

Semi-Supervised Learning of Visual Features by Non-Parametrically Predicting View Assignments with Support Samples

Mahmoud Assran, Mathilde Caron, Ishan Misra, Piotr Bojanowski, Armand Joulin, Nicolas Ballas*, and Michael Rabbat*
 Facebook AI Research

{massran, mathilde, imisra, bojanowski, ajoulin, ballasn, mikerabbat}@fb.com

Abstract

This paper proposes a novel method of learning by predicting view assignments with support samples (PAWS). The method trains a model to minimize a consistency loss, which ensures that different views of the same unlabeled instance are assigned similar pseudo-labels. The pseudo-labels are generated non-parametrically, by comparing the representations of the image views to those of a set of randomly sampled labeled images. The distance between the view representations and labeled representations is used to provide a weighting over class labels, which we interpret as a soft pseudo-label. By non-parametrically incorporating labeled samples in this way, PAWS extends the distance-metric loss used in self-supervised methods such as BYOL and SwAV to the semi-supervised setting. Despite the simplicity of the approach, PAWS outperforms other semi-supervised methods across architectures, setting a new state-of-the-art for a ResNet-50 on ImageNet trained with either 10% or 1% of the labels, reaching 75.5% and 66.5% top-1 respectively. PAWS requires $4\times$ to $12\times$ less training than the previous best methods.

1. Introduction

Learning with less labeled data has been a longstanding challenge of computer vision and machine learning research. One popular approach for learning with few labels is to first perform unsupervised pre-training on a large dataset followed by supervised fine-tuning on the small set of available labels. Self-supervised methods generally adhere to this paradigm (e.g., see [1] for an analysis in the context of semi-supervised learning), and they have demonstrated competitive performance on semi-supervised learning benchmarks across a wide range of self-supervised pre-training strategies [2, 1, 3, 4]. However, the self-supervised

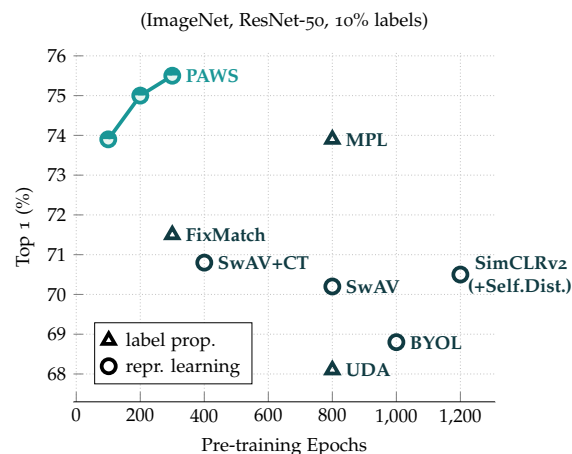


Figure 1: Training a ResNet-50 on ImageNet when only 10% of the training set is labeled. The figure shows top-1 validation accuracy as a function of the number of training epochs. The proposed method, PAWS, achieves higher accuracy than previous work while requiring significantly fewer training epochs. Concretely, 100 epochs of PAWS training takes less than 8.5 hours using 64 NVIDIA V100-16G GPUs.

paradigm also requires substantially more computational effort than other approaches and does not make use of labeled data when it is available.

An alternative line of work suggests to use available labeled data to generate pseudo-labels for the unlabeled data, and then train a model using the labeled and pseudo-labeled data [5, 6, 7, 8, 9, 10, 11]. This begs the question, can we get the best of both worlds, leveraging labeled data throughout training while also building on advances in self-supervised learning?

This paper proposes a novel method of learning by predicting view assignments with support samples (PAWS). The method trains a model to minimize a consistency loss, which ensures that different views of the same unlabeled instance are assigned similar pseudo-labels. The pseudo-labels are generated non-parametrically, by comparing the representations of the image views to those of a set of ran-

*Co-last author

Code: github.com/facebookresearch/suncet

domly sampled labeled images. The distance between the view representations and labeled representations is used to provide a weighting over class labels, which we interpret as a soft pseudo-label. By non-parametrically incorporating labeled samples in this way, PAWS extends the distance-metric loss in self-supervised methods such as BYOL [4] and SwAV [3] to the semi-supervised setting.

Despite the simplicity of the approach, PAWS outperforms other semi-supervised methods across architectures, setting a new state-of-the-art for a ResNet-50 trained on ImageNet with either 10% or 1% of the training instances labeled, achieving 75% and 66% top-1 respectively. Moreover, this is achieved with only 200 epochs of training, which is $4\times$ less than that of the previous best method. The same conclusion holds when training with wider ResNet architectures as well (i.e., ResNet-50 $2\times$ or $4\times$).

2. Related Work

Semi-supervised learning. One procedure to simultaneously learn with both labeled and unlabeled data is to combine a supervised loss on the labeled samples with an unsupervised loss on the unlabeled samples. For example, [12, 13, 14] train a model by adding an unsupervised regularization term to a supervised cross-entropy loss. Similarly, UDA [15] adds a supervised cross-entropy loss to an appropriately weighted unsupervised regularization term. Likewise, S4L [16] adds a supervised cross-entropy loss to a weighted mixture of self-supervised pretext loss terms. This idea of adding a supervised cross-entropy loss to an unsupervised instance-based loss has also been exploited to learn representations suitable for both image classification and instance recognition [17].

There is also a family of semi-supervised methods related to self-training [5] that explicitly generate pseudo-labels for the unlabeled samples and that optimize prediction accuracy on both the ground truth labels (for the labeled samples) and the pseudo-labels (for the unlabeled samples). For example Pseudo-Label [18] and earlier related methods [19, 20, 21] first train a model on the labeled samples, use this model to assign pseudo-labels to unlabeled samples, and then re-train the model using both the labeled and unlabeled samples. The MixMatch trilogy of work [10, 9, 11] operates similarly, but generates the pseudo-labels in an on-line fashion. Specifically, FixMatch [11] trains with a supervised cross-entropy loss while simultaneously making predictions on weakly augmented unlabeled images. When the unsupervised predictions are confident enough, they are used as pseudo-labels for strongly augmented views of those same unlabeled images.

