p := \operatorname{softmax}(f_{\theta}(\boldsymbol{x})/\tau_t)$ —as in [8, 9]. For a mini-batch of input images  $\boldsymbol{X} = \{\boldsymbol{x}\}_{i=1}^B$ , we denote the model probability  $\boldsymbol{P} = \{\boldsymbol{p}_i\}_{i=1}^B \subset \mathbb{R}^K$ . In our pseudo-labeling, for the given model probability  $\boldsymbol{P}$ , we want to assign pseudo-labels  $\boldsymbol{W}^* = \{\boldsymbol{w}^*\}_{i=1}^B$ 110 111 112 that will be used for training the model by predicting them. 113

We argue that such pseudo-labels should maximize the mutual information between themselves and 114 data while accounting for the model probability P. Let  $\mathcal{B} \in \{1, ..., B\}$  and  $\mathcal{Y}_{W} \in \{1, ..., K\}$  be the 115 random variables associated with the data index in mini-batch and label by probability distributions 116

<sup>&</sup>lt;sup>1</sup>In our setting, the model consists of an encoder, projection head, and classification (prototype) head as in [8, 9]; the encoder output will be used as a representation.

117  $W = \{w\}_{i=1}^{B}$ , respectively. Our online pseudo-label (cluster) assignment is determined by solving 118 the following optimization problem:

$$\boldsymbol{W}^* = \arg\min_{\boldsymbol{W} \subset \Delta_K} \frac{1}{B} \sum_{i=1}^B D_{\mathrm{KL}}(\boldsymbol{w}_i, \boldsymbol{p}_i) - \beta \hat{I}(\boldsymbol{\mathcal{Y}}_{\boldsymbol{W}}; \boldsymbol{\mathcal{B}}),$$
(1)

where  $\Delta_K := \{x \in \mathbb{R}_+^K \mid x^{\mathsf{T}}\mathbf{1}_K = 1\}$ ,  $\hat{I}$  indicates an empirical (Monte Carlo) estimates of MI, and  $\beta$  is a trade-off parameter. The problem consists of the (1) KL divergence term that makes pseudo-labels to be based on the model probability p and (2) MI term between the pseudo-labels and data to induce more information into the pseudo-labels. By combining these two terms, we provide a refined pseudo-label that take account of both the model probability and MI.

To make the optimization problem tractable, we substitute the MI term  $\hat{I}$  with the mini-batch estimates of the entropy  $\hat{H}(\mathcal{Y}_{W}|\mathcal{B})$  and marginal entropy  $\hat{H}(\mathcal{Y}_{W})$  in Eq. 2. We get:

$$\hat{I}(\mathcal{Y}_{W}; \mathcal{B}) = \hat{H}(\mathcal{Y}_{W}) - \hat{H}(\mathcal{Y}_{W}|\mathcal{B}) = -\sum_{j=1}^{K} \bar{w}_{j} \log \bar{w}_{j} + \frac{1}{B} \sum_{i=1}^{B} \sum_{j=1}^{K} w_{ij} \log w_{ij},$$
(2)

$$\boldsymbol{W}^* = \underset{\boldsymbol{W} \subset \Delta_K}{\operatorname{arg\,min}} - \frac{1}{B} \sum_{i=1}^{B} \sum_{j=1}^{K} w_{ij} \log p_{ij} + \frac{1-\beta}{B} \sum_{i=1}^{B} \sum_{j=1}^{K} w_{ij} \log w_{ij} + \beta \sum_{j=1}^{k} \overline{w}_j \log \overline{w}_j, \quad (3)$$

where  $\overline{w}_j = \frac{1}{B} \sum_{i=1}^{B} w_{ij}$  is the marginal probability of a cluster j with W. In practice, we find the optimal point  $W^*$  of the optimization problem Eq. 3.

### 128 3.2 Solving strategy

To solve efficiently, we propose a fixed-point iteration that guarantees convergence to the unique optimal solution  $W^*$  of our optimization problem. The method is based on the following proposition.

- **Proposition 1.** For  $\beta \in [0, 1)$ , the problem Eq. 3 is a strictly convex optimization problem; has a
- unique optimal point  $W^*$  that satisfies the following necessary and sufficient condition.

$$\forall (i,j) \in \{1,...,B\} \times \{1,...,K\}, \quad w^*{}_{ij} = \frac{\overline{w^*}_j^{-\frac{1}{p-\beta}} p_{ij}^{\frac{1}{1-\beta}}}{\sum_{k=1}^{K} \overline{w^*}_k^{-\frac{1}{1-\beta}} p_{ik}^{\frac{1}{1-\beta}}}.$$
(4)

