

Acknowledgments

RJ and SA acknowledge support from the U.S. National Science Foundation under award # 1931978. All authors are grateful for computing infrastructure support from the Tufts High-Performance Computing cluster, partially funded by NSF under grant OAC CC* 2018149.

Appendix Contents

A Code and Data Resources for Reproducibility	13
B Dataset Details	14
B.1. Example Images	14
B.2. Dataset Selection	14
B.3. Classification Task Description	14
C Additional Results	16
C.1. Impact of pretraining on accuracy-over-time profiles	16
C.2. Validation-set profiles of accuracy-over-time	17
C.3. Additional performance metrics: Profiles over time on AIROGS	18
C.4. Variability in Performance Across Trials	18
D Method Details	20
D.1. Algorithm : Unified training and hyperparameter tuning via random search on a budget	20
D.2. Semi-supervised method details	21
D.3. Self-supervised method details	21
E Additional Analysis and Discussion	22
E.1. Effectiveness of Hyperparameter Tuning	22
E.2. Differentiating Between Methods	22
E.3. Answers to Common Questions from Reviewers	23
F. Hyperparameter Details	24
F.1. Hyperparameter Tuning Strategy: Random Search Details	24
F.2. Hyperparameter transfer strategy	25

A. Code and Data Resources for Reproducibility

All code and data resources needed to reproduce our analysis, including information on exact splits we used for each of the 4 datasets (TissueMNIST, PathMNIST, TMED-2, AIROGS) can be found in our github repo:

<https://github.com/tufts-ml/SSL-vs-SSL-benchmark>

Primer on our codebase. Our codebase builds upon the open-source PyTorch repo by Suzuki [72]. Suzuki’s code was originally intended as a reimplementation in PyTorch of Oliver et al. [62]’s benchmark of semi-supervised learning (while Oliver et al’s original repo was in Tensorflow, we prefer PyTorch).

We added many additional algorithms (we added MixMatch, FixMatch, FlexMatch, and CoMatch, as well as all 7 self-supervised methods) and customized the experiments, especially providing a runtime-budgeted hyperparameter tuning strategy as outlined in App. D.

In a way, this makes our repo a “cousin” of the codebase of Su et al. [71]’s fine-grained classification benchmark, because their [github repo](#) also credits Suzuki’s repo as an ancestor.

B. Dataset Details

B.1. Example Images

Below we show a few examples for each dataset. For full details, please refer to the original papers.