Another closely related line of work in self-training uses an explicit teacher-student configuration. For example, Mean Teacher [8] and Noisy Student [7] use a teacher network to assign pseudo-labels to unlabeled samples, which

are then used to train a student network. Similarly, MPL [6] uses a teacher network to pseudo-label unlabeled images for a student network. The student then performs an update by minimizing its prediction error with respect to the teacher’s pseudo-label. Subsequently, the student is evaluated on a mini-batch of labeled samples, and the teacher network is updated using a meta-learning loss based on the student’s evaluation performance. In MPL, the overall teacher update consists of the combination of the student’s meta-learning loss plus a separate UDA loss. After self-training, the MPL student model is subsequently fine-tuned on the labeled samples using a standard cross-entropy loss.

There is also the Co-training framework [22] which bears a coarse resemblance to the self-training procedure, but possesses notable differences. Specifically, Co-training learns a separate feature extractor on each (conditionally independent) view of the data, combines the predictions of the different feature extractors, and alternates between pseudo-labeling a subset of the data and training on the generated pseudo-labels.

Few-shot learning. In few-shot classification, a network must be adapted to learn to recognize new classes when given only a few labeled examples of these classes [23, 24, 25, 26]. One common approach, which is adopted by Matching Networks [23] and Prototypical Networks [24], is to learn a metric space to embed the data. A differentiable nearest-neighbour classifier is then used in this space to predict the class of a query point given some labeled data-points in the support set [23, 24]. Although there are few-shot approaches that learn entirely from unsupervised data [27], the majority train using labeled data, which is in contrast to the self-supervised approaches discussed next.

Self-supervised learning. Major advances have been made in learning useful image representations from unlabeled data. Some methods take the approach of incorporating domain-specific knowledge in the form of specific pre-training tasks, such as solving jigsaws [28]. More recent success has been achieved by contrasting multiple views of an image [2, 29, 30], where the views come from different random augmentations. Such methods aim to learn a mapping from images to a representation space such that different views of the same image have similar representations. Various approaches have been proposed to avoid the trivial solution of collapsing all images to the same point, including contrasting negative samples [2] and using Sinkhorn-Knopp normalization [31, 3].

It has been demonstrated that self-supervised pre-training produces image representations that can be leveraged effectively for semi-supervised learning [1]. Contrastive self-supervised pre-training generally benefits from training with very large batch sizes, containing sufficiently

many positive and negative examples, and consequently is very computationally expensive, e.g., requiring between 800–1000 epochs of pre-training to learn state-of-the-art representations on ImageNet. Some recent works have demonstrated that the batch-size requirements can be reduced at the expense of maintaining an additional memory bank [32, 29, 30, 4, 3]. Further performance benefits have been obtained by distilling very large pre-trained teacher models to smaller student models [1]. In contrast, PAWS only trains with positive examples, and leverages available annotated data during pre-training to significantly reduce the amount of pre-training required.

3. Methodology

We consider a large dataset of unlabeled images $\mathcal{D} = (\mathbf{x}_i)_{i \in [1, N]}$ and a small support dataset of annotated images $\mathcal{S} = (\mathbf{x}_{s_i}, y_i)_{i \in [1, M]}$, with $M \ll N$.¹ Our goal is to learn image representations by leveraging both \mathcal{D} and \mathcal{S} during pre-training. After pre-training with \mathcal{D} and \mathcal{S} , we fine-tune the learned representations using only the labeled set \mathcal{S} .

3.1. High-level Description

A schematic of the high-level pre-training approach is shown in Figure 2. Given an image \mathbf{x}_i from \mathcal{D} , we use a random set of data augmentations to generate two views, an anchor view $\hat{\mathbf{x}}_i$, and an associated positive view $\hat{\mathbf{x}}_i^+$. Learning proceeds by non-parametrically assigning soft pseudo-labels to the anchor and positive view and subsequently minimizing the cross-entropy $H(\cdot, \cdot)$ between them.

The soft pseudo-labels are generated using a differentiable similarity-based classifier π_d that measures the similarity of a given representation to those of a mini-batch of labeled samples from the support set \mathcal{S} , and outputs a (soft) class label. We use a simple Soft Nearest Neighbours strategy [33] for the similarity classifier π_d .

Connection to few-shot learning. The mini-batch of labeled samples is obtained by first sampling a subset of classes and then sampling a few instances of each class. This, along with the use of a soft nearest-neighbours strategy is similar to approaches previously used for few-shot classification [23]. However, unlike [23], we do not use LSTMs or other mechanisms for encoding or accessing elements of the support set, and furthermore, we never seek to directly predict the labels of elements of the support set. Rather, the support set is only used to assign pseudo-labels to unlabeled image views, and the loss is only evaluated with respect to the pseudo-labels assigned to the unlabeled image views.

¹Note that the images in the support set \mathcal{S} may overlap with the images in the dataset \mathcal{D} .

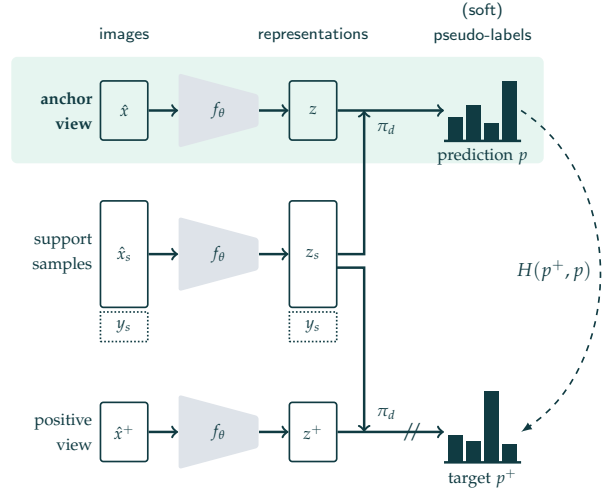


Figure 2: PAWS. The method assigns soft pseudo-labels to an anchor view of an image and an associated positive view, and subsequently minimizes the cross-entropy H between them. The soft pseudo-labels are generated using a differentiable similarity classifier π_d that measures the similarity to a mini-batch of labeled support samples, and outputs a soft class distribution. Positive views are created using data-augmentations of the anchor view. Since the trivial collapse of all representations to a single vector would lead to high-entropy predictions by the similarity classifier, sharpening the target pseudo-labels is sufficient to eliminate all trivial solutions.