The proposition is driven by proving the strict convexity and then applying the Karush–Kuhn–Tucker (KKT) condition. By substituting the necessary and sufficient condition (Eq. 4) of proposition 1 into  $\overline{w}_j = \frac{1}{B} \sum_{i=1}^{B} w_{ij}$ , we get the necessary and sufficient condition with  $\overline{w^*}$ :

$$\overline{w^*}_j = \overline{w^*}_j^{-\frac{\beta}{1-\beta}} \frac{1}{B} \sum_{i=1}^B \frac{p_{ij}^{\frac{1}{1-\beta}}}{\sum_{k=1}^K \overline{w^*}_k^{-\frac{\beta}{1-\beta}} p_{ik}^{\frac{1}{1-\beta}}} \Leftrightarrow \overline{w^*}_j = \left[\frac{1}{B} \sum_{i=1}^B \frac{p_{ij}^{\frac{1}{1-\beta}}}{\sum_{k=1}^K \overline{w^*}_k^{-\frac{\beta}{1-\beta}} p_{ik}^{\frac{1}{1-\beta}}}\right]^{1-\beta}.$$
 (5)

Based on Eq. 5, we propose the following update rule for  $\{u_i^{(n)}\}_{i=1}^K \subset \mathbb{R}_+$ :

$$\forall j \in \{1, ..., K\}, \quad u_j^{(n+1)} = \left[\frac{1}{B} \sum_{i=1}^B \frac{p_{ij}^{\frac{1}{1-\beta}}}{\sum_{k=1}^K (u_k^{(n)})^{-\frac{\beta}{1-\beta}} p_{ik}^{\frac{1}{1-\beta}}}\right]^{1-\beta}, \tag{6}$$

where  $u_j^{(n)}$  converges to  $\overline{w^*}_j$  as  $n \to \infty$ . We can easily get  $w^*_{ij}$  by Eq. 4 when the marginal probability  $\overline{w^*}_j$  is given. The proof of the proposition and convergence is in the Appendix.

By using the iterative updates of Eq. 6, we get our desirable pseudo-labels. This requires a few lines of code that are simple to implement. We find that a few steps of iterations are enough for training. This is supported by the convergence analysis in Sec. 4.3. We use this fixed point iteration for pseudo-labeling and name the method–Mutual Information Regularized Assignment (MIRA) since it finds the pseudo-labels that are regularized by the mutual information.

### 144 3.3 Representation learning with MIRA

We explain how our pseudo-labeling algorithm is applied to representation learning. We integrate the computed pseudo-labels with swapped prediction loss [8]. Specifically, given the two mini-batches of

differently augmented views  $X^{(1)}, X^{(2)}$ , MIRA outputs the pseudo-labels  $U^{(1)}, U^{(2)}$  for each mini-batch independently. In parallel, model  $f_{\theta}$  provides the temperature  $\tau_s$  scaled softmax predictions 147

148

 $Q^{(1)}, Q^{(2)}$  of each mini-batch. The swapped prediction loss is given as follows: 149

$$L(\mathbf{X}^{(1)}, \mathbf{X}^{(2)}) = \ell(\mathbf{U}^{(1)}, \mathbf{Q}^{(2)}) + \ell(\mathbf{U}^{(2)}, \mathbf{Q}^{(1)})$$
  
=  $-\frac{1}{B} \sum_{i=1}^{B} \sum_{j=1}^{K} u_{ij}^{(1)} \log q_{ij}^{(2)} - \frac{1}{B} \sum_{i=1}^{B} \sum_{j=1}^{K} u_{ij}^{(2)} \log q_{ij}^{(1)}.$  (7)

This loss function (Eq. 7) is minimized with respect to the parameters  $\theta$  of the model  $f_{\theta}$  used to 150 produce the predictions  $Q^{(1)}, Q^{(2)}$ . For more detailed information about swapped prediction loss, 151 please refer to [8]. 152

In this paper, we verify our pseudo-labeling algorithm MIRA for a representation learning purpose 153 with Eq. 7. For convenience, in the rest of this paper, we call the representation learning with 154 MIRA also as MIRA. We note that MIRA can integrate recently suggested SSL components such as 155 exponential moving average (EMA) or multi-crop augmentation strategy following the baselines [14, 156 8, 9]. The pseudo-code for MIRA is provided in the Appendix. We discuss some further details as 157 follows: 158

**Preventing collapse** The MI term in Eq. 3 takes a minimum value when collapsing happens. MIRA 159 naturally avoids collapsed solution via penalizing assignment that exhibits low MI. To be more 160 specific, unless starting from the collapsed state, MIRA finds MI-maximizing points around the 161 model prediction; will not choose collapsed pseudo-labels. Hence, the iterative training to predict 162 such labels will not lead to collapsing whenever the prediction of pseudo-labels is achievable. Our 163 empirical results verify that MIRA doesn't require extra training techniques or artificial constraints to 164 address collapsing. 165