3.2. Detailed Methodology

Let $\mathbf{x} \in \mathbb{R}^{n \times (3 \times H \times W)}$ denote a mini-batch of n anchor image views, and let $\mathbf{x}^+ \in \mathbb{R}^{n \times (3 \times H \times W)}$ denote the associated n positive image views. Similarly, let $\mathbf{x}_S \in \mathbb{R}^{m \times (3 \times H \times W)}$ denote a mini-batch of m support images drawn from \mathcal{S} with one-hot class labels $\mathbf{y}_S \in \mathbb{R}^{m \times K}$, where K is the number of classes.

Encoder. Given a parameterized encoder, denoted by $f_\theta : \mathbb{R}^{3 \times H \times W} \rightarrow \mathbb{R}^d$, let $\mathbf{z} \in \mathbb{R}^{n \times d}$ and $\mathbf{z}^+ \in \mathbb{R}^{n \times d}$ denote the representations computed from \mathbf{x} and \mathbf{x}^+ , respectively, and let $\mathbf{z}_S \in \mathbb{R}^{m \times d}$ denote the m support representations computed from \mathbf{x}_S . In our experiments below, the encoder will be the trunk of a deep residual network [34]. The i^{th} representation in the mini-batch \mathbf{z} is written as a row-vector $z_i \in \mathbb{R}^{1 \times d}$, and its associated positive view in the mini-batch is denoted z_i^+ ; i.e., $z_i = f_\theta(\mathbf{x}_i)$ and $z_i^+ = f_\theta(\mathbf{x}_i^+)$. For a scalar-valued similarity function $d(\cdot, \cdot) \geq 0$, the similarity classifier $\pi_d(\cdot, \cdot)$ is given by

$$\pi_d(z_i, \mathbf{z}_S) = \sum_{(z_{s_j}, y_j) \in \mathbf{z}_S} \left(\frac{d(z_i, z_{s_j})}{\sum_{z_{s_k} \in \mathbf{z}_S} d(z_i, z_{s_k})} \right) y_j$$

where y_j is the one-hot ground truth label vector associated with the j^{th} row vector z_{s_j} from \mathbf{z}_S .

Similarity metric and predictions. In this work, we take the similarity metric $d(a, b)$ to be $\exp(a^T b / (\|a\| \|b\| \tau))$, the exponential temperature-scaled cosine. For L2-normalized representations, the similarity classifier $\pi_d(\cdot, \cdot)$ can be concisely written as

$$p_i := \pi_d(z_i, \mathbf{z}_S) = \sigma_\tau(z_i \mathbf{z}_S^\top) \mathbf{y}_S,$$

where $\sigma_\tau(\cdot)$ is the softmax with temperature $\tau > 0$, and $p_i \in [0, 1]^K$ is the prediction for representation z_i .² The positive view predictions p_i^+ are calculated similarly from representations z_i^+ .

To avoid representation collapse, rather than contrast negative samples or incorporate Sinkhorn-Knopp normalization, we compare the prediction of one view with the sharpened prediction of the other view. We define the sharpening function $\rho(\cdot)$ with temperature $T > 0$ as

$$[\rho(p_i)]_k := \frac{[p_i]_k^{1/T}}{\sum_{j=1}^K [p_i]_j^{1/T}}, \quad k = 1, \dots, K.$$

Sharpening the targets encourages the network to produce confident predictions. As will be clear in Section 4, sharpening the targets is provably sufficient to eliminate collapsing solutions in the PAWS framework. Empirically, we have observed that training without sharpening can result in collapsing solutions.

Note that in the case where the support set contains only one instance per sampled class, sharpening the target predictions is equivalent to using a lower temperature in the cosine similarity between the unlabeled representation and support representations. However, when the sampled support set contains more than one instance per sampled class, then sharpening the target predictions is actually different from adjusting the cosine temperature. In this case, it is preferable to sharpen the target predictions rather than use a different temperature in the cosine similarity, since changing the cosine temperature can significantly affect the accuracy of the similarity classifier π_d .

Training objective. To train the encoder, we penalize when the predictions p_i and p_i^+ of two views of the same image are different. As mentioned above, we compare the prediction of one view with the sharpened prediction of the other view; i.e., $H(\rho(p_i), p_i^+) + H(\rho(p_i^+), p_i)$.

We also incorporate a regularization term to encourage the image view representations to utilize the full set of classes represented in the support set. Let $\bar{p} := \frac{1}{2n} \sum_{i=1}^n (\rho(p_i) + \rho(p_i^+))$ denote the average of the sharpened predictions across all unlabeled representations. The

²Specifically, given a vector $a \in \mathbb{R}^K$, the softmax $\sigma_\tau(a) \in [0, 1]^K$ is defined as $[\sigma_\tau(a)]_k := \frac{\exp(a_k/\tau)}{\sum_{j=1}^K \exp(a_j/\tau)}$ for $k = 1, \dots, K$.

regularization term, which we refer to as *mean entropy maximization* (ME-MAX), seeks to maximize the entropy of \bar{p} , denoted $H(\bar{p})$. That is, while the individual predictions are encouraged to be confident, the average prediction is encouraged to be close to the uniform distribution. The ME-MAX regularizer has previously been used in the discriminative unsupervised clustering community for balancing learned cluster sizes (see, e.g., [35]).

Thus, the overall objective to be minimized when training the parameters θ of the encoder f_θ is

$$\frac{1}{2n} \sum_{i=1}^n (H(\rho(p_i^+), p_i) + H(\rho(p_i), p_i^+)) - H(\bar{p})). \quad (1)$$

Note that we only differentiate the cross-entropy loss terms with respect to the predictions p_i and p_i^+ , and not the sharpened targets $\rho(p_i)$ and $\rho(p_i^+)$.

The discussion so far has assumed that we only generate two views for each unlabeled image. One could generate more than two views, in which case we sum the loss over all views and take the target to be the average prediction across the other views of the same image.

The proposed approach seeks to improve on existing self-supervised approaches for semi-supervised learning by: (i) efficiently using available task information, and (ii) addressing representation collapse. On the first issue, since the similarity classifier is differentiable, we evaluate gradients with respect to the labeled samples, but do not directly optimize prediction accuracy on the ground truth labels to avoid overfitting. On the second issue, since the trivial collapse of all representations to a single vector would lead to high-entropy predictions by the similarity classifier, sharpening the target pseudo-labels is sufficient to eliminate all trivial solutions as we will demonstrate in Section 4.

Neural architectures with external memory. PAWS can be interpreted as a neural network architecture with an external memory. Typically, in those architectures, a differentiable neural attention mechanism is used to read and access a memory space which contains a set of elements that are relevant to the task at hand. In PAWS, the support representations \mathbf{z}_S of labeled images characterize the external memory of the network, while the non-parameteric classifier π_d corresponds to the soft-attention operation that retrieves memory elements given a query z_i . From this perspective, PAWS optimizes an encoder network such that two views of the same image activate the same elements in the memory. Moreover, by randomly sampling a subset of labeled images to use as the support set at each iteration, PAWS avoids developing a strong dependence on any particular elements in the memory.

Assimilation & Accommodation. PAWS also has connections to Piaget’s Constructivist learning theory of *assim-*

ilation & accommodation [36], which provided grounding for work in cybernetics [37, Chapter VII].³ At the heart of Constructivism is the idea that every individual possesses representations relating to distinct semantic concepts that are updated through the process of *assimilation and accommodation*. During assimilation, the mind adapts its representations of new observations to fit its *past* observations, while during accommodation, the representations of *past* observations are updated to account for the new observations (cf. Appendix G). In the PAWS procedure, backpropagating with respect to the image views can be seen as a process of assimilation, ensuring that new observations (the image views) are consistent with the current schemata (the support representations). Similarly, backpropagating with respect to the support samples can be seen as a process of accommodation, ensuring that the current schemata (the support representations) are effective at describing the new observations (the image views).

4. Theoretical Guarantees

Next we show that PAWS is guaranteed to avoid the trivial collapse of representations under the following assumptions.

Assumption 1 (Class Balanced Sampling). Each mini-batch of labeled support samples contains an equal number of instances from each of the sampled classes.

Assumption 2 (Target Sharpening). The target p^+ is sharpened, such that it is not equal to the uniform distribution.

Proposition 1 (Non-Collapsing Representations). Suppose Assumptions 1 and 2 hold. If f_θ is such that the representations collapse, i.e., $z_i = z$ for all $z_i \in \mathcal{S}$, then $\|\nabla_\theta H(p^+, p)\| > 0$.

Proof. Since $z = z_i$ for all $z_i \in \mathcal{S}$, it holds that $d(z, z_i) = d(z, z_j)$ for all $z_i, z_j \in \mathcal{S}$. Therefore $p := \pi_d(z, \mathcal{S}) = \frac{1}{n} \sum_{(z_i, y_i)} y_i$, where y_i is the one-hot class label for the representation z_i . Let K denote the number of classes represented in the mini-batch of support samples. By Assumption 1, since the mini-batch of support samples contains an equal number of instances from each sampled class, it follows that there are n/K instances for each of the K represented classes. Therefore, the prediction p further simplifies to $\frac{1}{n} (\mathbf{1}_K \frac{n}{K}) = \frac{1}{K} \mathbf{1}_K$, the uniform distribution over the K classes. However, by Assumption 2, the targets p^+ are sharpened such that they are not equal to the uniform distribution. Therefore, $p \neq p^+$, from which it follows that $\|\nabla H(p^+, p)\| > 0$. ■

³This connection did not readily carry-over to Artificial Intelligence (AI) in the 70's due to the largely symbolic nature of AI approaches at the time; e.g., it was not obvious how to represent the near infinite variations of a hand-drawn curve in a single concise representation; an issue which is now largely resolved by gradient-based learning and modern neural network architectures.

Proposition 1 provides a theoretical guarantee that the proposed method is immune to the trivial collapse of representations. It is also straightforward to extend Proposition 1 to accommodate certain popular transformations of the labels y_i , such as label smoothing. In short, the underlying principle is that collapsing representations result in high entropy predictions under the non-parametric similarity classifier, but the targets are always low-entropy (because we sharpen them), and so collapsing all representations to a single vector is not a stationary point of the training dynamics.

Note that the sharpening function defined in Section 3 may not always satisfy Assumption 2, unless one introduces a simple tie-breaking mechanism. However, in practice, such a mechanism is not necessary as the targets never become uniform (since we apply sharpening from the start of the training). There are also alternative strategies to guarantee the non-collapse of representations without making the target-sharpening assumption, such as by directly using the available class labels for prediction or adding an entropy-minimization term; see Appendix E for more details.

5. Implementation Details

We first pre-train a network using PAWS, and then fine-tune the learned representations for the classification task using only the labeled samples. We also report results using the pre-trained representations directly in a nearest-neighbour classifier.

We adopt similar hyper-parameter settings that have previously been reported in the self-supervised literature [1, 2, 32, 3, 4]. Specifically, for pre-training, we use the LARS optimizer [38] with a momentum value of 0.9, weight decay 10^{-6} , cosine-similarity temperature of $\tau = 0.1$, and batch-size of 4096. We linearly warm-up the learning-rate from 0.3 to 6.4 during the first 10 epochs of pre-training, and decay it following a cosine schedule [39] thereafter.

To construct the different image views, we use the multi-crop strategy from SwAV [3], generating two large crops (224×224), and six small crops (96×96) of each unlabeled image. Each small crop has two positive views (the two large crops), while each large crop has only one positive view (the other large crop).⁴ To construct the support mini-batch at each iteration, we also randomly sample 6720 images, comprising 960 classes and 7 images per class, from the labeled set. For all sampled images (both unlabeled images and support images), we apply the SimCLR data-augmentations [2, 1], specifically random crop, horizontal flip, color distortion, and Gaussian blur. For the sampled support images, we also apply label smoothing with a smoothing factor of 0.1. Lastly, for the target sharpening, we use a temperature of $T = 0.25$.