**Comparison to SwAV and SeLa** SeLa [2] and SwAV [8] assume the equipartition of data into 166 clusters. They formulate their pseudo-labeling process into optimal transport (OT) problem; solving 167 it with the iterative Sinkhorn-Knopp (SK) algorithm [16]. Mathematically, the difference to MIRA 168 is in how to deal with the marginal entropy. SeLa and SwAV constrain the marginal entropy to 169 maximum value-equipartition while MIRA decides marginal entropy by MI regularization<sup>2</sup>. Asano 170 et al. [2] argue that their pseudo-labels with OT problem maximize the information between labels 171 and data indices under the equipartition constraints. However, it more resembles assuming MI 172 maximization and finding the assignments that are OT to the model probability. In contrast, MIRA 173 directly maximizes the MI without artificial constraints. While SwAV performs better than SeLa in 174 most self-supervised benchmarks, we empirically verify that MIRA improves over SwAV in various 175 downstream tasks. 176

#### **Experiments** 4 177

In this section, we evaluate the representation quality learned via MIRA. We first provide the 178 implementation details of our representation learning with MIRA (Sec. 4.1). We present our main 179 results on linear, k-NN, semi-supervised learning, and transfer learning benchmarks in comparison to 180 other self-supervised baselines (Sec. 4.2). Finally, we conduct an analysis of MIRA (Sec. 4.3). 181

### 4.1 Implementation details 182

We mostly follow the implementation details from our baselines [8, 9, 50]. More training details 183 about evaluation procedures and analysis are described in the Appendix. 184

Architecture The training model (network) consists of an encoder, projection head, and clas-185 sification head. We use a widely used ResNet50 [28] as our base encoder and use the output of 186 average-pooled 2048d embedding as our representation for both representation training and down-187 stream evaluations. The projection head is a 3-layer fully connected MLP of sizes [2048, 2048, d]; 188 hidden layers are followed by batch normalization [32] and ReLU. The classification head is used to 189 predict the pseudo-labels; is composed of an L2-normalization layer and a weight-normalized layer 190 of the size  $d \times K$  as in [8, 9]. We use d = 256 and K = 3000. 191

<sup>&</sup>lt;sup>2</sup>Adding the equipartition constraints, our optimization problem converts to the OT problem of SwAV.

Table 1: Linear evaluation with respect to train- Table 2: Linear evaluation on ImageNet. Comparison trained on training set of ImageNet. † are results from denotes for self-labeling by [25]. Results style: best [12]. Results style: best, second best

ing epochs. All models use a ResNet-50 encoder and with other self-supervised methods on ImageNet. SL

					Method	Arch.	Epochs	Top-1	Top-5
	Epochs           Method         100         200         400         800		Supervised	R50	-	-	-		
Method			PCL [37]	R 50	200	67.6			
without multi-crop	augmei	ntations	5		SimSiam [12]	R50	800	71.3	-
SimCLR <sup>†</sup> [10]	66.5	68.3	69.8	70.4	SimCLR-v2 [11]	R50	800	71.7	-
BYOL† [27]	66.5	70.6	73.2	74.3	InfoMin [44]	R50	800	73	91.1
SimSiam <sup>†</sup> [29]	68.1	70.0	70.8	71.3	BarlowTwins [50]	R50	1000	73.2	91.0
MoCo-v3 [14]	68.9	-	-	73.8	VicReg [3]	R50	1000	73.2	91.1
$D_{aan}Cluster v2[9]$			70.2		SelfClassifier [1]	R50	800	74.1	-
DeepCluster-v2[o]	-	-	70.2	- 71.0	TWIST w/o SL [25]	R50	800	74.1	-
SWAV   [0] TWIET [25]	00.5 70.4	70.0	71.0	71.0	BYOL [27]	R50	1000	74.3	91.6
1 WISI [23]	7 <b>0.4</b>	$\frac{70.9}{72.1}$	72.0	72.0	MoCo-v3 [14]	R50	1000	74.6	-
MIKA	<u>09.4</u>	/2.1	12.9	/5.8	DeepCluster-v2 [8]	R50	800	75.2	-
with multi-crop aug	gmenta	tions			SwAV [8]	R50	800	75.3	-
DeepCluster-v2 [8]	-	-	-	75.2	DINO [8]	R50	800	75.3	-
SwAV [8]	72.1	73.9	74.6	75.3	TWIST w/ SL [25]	R50	450	75.5	-
TWIST [25]	72.9	73.7	74.4	74.1		D50	400	75.5	02.5
MIRA	73.5	74.8	75.5	-	MIKA	к30	400	/3.5	92.5

**Training details** We train our model on the training set of the ImageNet-1k ILSVRC-2012 dataset 192 [18] without using class labels. We use the same data augmentation scheme (color jittering, Gaussian 193 blur, and solarization) and multi-crop strategy (two  $224 \times 224$  and six  $96 \times 96$ ) used in [9]. We use 194 a batch size of 4096 and employ the LARS optimizer [49] with a weight decay of  $10^{-6}$ . We use 195 linearly scaled learning rate of  $lr \times$  batch size/256 [26] with a base learning rate of 0.3.<sup>3</sup> We adjust 196 the learning rate with 10 epochs of a linear warmup followed by cosine scheduling. We also use EMA 197 network by default. When the EMA is used, we set the momentum update parameter to start from 198 0.99 and increase to 1 by cosine scheduling. We use temperature scales of  $\tau_s = 0.1$ ,  $\tau_t = 0.225$  with 199 trade-off coefficient  $\beta = 2/3$ . We assign soft pseudo-labels after 30 steps of the fixed point iteration. 200 We further discuss this choice in Sec. 4.3. Otherwise stated, we use the encoder model trained by 201 MIRA with 400 epochs training and multi-crop augmentations for the evaluations in this section. 202

### 4.2 Main results 203

**Linear evaluation** Tables 1 and 2 report linear evaluation results. We follow the linear evaluation 204 settings in [27, 10]. We train a linear classifier on the top of the frozen trained backbone with the 205 labeled training set of ImageNet. We train for 100 epochs using a LARS optimizer with a batch 206 size of 1024. We use a base learning rate of 0.1 and adjust the learning rate by cosine annealing 207 schedule. We apply random-resized-crop and horizontal flip augmentations for training. We evaluate 208 the representation quality by the linear classifier's performance on the validation set of ImageNet. 209

Table 1 shows linear evaluation performance in top-1 accuracy for different un-/self-supervised 210 representation training epochs. We train and evaluate MIRA with and without multi-crop augmen-211 tations. With multi-crop augmentations, MIRA consistently outperforms baselines while achieving 212 75.5% top-1 accuracy with only 400 epochs of training. We also report that 200 epochs of training 213 with MIRA can outperform the 800 epochs results of other baselines that don't use multi-crops. 214 Without multi-crop augmentations, MIRA is comparable to MoCo-v3 [14] and performs slightly 215 worse than BYOL [27]. However, MIRA performs the best among the clustering-based [6, 8] and 216 information-driven [50, 25] methods. 217

In Table 2, we compare MIRA to other self-supervised methods with the final performance. MIRA 218 achieves the state-of-the-art performance on linear evaluation of ImageNet with only 400 epochs of 219 training. While TWIST can achieve similar performance to MIRA within 450 epochs, they require an 220 extra training stage with self-labeling; without it, they achieve 74.1% accuracy with 800 epochs of 221 training. In contrast, MIRA doesn't require additional training. 222

<sup>&</sup>lt;sup>3</sup>Otherwise stated, we also use linearly scaled learning rate for evaluation training.

Table 3: k-NN classification results on ImageNet Table 4: Semi-supervised learning results on Imawe evaluate the baselines by models of official codes. best, second best Other baseline results are from [9]. Results style: best

with respect to subsets. For 1% and 10% results, geNet. The baselines results are from [50]. Results style:

					1	%	10	)%
	Imag	eNet su	bset		Top-1	Top-5	Top-1	Top-5
Method	100%	10%	1%	Supervised	25.4	48.4	56.4	80.4
BYOL [27] SwAV [8] BarlowTwins [50] DeepCluster-v2 [8] DINO [9]	64.8 65.7 66.0 67.1 67.5	57.4 57.4 59.0 59.2 59.3	45.2 44.3 47.7 46.5 47.2	SimCLR [10] BYOL [27] SwAV [8] BarlowTwins [50]	48.3 53.2 53.9 <u>55</u>	75.5 78.4 78.5 79.2	65.6 68.8 <b>70.2</b> 69.7	87.8 89 <u>89.9</u> 89.3
MIRA	68.7	60.7	47.8	MIRA	55.5	80.3	<u>69.9</u>	90.0

Table 5: Linear evaluation results on the transfer learning datasets. Following [21], we report top-1 accuracy on Food, CIFAR-10/100, SUN397, Cars, DTD; mean-per-class accuracy on Aircraft, Pets, Caltech-101, Flowers; 11-point mAP metric on VOC2007. Results style: best