⁴The target for the small crops is the average of the large crop predictions.

Following previous self-supervised methods, the encoder f_θ in our experiments is a ResNet trunk with a 3-layer MLP projection head [1, 4]. To facilitate comparison with BYOL [4], we also include a 2-layer MLP prediction head, g_ζ , after f_θ , before computing the anchor predictions. Specifically, the representations z and z_S are fed into g_ζ before computing their cosine similarity. While this prediction head is included in our default setup for consistency with previous work, the ablation experiments below (see Table 6), show that PAWS also works well without it. Similar to previous self-supervised methods [1, 2, 4], we also use global batch normalization during pre-training, and exclude the bias and batch-norm parameters from weight decay and LARS adaptation.

After pre-training, we fine-tune a linear classifier from the first layer of the projection head in the encoder f_θ , and follow the evaluation protocol of BYOL [4]. Specifically, we simultaneously fine-tune the encoder/classifier weights using the available labeled samples and a standard supervised cross-entropy loss. See Appendix A for more details, and Section 7 for ablation experiments.

We also report the results of using the pre-trained representations directly in a nearest-neighbour classifier (i.e., without fine-tuning). Specifically, the nearest-neighbour classifier compares the representations of new query images to those of the available labeled data. We refer to this approach as PAWS-NN.

6. Main Results

In this section we analyze the features learned by PAWS on ImageNet [41]. The standard procedure for evaluating semi-supervised methods on ImageNet is to assume that some percentage of the data is labeled, and treat the rest of the data as unlabeled. For reproducibility, we use the same 1% and 10% data splits used in previous works [2, 1].

While we assume that the overall support set contains all relevant labels for the downstream task, we believe this is reasonable since the overall (labeled) support set is small and can be more easily curated. Exploring performance in settings with class imbalance or partial coverage are beyond the scope of this paper and are left as future work.

Baselines. We focus on comparing PAWS to other methods in the literature that train using the same architectures to make a fair comparison. We do not include comparisons with results that first train a larger teacher model and then distill it to a smaller student [1]. For reference, the best reported result in the literature for a ResNet-50 and 1% or 10% labeled data are 73.9% and 77.5% top-1, achieved by distilling from a ResNet-152 with $3\times$ wider channels and selective kernels [1]. We impose this constraint on the baselines to provide a fair comparison and better isolate what factors contribute to performance improvements. It is know

ResNet-50

Method	Epochs	Top 1	
		1%	10%
<i>Methods using label propagation:</i>			
UDA [15]	800	–	68.1
FixMatch [11]	300	–	71.5
MPL [6]	*800	–	73.9
<i>Methods using only representation learning:</i>			
BYOL [4]	1000	53.2	68.8
SwAV [3]	800	53.9	70.2
SwAV+CT [40]	400	–	70.8
SimCLRv2 [1]	800	57.9	68.4
SimCLRv2 (+Self.Dist.) [1]	1200	60.0	70.5
PAWS	100	63.8	73.9
PAWS	200	66.1	75.0
PAWS	300	66.5	75.5
<i>Non-parametric classification (no fine-tuning):</i>			
PAWS-NN	100	61.5	71.0
PAWS-NN	200	63.2	71.9
PAWS-NN	300	64.2	73.1

Table 1: **(ResNet-50, ImageNet)** *For label propagation methods, the number of epochs is counted with respect to the unsupervised mini-batches. *For Meta Pseudo-Labels (MPL), the number of epochs only includes the student-network updates, and does not count the additional 500,000 teacher-network updates (computationally equivalent to an additional 800 epochs) that must happen sequentially (not in parallel) with the student updates. PAWS-NN refers to performing nearest-neighbour classification directly using the PAWS-pretrained representations, with the labeled training samples as support, while PAWS refers to fine-tuning a classifier using the available labeled data after PAWS-pretraining.

Additional ResNet Architectures

Method	Architecture	Epochs	Top 1	
			1%	10%
BYOL [4]	ResNet-50 (2 \times)	1000	62.2	73.5
SimCLRv2 [1]	ResNet-50 (2 \times)	800	66.3	73.9
PAWS	ResNet-50 (2 \times)	100	68.2	77.0
PAWS	ResNet-50 (2 \times)	200	69.6	77.8
SimCLR [2]	ResNet-50 (4 \times)	1000	63.0	74.4
BYOL [4]	ResNet-50 (4 \times)	1000	69.1	75.7
PAWS	ResNet-50 (4 \times)	100	69.8	78.5
PAWS	ResNet-50 (4 \times)	200	69.9	79.0

Table 2: Semi-supervised classification results on ImageNet when training with larger ResNet architectures.

that using distillation in conjunction with larger architectures can result in improvements for any method, and we leave further investigation of distilling larger models pre-trained with PAWS for future work.

Comparison to self-supervised pre-training. We compare PAWS to other self-supervised pre-training ap

proaches, namely SimCLRv2 [1], BYOL [4], SwAV [3], and SwAV+CT [40], which simply adds a supervised contrastive-task loss to SwAV pre-training. Results are reported in Table 1 for a ResNet-50 encoder network and in Figure 1. PAWS outperforms all other self-supervised representation learning approaches while using roughly $10\times$ fewer pre-training epochs. Specifically, with just 100 epochs of pre-training, PAWS surpasses the state-of-the-art in self-supervised representation learning. With 200 epochs of pre-training, PAWS further improves upon this result and achieves 75% top-1 accuracy in the 10% label setting and 66% top-1 in the 1% label setting, setting a new state-of-the-art for a ResNet-50. Using the pre-trained representations directly in a nearest-neighbour classifier (PAWS-NN) also performs surprisingly well—surpassing all other self-supervised representation learning methods—although fine-tuning increases top-1 accuracy by 1–3%. Because PAWS with fine-tuning consistently achieves superior results compared to PAWS-NN, we only report results for PAWS for the remainder of the paper.

By reducing the number of pre-training epochs, PAWS can obtain significant computational savings compared to other approaches. We illustrate this observation by comparing PAWS training time on 64 NVIDIA V100-16G GPUs to the self-supervised SwAV method trained on identical hardware [3]. Pre-training with SwAV for 800 epochs requires 49.6 hours, while pre-training with PAWS for 100 epochs only requires 8.2 hours, and results in a +9.9% improvement in top-1 accuracy in the 1% label setting, and a +3.7% improvement in top-1 accuracy in the 10% label setting. In contrast to SimCLRv2 and BYOL, the PAWS method does not use an additional momentum encoder or a memory buffer, and thereby avoids this added computational and memory overhead, but may also benefit (in terms of final model accuracy) by incorporating such innovations.

Comparison to semi-supervised methods. We also compare PAWS to other semi-supervised learning methods, namely UDA [15], FixMatch [11] and MPL [6]. Results are reported in Table 1 for a ResNet-50 encoder network in the 10% label setting. MPL holds the current state-of-the-art in semi-supervised learning, and simultaneously trains a student and teacher network for 800 epochs by adding a meta-learning loss and a teacher network to the UDA objective. PAWS outperforms MPL, the state-of-the-art semi-supervised learning approach, while requiring significantly fewer training epochs.

Impact of larger architectures. We examine the impact of training larger encoder networks with PAWS pre-training. Specifically, we pre-train ResNet-50 encoders with width multipliers of $2\times$ and $4\times$ in Table 2. As expected, increasing the model capacity improves semi-supervised perfor-

mance. Specifically, pre-training a Resnet-50 ($4\times$) for 200 epochs with PAWS achieves 69.9% top-1 accuracy in the 1% label setting and 79.0% top-1 accuracy in the 10% label setting. We expect increasing the model capacity further to yield additional performance improvements. In general, results with the larger models are consistent with previous observations; PAWS pre-training outperforms other methods using similar architectures, while requiring significantly fewer pre-training epochs.

7. Ablation Study

Learning during pre-training. To further examine the behaviour of PAWS, we examine some metrics related to model quality during pre-training in Figure 3. Figure 3a shows the training cross-entropy loss when pre-training for 100 epochs. As expected, this loss decreases during training, indicating that the model is learning to assign similar pseudo-labels to different views of the same image.

Figure 3b shows two additional losses computed using the sampled mini-batch and support set during training. Here, the instance discrimination loss is the normalized temperature-scaled cross-entropy loss [2] computed using only unlabeled samples in the minibatch, and the classification loss is supervised noise-contrastive estimation loss [40, 42] computed using only labeled samples in the support set. Note that these losses are only computed and reported to better understand PAWS pre-training, and they are not used to train the model. The decreasing instance discrimination loss (top) indicates that the model is learning representations that are invariant to the data augmentations used to construct different views. The decreasing classification loss (bottom) also indicates that the model is learning to correctly classify labeled examples in the support set, despite not directly using labeled examples as targets.

Support set construction. PAWS pre-training requires specifying how to sample a support set. At each iteration, a support set is sampled by first sampling a subset of the K classes, and then sampling a certain number of images per class. We ablate the effect of these two parameters in Table 3. Since we experiment with ImageNet, we can sample up to 1000 classes. Overall, we observe that using a larger support set consistently improves performance. Sampling more classes and fewer samples per class is better than the contrary (cf. bottom two rows). Note that no result is reported for 1000 classes and 16 images per class for the case of 1% labeled data, since in that case there are only 12811 labeled images in total.

Small batch training. Our default PAWS implementation runs on 64 GPUs, with a batch-size of 4096 unlabeled images and a supervised support mini-batch of 6720 images,

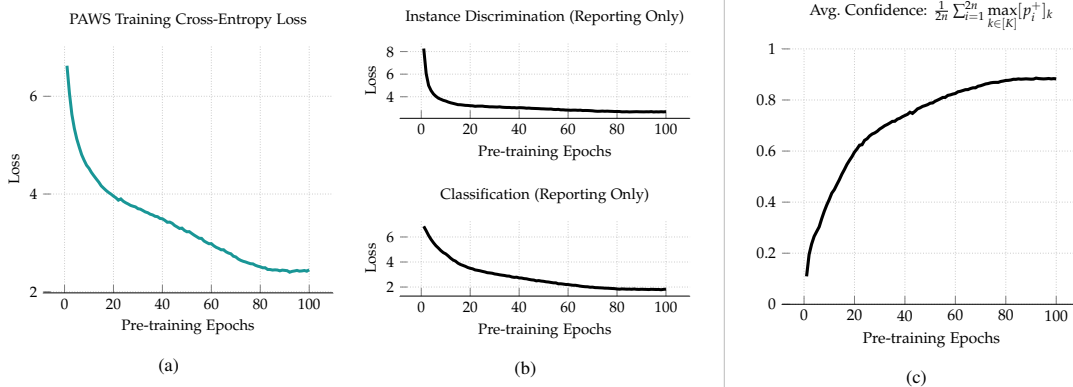


Figure 3: Reporting various metric during training of a ResNet-50 on ImageNet, when 10% of the data is labeled. Fig.3a Cross-entropy loss between anchor view and (target) positive view during training. As expected, this loss decreases during training, indicating that the model is learning to assign similar pseudo-labels to different views of the same image. Fig.3b Additional losses computed with the sampled mini-batch and support-set during training for *reporting purposes only*. Specifically, no gradient is computed with respect to these losses. The decrease in the instance discrimination loss during training suggests that the model is learning representations that are invariant to the data-augmentations used for training. The decrease in the classification loss indicates that the model is learning to correctly classify the labeled support samples. Fig.3c The average confidence of the argmax target prediction during training. As training progresses, the model’s target predictions become increasingly confident.