	Aircraft	Caltech101	Cars	CIFAR10	CIFAR100	DTD	Flowers	Food	Pets	SUN397	VOC2007	avg.
Supervised	43.59	90.18	44.92	91.42	73.90	72.23	89.93	69.49	91.45	60.49	83.6	73.75
InfoMin [44]	38.58	87.84	41.04	91.49	73.43	74.73	87.18	69.53	86.24	61.00	83.24	72.21
MoCo-v2 [13]	41.79	87.92	39.31	92.28	74.90	73.88	90.07	68.95	83.3	60.32	82.69	72.31
SimCLR-v2 [11]	46.38	89.63	50.37	92.53	76.78	76.38	92.9	73.08	84.72	61.47	81.57	75.07
BYOL [27]	53.87	91.46	56.4	93.26	77.86	76.91	94.5	73.01	89.1	59.99	81.14	77.05
DeepCluster-v2 [8]	54.49	91.33	58.6	94.02	79.61	78.62	94.72	77.94	89.36	65.48	83.94	78.92
SwAV [8]	54.04	90.84	54.06	93.99	79.58	77.02	94.62	76.62	87.6	65.58	83.68	77.97
MIRA	59.06	92.21	61.05	94.20	79.51	77.66	96.07	78.76	89.95	65.84	84.10	79.86

**Semi-supervised learning** In Table 4, we evaluate the trained model on the semi-supervised 223 learning benchmark of ImageNet. Following the evaluation protocol in [27, 10], we add a linear 224 classifier on top of the trained backbone and fine-tune the model with ImageNet 1% and 10% subsets. 225 We report top-1 and top-5 accuracies on the validation set of ImageNet. For the 1% subset, MIRA 226 outperforms the baselines; both the top-1 and top-5 accuracies achieve the best. For the 10% subset, 227 MIRA is comparable to SwAV [8]. 228

*k*-NN evaluation We further evaluate the quality of learned representation via the nearest neighbor 229 classifier. We follow the procedures of [9]. First, representations of the labeled training data are stored. 230 Then, the label of the new validation data is predicted with the majority vote of k-nearest stored 231 representations. We use the same evaluation settings in [9] with 20 nearest neighbors, temperature 232 scaling<sup>4</sup> of 0.07, and cosine distance metric. 233

Table 3 shows the k-NN classification accuracies on the validation set of ImageNet. We use 1/10/100%234 subsets of ImageNet training dataset to produce labeled representations. For ImageNet 1% and 10% 235 subsets, we use the same subsets of semi-supervised learning evaluation. The results show that our 236 method achieves state-of-the-art k-NN evaluation performance with ResNet50. To be more specific, 237 our method outperforms the previous state-of-the-art DINO [9] on 100% and 10% subset evaluation 238 by  $1.2 \sim 1.4\%$ . We note that BarlowTwins [50], a method also motivated by information-maximization, 239 shows a strong performance of 47.7% in the 1% subset evaluation. 240

**Transfer learning** We further evaluate the representation learned by MIRA on the transfer learning 241 benchmark following [21] that includes FGVC aircraft [39], Caltech-101 [24], Standford Cars [34], 242 CIFAR-10/100 [35], DTD [15], Oxford 102 Flowers [40], Food-101 [5], Oxford-IIIT Pets [41], 243 SUN397 [45], and Pascal VOC2007 [23] datasets. We follow the linear evaluation procedure in [21] 244 that fits a multinomial logistic regression model on the extracted representations of 2048d from the 245 trained backbone. First, we perform a hyperparameter search on the L2-normalization coefficient of 246 the logistic regression model; then the final performance is evaluated on the model that is retrained 247 on all training and validation sets with the found coefficient. 248

<sup>&</sup>lt;sup>4</sup>The temperature scaling  $\tau$  is used to calculate contributions  $\alpha_i \sim \exp(\text{distance}_i/\tau)$  and voting is weighted by the contributions of the nearest neighbors.



Figure 2: Convergence analysis of MIRA and Sinkhorn-Knopp. We observe the converging behavior of MIRA (blue) and Sinkhorn-Knopp (yellow). We experiment with trained models of MIRA (left) and SwAV (right). Since both methods are proven to converge, we iterate each method 1000 steps and regard the results as ground truth. We report the sum-squared error (SSE) with respect to the converging point in the log scale.

Table 5 shows the performance of our algorithm compared to other baselines in 11 datasets. MIRA outperforms supervised representation on 10 out of 11 datasets. Compared to the other self-supervised methods, representation learned by MIRA achieves the best performance in 9 out of 11 datasets and improves 0.9% over the second-best baseline method on average. The results confirm that the representation trained with MIRA has a strong generalization ability for classification.

### 254 4.3 Analysis

Convergence of pseudo-label assignment We study the speed of convergence of the proposed fixed-point iteration in MIRA. We also experiment with the Sinkhorn-Knopp (SK) algorithm [16] used in SwAV [8] as a baseline. We experiment with both methods on the ImageNet with a batch size of 512. We observe the converging behavior with the pre-trained models from MIRA and SwAV. Results are averaged over 1000 randomly sampled batches.

Figure 2 shows the result of the converging behavior of our method (**blue**) and SK algorithm (**yellow**) on trained models of MIRA (**left**) and SwAV (**right**). Our fixed-point iteration converges faster than the SK algorithm in both pre-trained models. Especially our default setting of 30 steps of updates are sufficient for our fixed point iteration.