Classes	Imgs. per Class	Top 1	
		1%	10%
1000	16	–	74.5
1000	12	63.9	74.2
960	7	63.8	73.9
960	4	63.7	72.0
448	8	61.8	70.1

Table 3: **Support Set.** Ablating the composition of the sampled support mini-batches when training a ResNet-50 on ImageNet for 100 epochs. Our default setup is shaded in green. Increasing the size of the support set improves performance. However, when sampling a fixed number of instances, it is preferable to sample many classes with a few images per class, rather than few classes with many images per class.

comprising 960 classes and 7 images per class. We observe that PAWS can also be effectively trained with small batch sizes as well. Table 4 ablates the effect of the batch size when training on 8 NVIDIA V100-16G GPUs, when 10% of the training set is labeled. For this small-batch experiment, we set the unsupervised batch size to 256 and attempt to use as large a support set as is possible on 8 GPUs, since the ablation in Table 3 shows that larger supports lead to better performance. Following a roughly square-root scaling of the learning-rate (relative to the large-batch default setup), we linearly warmup the learning-rate from 0.3 to 1.2 during the first 10 epochs of pre-training, and decay it following a cosine schedule thereafter. We also disable ME-MAX regularization for the small batch experiment, since it is not obvious, a priori, that such regularization will be effective for small batches. All other settings are kept fixed. Table 4 demonstrates that PAWS can still achieve good performance with small batches after only 100 epochs of pre-training on 8 GPUs.

GPUs	Batch Size	Support Set		
		Classes	Imgs. per Class	Top 1
8 V100	256	560	3	70.2
64 V100	4096	448	8	70.1
64 V100	4096	960	7	73.9

Table 4: **Batch Size.** Examining the effect of the batch size when training a ResNet-50 on ImageNet for 100 epochs and 10% of the training set is labeled. PAWS still achieves good performance after only 100 epochs of pre-training with small batch sizes on 8 NVIDIA V100-16G GPUs.

8. Discussion

By leveraging a small labeled support set during pre-training, PAWS achieves competitive classification accuracy for semi-supervised problems and requires significantly less training than previous works. PAWS also provably avoids collapsing solutions, a common challenge in self-supervised approaches.

PAWS can be interpreted as a neural network architecture with an external memory that is trained using the *assimilation & accommodation* principle [36]. During assimilation, PAWS updates the representations of new observations so that they are easily described by its external memory (or schemata), while during accommodation, PAWS updates its external memory to account for the new observations.

The use of a supervised support set has some practical advantages as well, since it enables the model to learn efficiently. However, it remains an interesting question to see if one can learn competitive representations in this framework using only instance supervision and more flexible memory representations. We plan to investigate those directions in future work.

References

- [1] T. Chen, S. Kornblith, K. Swersky, M. Norouzi, and G. Hinton, “Big self-supervised models are strong semi-supervised learners,” *arXiv preprint arXiv:2006.10029*, 2020. [1](#), [2](#), [3](#), [5](#), [6](#), [7](#), [11](#), [15](#), [16](#)
- [2] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A simple framework for contrastive learning of visual representations,” *preprint arXiv:2002.05709*, 2020. [1](#), [2](#), [5](#), [6](#), [7](#), [14](#)
- [3] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, “Unsupervised learning of visual features by contrasting cluster assignments,” *arXiv preprint arXiv:2006.09882*, 2020. [1](#), [2](#), [3](#), [5](#), [6](#), [7](#), [12](#)
- [4] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. A. Pires, Z. D. Guo, M. G. Azar, *et al.*, “Bootstrap your own latent: A new approach to self-supervised learning,” *arXiv preprint arXiv:2006.07733*, 2020. [1](#), [2](#), [3](#), [5](#), [6](#), [7](#), [11](#), [13](#)
- [5] B. Zoph, G. Ghiasi, T.-Y. Lin, Y. Cui, H. Liu, E. D. Cubuk, and Q. V. Le, “Rethinking pre-training and self-training,” *arXiv preprint arXiv:2006.06882*, 2020. [1](#), [2](#)
- [6] H. Pham, Q. Xie, Z. Dai, and Q. V. Le, “Meta pseudo labels,” *arXiv preprint arXiv:2003.10580*, 2020. [1](#), [2](#), [6](#), [7](#), [15](#)
- [7] Q. Xie, M.-T. Luong, E. Hovy, and Q. V. Le, “Self-training with noisy student improves imagenet classification,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10687–10698, 2020. [1](#), [2](#)
- [8] A. Tarvainen and H. Valpola, “Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results,” *arXiv preprint arXiv:1703.01780*, 2017. [1](#), [2](#), [15](#)
- [9] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. A. Raffel, “Mixmatch: A holistic approach to semi-supervised learning,” in *Advances in Neural Information Processing Systems*, pp. 5050–5060, 2019. [1](#), [2](#), [15](#)
- [10] D. Berthelot, N. Carlini, E. D. Cubuk, A. Kurakin, K. Sohn, H. Zhang, and C. Raffel, “Remixmatch: Semi-supervised learning with distribution alignment and augmentation anchoring,” *arXiv preprint arXiv:1911.09785*, 2019. [1](#), [2](#), [15](#)
- [11] K. Sohn, D. Berthelot, C.-L. Li, Z. Zhang, N. Carlini, E. D. Cubuk, A. Kurakin, H. Zhang, and C. Raffel, “Fixmatch: Simplifying semi-supervised learning with consistency and confidence,” *arXiv preprint arXiv:2001.07685*, 2020. [1](#), [2](#), [6](#), [7](#), [15](#)
- [12] Y. Grandvalet and Y. Bengio, “Entropy regularization,” *Semi-supervised learning*, pp. 151–168, 2006. [2](#), [16](#)
- [13] T. Miyato, S.-i. Maeda, M. Koyama, and S. Ishii, “Virtual adversarial training: a regularization method for supervised and semi-supervised learning,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 8, pp. 1979–1993, 2018. [2](#), [15](#)
- [14] V. Verma, K. Kawaguchi, A. Lamb, J. Kannala, Y. Bengio, and D. Lopez-Paz, “Interpolation consistency training for semi-supervised learning,” *arXiv preprint arXiv:1903.03825*, 2019. [2](#), [15](#)
- [15] Q. Xie, Z. Dai, E. Hovy, M.-T. Luong, and Q. V. Le, “Unsupervised data augmentation,” *arXiv preprint arXiv:1904.12848*, 2019. [2](#), [6](#), [7](#), [15](#)
- [16] X. Zhai, A. Oliver, A. Kolesnikov, and L. Beyer, “S4L: Self-supervised semi-supervised learning,” in *Proceedings of the IEEE international conference on computer vision*, pp. 1476–1485, 2019. [2](#)
- [17] M. Berman, H. Jégou, A. Vedaldi, I. Kokkinos, and M. Douze, “Multigrain: a unified image embedding for classes and instances,” *arXiv preprint arXiv:1902.05509*, 2019. [2](#)
- [18] D.-H. Lee, “Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks,” in *In International Conference on Machine Learning Workshop*, 2013. [2](#)
- [19] H. Scudder, “Probability of error of some adaptive pattern-recognition machines,” *IEEE Transactions on Information Theory*, vol. 11, no. 3, 1965. [2](#)
- [20] D. Yarowsky, “Unsupervised word sense disambiguation rivaling supervised methods,” in *In 33rd Annual Meeting of the Association for Computational Linguistics*, 1995. [2](#)
- [21] E. Riloff, “Automatically generating extraction patterns from untagged text,” in *In Proceedings of the National Conference on Artificial Intelligence*, 1996. [2](#)
- [22] A. Blum and T. Mitchell, “Combining labeled and unlabeled data with co-training,” in *Proceedings of the eleventh annual conference on Computational learning theory*, pp. 92–100, 1998. [2](#)
- [23] O. Vinyals, C. Blundell, T. Lillicrap, K. Kavukcuoglu, and D. Wierstra, “Matching networks for one shot learning,” *arXiv preprint arXiv:1606.04080*, 2016. [2](#), [3](#)
- [24] J. Snell, K. Swersky, and R. S. Zemel, “Prototypical networks for few-shot learning,” *arXiv preprint arXiv:1703.05175*, 2017. [2](#)
- [25] S. Ravi and H. Larochelle, “Optimization as a model for few-shot learning,” 2016. [2](#)
- [26] B. M. Lake, T. D. Ullman, J. B. Tenenbaum, and S. J. Gershman, “Building machines that learn and think like people,” *Behavioral and brain sciences*, vol. 40, 2017. [2](#)
- [27] K. Hsu, S. Levine, and C. Finn, “Unsupervised learning via meta-learning,” *arXiv preprint arXiv:1810.02334*, 2018. [2](#)
- [28] I. Misra and L. van der Maaten, “Self-supervised learning of pretext-invariant representations,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6707–6717, 2020. [2](#)
- [29] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, “Momentum contrast for unsupervised visual representation learning,” *arXiv preprint arXiv:1911.05722*, 2019. [2](#), [3](#)
- [30] X. Chen, H. Fan, R. Girshick, and K. He, “Improved baselines with momentum contrastive learning,” *arXiv preprint arXiv:2003.04297*, 2020. [2](#), [3](#)
- [31] Y. M. Asano, C. Rupprecht, and A. Vedaldi, “Self-labelling via simultaneous clustering and representation learning,” *arXiv preprint arXiv:1911.05371*, 2019. [2](#)

- [32] X. Chen and K. He, “Exploring simple siamese representation learning,” *arXiv preprint arXiv:2011.10566*, 2020. **3, 5, 11, 13**
- [33] R. Salakhutdinov and G. Hinton, “Learning a nonlinear embedding by preserving class neighbourhood structure,” in *Artificial Intelligence and Statistics*, pp. 412–419, PMLR, 2007. **3**
- [34] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016. **3**
- [35] A. Joulin and F. Bach, “A convex relaxation for weakly supervised classifiers,” *arXiv preprint arXiv:1206.6413*, 2012. **4**
- [36] J. Piaget, “Cognitive development in children: Piaget,” *Journal of research in science teaching*, vol. 2, no. 3, pp. 176–186, 1964. **5, 8**
- [37] M. A. Boden, *Jean Piaget*. Viking Adult, 1980. **5**
- [38] Y. You, I. Gitman, and B. Ginsburg, “Large batch training of convolutional networks,” *arXiv preprint arXiv:1708.03888*, 2017. **5**
- [39] I. Loshchilov and F. Hutter, “SGDR: Stochastic gradient descent with warm restarts,” *arXiv preprint arXiv:1608.03983*, 2016. **5**
- [40] M. Assran, N. Ballas, L. Castrejon, and M. Rabbat, “Recovering petaflops in contrastive semi-supervised learning of visual representations,” *arXiv preprint arXiv:2006.10803*, 2020. **6, 7, 15**
- [41] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, “Imagenet large scale visual recognition challenge,” *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015. **6, 13**
- [42] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan, “Supervised contrastive learning,” *arXiv preprint arXiv:2004.11362*, 2020. **7**
- [43] A. Krizhevsky, G. Hinton, *et al.*, “Learning multiple layers of features from tiny images,” 2009. **13, 15**
- [44] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, “mixup: Beyond empirical risk minimization,” *arXiv preprint arXiv:1710.09412*, 2017. **14**
- [45] S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, and Y. Yoo, “Cutmix: Regularization strategy to train strong classifiers with localizable features,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6023–6032, 2019. **14**
- [46] E. D. Cubuk, B. Zoph, D. Mane, V. Vasudevan, and Q. V. Le, “Autoaugment: Learning augmentation strategies from data,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 113–123, 2019. **14**
- [47] S. Laine and T. Aila, “Temporal ensembling for semi-supervised learning,” *arXiv preprint arXiv:1610.02242*, 2016. **15**
- [48] J. Jackson and J. Schulman, “Semi-supervised learning by label gradient alignment,” *arXiv preprint arXiv:1902.02336*, 2019. **15**
- [49] X. Wang, D. Kihara, J. Luo, and G.-J. Qi, “Enact: Self-trained ensemble autoencoding transformations for semi-supervised learning,” *arXiv preprint arXiv:1911.09265*, 2019. **15**
- [50] S. Zagoruyko and N. Komodakis, “Wide residual networks,” *arXiv preprint arXiv:1605.07146*, 2016. **15**
- [51] M. A. Boden, “Artificial intelligence and piagetian theory,” *Synthese*, pp. 389–414, 1978. **17**
- [52] J. Piaget, “Biology and knowledge: An essay on the relations between organic regulations and cognitive processes.,” 1971. **17**
- [53] J. S. Bruner, “Reply to individual and collective problems in the study of thinking,” *Annals of the New York Academy of Sciences*, vol. 91, no. 1, pp. 22–37, 1961. **17**