Multi-crop and EMA Table 6 reports an ablation study on how EMA and multi-crop augmentation affects our representation quality. We train a model for 200 epochs in the settings with and without EMA or Multi-crop. Both EMA and Multi-crop augmentations greatly improve the linear evaluation performance as in [8, 9]. We take a further comparison with baselines that are in the same setups. With the only difference in the pseudo-labeling algorithm, our method outperforms SwAV [8] by 1.3% in top-1 accuracy. While DINO [9] also uses both multi-crop and EMA, our method outperforms DINO with fewer training epochs. The results validate the effectiveness of our pseudo-labeling algorithm. These results validate the effectiveness of our pseudo-labeling algorithm.

Table 6: Ablation study about EMA and multi-crop augmentation. We report top-1 accuracy with linear evaluation on validation set of ImageNet. The results of SwAV is from [29].

Method	Multi-Crop	EMA	Epochs	Top-1
SwAV	×	×	200	69.1
DINO	✓	✓	300	74.5
MIRA	×	×	200	70.4
	×	✓	200	72.1
	✓	✓	200	74.8

271

**Scalability** We further validate MIRA's scalability on the small-and medium-scaled datasets. ResNet-18 is used as a base encoder throughout the experiments. While changing the base encoder, other architectural details remain the same as in ImageNet-1k. We do not apply multi-crop augmentations while using the EMA. We use image sizes of  $32 \times 32$  and  $256 \times 256$  for small and medium datasets, respectively. Following the procedures in [17], we report the linear evaluation performance

- 277 on the validation set. More experimental details about the optimizer, batch size, augmentations, etc.,
- are provided in the Appendix.

Table 7: Linear evaluation performance in small-and medium-scaled datasets. We report top-1 and top-5 accuracies of linear evaluation on validation dataset. The training results are based on 1000 and 400 epochs of training on CIFAR-10/100 and ImageNet-100, respectively. Results style: best, second best

Method	Arch	CIFAR-10		CIFA	R-100	ImageNet-100	
		Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
BarlowTwins [50]	R18	92.10	99.73	70.90	91.91	80.16	95.14
BYOL [27]	R18	92.58	99.79	70.46	91.96	80.32	94.94
DeepCluster-v2 [8]	R18	88.85	99.58	63.61	88.09	75.36	93.10
DINO [8]	R18	89.52	99.71	66.76	90.34	74.90	92.78
SwAV [8]	R18	89.17	99.68	64.88	88.78	77.83	95.06
MIRA	R18	93.02	<b>99.8</b> 7	70.65	92.23	81.00	95.56

<sup>279</sup> The results are in Table 7. In CIFAR-10 and ImageNet-100, our method outperforms other self-

supervised baselines by 0.4% and 0.7% in top-1 accuracy, respectively. For CIFAR-100, our method

is comparable to the best performing baseline–BarlowTwins; MIRA performs better in top-5 accuracy.

**Training with small batch** Throughout the experiments in Sec. 4.2, we use a batch size of 4096. While such batch size is commonly used in self-supervised methods, large amounts of GPU memory are required; hence limiting the accessibility. In Table 8, we test our method with a smaller batch size of 512 that can be used in an 8 GPU machine with 96GB memory. In this setting, we use the SGD optimizer with a weight decay of  $10^{-4}$ . We also test the robustness of pseudo-labeling with the Sinkhorn-Knopp algorithm in SwAV [8] reproduced by us and compare the results.

We report a top-1 linear evaluation performance of both methods after 100 epochs of training. In the result, the performance gap between our method and SwAV is amplified from 2.9% to 6% in the reduced batch size of 512. One possible explanation is that since SwAV is based on the equipartition constraint, the performance of SwAV harshly degrades when the batch size is not enough to match

the number of clusters.

Table 8: Linear evaluation performance with smaller batch size. All results are based on ImageNet training. We also report the GPU memory usage and time spent for one epoch training. † is result by us.

Method	Batch size	Epochs	GPU	GPU memory	Time per Epoch	Top-1
SwAV† MIRA w/o EMA MIRA	512 512 512	100 100 100	8 × TITAN V 8 × TITAN V 8 × TITAN V	71 GB 71 GB 73 GB	23 min 23 min 29 min	62.3 66.3 68.3
SwAV [12] MIRA w/o EMA MIRA	4096 4096 4096	100 100 100	16 × A100 16 × A100	486 GB 504 GB	9 min 9 min	66.5 68.7 69.4

292

# 293 5 Discussion

**Conclusion** This paper proposes the mutual information maximization inspired pseudo-labeling algorithm MIRA. We formulate pseudo-labeling into an optimization problem and solve it in a principled way. We apply MIRA to representation learning and demonstrate its effectiveness in self-supervised learning benchmarks. We hope that our simple yet theoretically guaranteed approach to information maximization will guide many future applications.

**Limitation and negative social impact** Our information maximization perspective pseudo-labeling seems applicable to various tasks and domains, e.g., semi-supervised training [36]. We validate the effectiveness only in self-supervised visual representation learning. Furthermore, despite our improved training efficiency, the self-supervised learning methods still require a huge amount of computations compared to supervised learning. Such computational requirements may accelerate the environmental problems of global warming.

## 305 **References**

- [1] E. Amrani, L. Karlinsky, and A. Bronstein. Self-supervised classification network. *arXiv preprint arXiv:2103.10994*, 2021.
- Y. M. Asano, C. Rupprecht, and A. Velaldi. Self-labelling via simultaneous clustering and
   representation learning. In *ICLR*, 2020.
- [3] A. Bardes, J. Ponce, and Y. LeCun. Vicreg: variance-invariance-covariance regularization for self-supervised learning. In *ICLR*, 2022.
- [4] M. A. Bautista, A. Sanakoyeu, E. Sutter, and B. Ommer. Cliquecnn: Deep unsupervised exemplar learning. In *NIPS*, 2016.
- [5] L. Bossard, M. Guilaumin, and L. V. Gool. Food-101 mining discriminative components with random forests. In *ECCV*, 2014.
- [6] M. Caron, P. Bojanowski, A. Joulin, and M. Douze. Deep clustering for unsupervised learning of visual features. In *ECCV*, 2018.
- [7] M. Caron, P. Bojanowski, J. Mairal, and A. Joulin. Unsupervised pre-training of image features
   on non-curated data. In *ICCV*, 2019.
- [8] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In *NeurIPS*, 2020.
- [9] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging
   properties in self-supervised vision transformers. In *ICCV*, 2021.
- [10] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020.
- [11] T. Chen, S. Kornblith, K. Swersky, M. Norouzi, and G. Hinton. Big self-supervised models are
   strong semi-supervised learners. In *NeurIPS*, 2020.
- [12] X. Chen and K. He. Exploring simple siamese representation learning. In CVPR, 2021.
- [13] X. Chen, H. Fan, R. Girshick, and K. He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.
- [14] X. Chen, S. Xie, and K. He. An empirical study of training self-supervised vision transformers.
   In *ICCV*, 2021.
- [15] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi. Describing textures in the wild.
   In *CVPR*, 2014.
- [16] M. Cuturi. Lightspeed computation of optimal transport. In *NeurIPS*, 2013.
- [17] V. G. T. da Costa, E. Fini, M. Nabi, N. Sebe, and E. Ricci. solo-learn: a library of self-supervised
   methods for visual representation learning. *Journal of machine learning research*, 2022.
- [18] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: a large-scale hierarchical
   image database. In *CVPR*, 2009.
- [19] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert:pre-training of deep bidirectional
   transformers for language understanding. In *NAACL*, 2019.
- J. Du, E. Grave, B. Gunel, V. Chaudhary, O. Celebi, M. Auli, V. Stoyanov, and A. Conneau.
   Self-training improves pre-training for natural language understanding. In *NAACL*, 2021.
- <sup>344</sup> [21] L. Ericsson, H. Gouk, and T. M. Hospedales. How well do self-supervised models transfer? In <sup>345</sup> *CVPR*, 2021.
- [22] A. Ermolov, A. Siarohin, E. Sangineto, and N. Sebe. Whitening for self-supervised representation learning. In *ICML*, 2021.

- [23] M. Everingham, L. V. Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual
   object classes (voc) challenge. *International journal of computer vision*, 2010.
- [24] L. Fei-Fei, R. Fergus, and P. Perona. Learning generative visual models from few training examples: an incremental bayesian approach tested on 101 object categories. In *CVPR workshops*, 2004.
- W. Feng, K. Tao, Z. Rufeng, L. Huaping, and L. Hang. Self-supervised learning by estimating
   twin class distributions. *arXiv preprint arXiv:2110.07402*, 2021.
- P. Goyal, P. Dollár, R. Girshick, P. Noord-huis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia,
   and K. He. Accurate, large minibatch sgd: training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- [27] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. Á.
   Pires, Z. D. Guo, M. G. Azar, B. Piot, K. Kavukcuoglu, R. Munos, and M. Valko. Bootstrap
   your own latent: a new approach to self-supervised learning. In *NeurIPS*, 2020.
- [28] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [29] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual
   representation learning. In *CVPR*, 2020.
- [30] R. D. Hjelm, A. Fedorov, S. L-Marchildron, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio.
   Learning deep representations by mutual information estimation and maximization. In *ICLR*, 2019.
- [31] W. Hu, T. Miyato, S. Tokui, E. Matsumoto, and M. Sugiyama. Learning discrete representations
   via information maximizing self-augmented training. In *ICML*, 2017.
- [32] S. Ioeffe and C. Szegedy. Batch normalization: accelerating deep network training by reducing
   internal covariate shift. In *ICML*, 2015.
- [33] X. Ji, J. F. Henriques, and A. Veldaldi. Invariant information clustering for unsupervised image classification and segmentation. In *ICCV*, 2019.
- [34] J. Krause, J. Deng, M. Stark, and L. Fei-Fei. Collecting a large-scale dataset of fine-grained cars. In *Second workshop on fine-grained visual categorization*, 2013.
- [35] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *Master's thesis, Department of Computer Science, University of Toronto*, 2009.
- [36] D.-H. Lee. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *ICML Workshop on challenges in representation learning*, 2013.
- [37] J. Li, P. Zhou, C. Xiong, and S. Hoi. Prototypical contrastive learning of unsupervised representations. In *ICLR*, 2021.
- [38] R. Linsker. Self-organization in a perceptual network. *IEEE Computer*, 21:105–117, 1988.
- [39] S. Maji, E. Rahtu, J. Kannala, M. Blaschko, and A. Vedaldi. Fine-grained visual classification
   of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- [40] M. E. Nilsback and A. Zisserman. Automated flower classification over a large number of
   classes. In *Indian conference on computer vision, graphics and image processing (IVCGIP)*,
   2008.
- [41] O. M. Parkhi, A. Vedaldi, A. Zisserman, and C. V. Jawahar. Cats and dogs. In CVPR, 2012.
- [42] T. Piotrowski and R. L. G. Cavalcante. The fixed point iteration of positive concave mappings converges geometrically if a fixed point exists. *arXiv preprint arXiv:2110.11055*, 2022.
- [43] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
   P. Mishkin, J. Clark, G. Kruger, and I. Sutskever. Learning transferable visual models from
   natural language supervision. In *ICML*, 2021.

- [44] Y. Tian, C. Sun, B. Poole, D. Krishnan, C. Schmid, and P. Isola. What makes for good views for
   contrastive learning? In *NeurIPS*, 2020.
- [45] J. Xiao, J. Hays, S. X. Yu, and D. Lin. Sun database: Large-scale scene recognition from abbey to zoo. In *CVPR*, 2010.
- [46] J. Xie, R. Girshick, and A. Farhadi. Unsupervised deep embedding for clustering analysis. In
   *ICML*, 2016.
- [47] X. Yan, I. Misra, A. Gupta, D. Ghadiyaram, and D. Mahajan. Clusterfit: Improving generaliza tion of visual representations. In *CVPR*, 2020.
- [48] J. Yang, D. Parikh, and D. Batra. Joint unsupervised learning of deep representations and image
   clusters. In *CVPR*, 2016.
- [49] Y. You, I. Gitman, and B. Ginsburg. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017.
- I. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow twins: self-supervised learning via
   redundancy reduction. In *ICML*, 2021.
- <sup>408</sup> [51] M. Zheng, S. You, F. Wang, C. Qian, C. Zhang, X. Wang, and C. Xu. Ressl: relational self-<sup>409</sup> supervised learning with weak augmentation. In *NeurIPS*, 2021.

### 410 Checklist

1. For all authors... 411 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's 412 contributions and scope? [Yes] 413 (b) Did you describe the limitations of your work? [Yes] See Section 5 Discussion -414 Limitation and negative social impact 415 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See 416 Section 5 Discussion - Limitation and negative social impact 417 (d) Have you read the ethics review guidelines and ensured that your paper conforms to 418 them? [Yes] 419 2. If you are including theoretical results... 420 (a) Did you state the full set of assumptions of all theoretical results? [Yes] in the Appendix 421 (b) Did you include complete proofs of all theoretical results? [Yes] in the Appendix 422 3. If you ran experiments... 423 (a) Did you include the code, data, and instructions needed to reproduce the main experi-424 425 mental results (either in the supplemental material or as a URL)? [Yes] in supplemental material 426 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they 427 were chosen)? [Yes] in Appendix 428 (c) Did you report error bars (e.g., with respect to the random seed after running experi-429 ments multiple times)? [Yes] in Appendix 430 (d) Did you include the total amount of compute and the type of resources used (e.g., type 431 of GPUs, internal cluster, or cloud provider)? [Yes] 432 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets... 433 (a) If your work uses existing assets, did you cite the creators? [Yes] 434 (b) Did you mention the license of the assets? [Yes] in Appendix 435 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] 436 code in supplemental material. 437 (d) Did you discuss whether and how consent was obtained from people whose data you're 438 using/curating? [N/A] 439

440 441	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
442	5. If you used crowdsourcing or conducted research with human subjects
443 444	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
445 446	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
447 448	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]