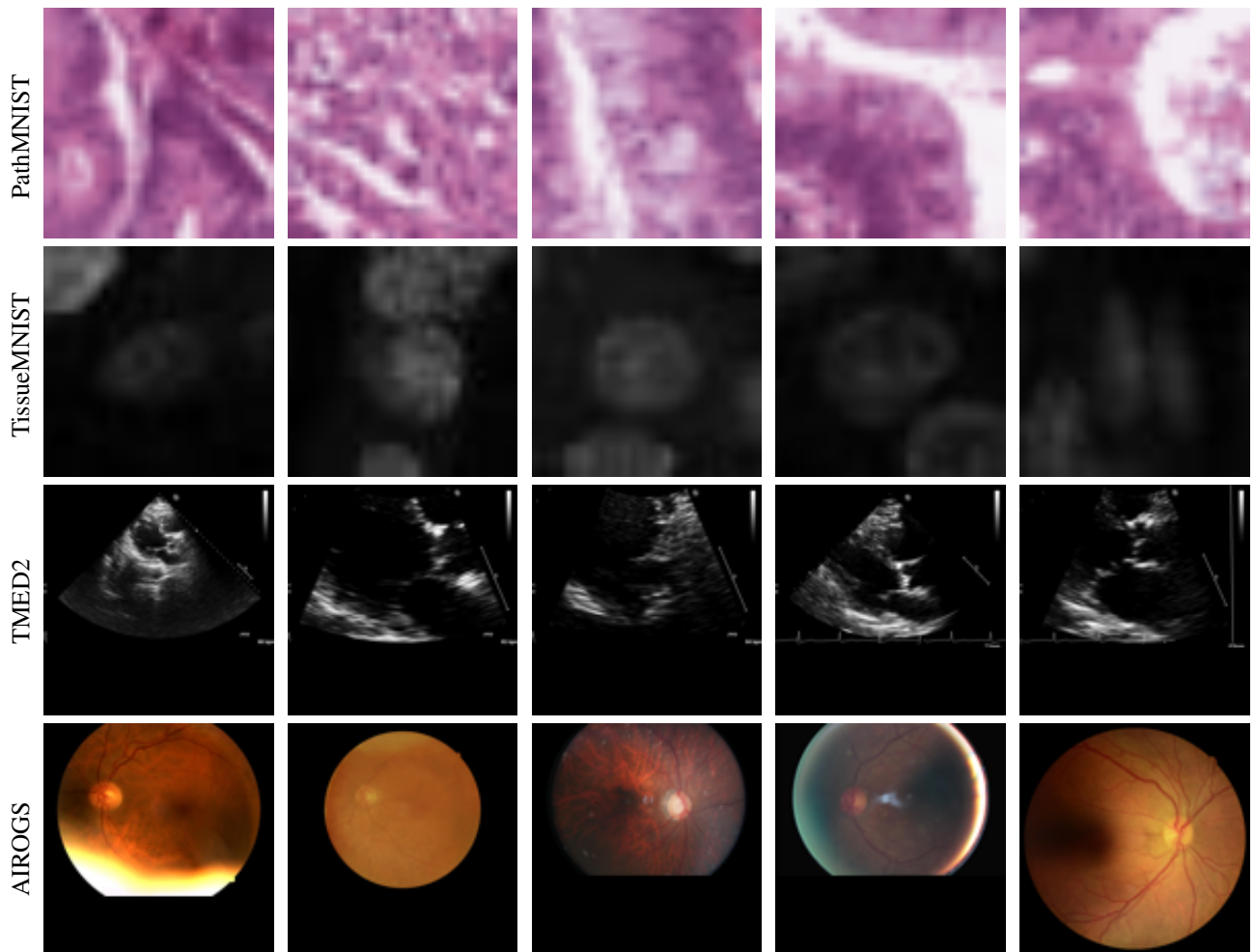


Figure B.1. Showing 5 random examples for each dataset.

B.2. Dataset Selection

We selected PathMNIST and TissueMNIST from 12 candidate datasets in the MedMNIST collections [79, 80] by matching two criteria: (i) contains at least 5 imbalanced classes; (ii) can build a large unlabeled set (at least 50000 images). Prior experiments from dataset creator Yang et al. [80] suggest 28x28 resolution is a reasonable choice. They report that a larger resolution (224x224) does not yield much more accurate classifiers for these two datasets.

B.3. Classification Task Description

TissueMNIST contains 28x28 images of human kidney cortex cells. The dataset contains 8 classes. See [80] for details.

Class ID	Abbreviation	Description
0	CD/CT	Collecting Duct, Connecting Tubule
1	DCT	Distal Convoluted Tubule
2	GE	Glomerular endothelial cells
3	IE	Interstitial endothelial cells
4	LEU	Leukocytes
5	POD	Podocytes
6	PT	Proximal Tubule Segments
7	TAL	Thick Ascending Limb

PathMNIST contains 28x28 patches from colorectal cancer histology slides that comprise 9 tissue types. See [48, 80] for details.

Class ID	Abbreviation	Description
0	ADI	adipose
1	BACK	background
2	DEB	debris
3	LYM	lymphocytes
4	MUC	mucus
5	MUS	smooth muscle
6	NORM	normal colon mucosa
7	STR	cancer-associated stroma
8	TUM	colorectal adenocarcinoma epithelium

TMED-2 contains 112x112 2D grayscale images captured from routine echocardiogram scans (ultrasound images of the heart). In this study, we adopt the view classification task from [39]. For more detail please see [39, 40]

Class ID	Abbreviation	Description
0	PLAX	parasternal long axis
1	PSAX	parasternal short axis
2	A2C	apical 2-chamber
3	A4C	apical 4-chamber

AIROGS is a dataset of color fundus photographs of the retina. The binary classification task is to detect evidence of referable glaucoma [24]. We use 384x384 resolution, as suggested by several challenge participants.

Class ID	Abbreviation	Description
0	No Glauc.	no referable glaucoma
1	Glaucoma	referable glaucoma (signs associated with visual field defects on standard automated perimetry)

C. Additional Results

C.1. Impact of pretraining on accuracy-over-time profiles

To study the impact of pretraining, we compare the accuracy-over-time profiles of TissueMNIST and PathMNIST based on the two different initialization strategy. Fig. C.1 shows balanced-accuracy-over-time profiles for initialization of neural net parameters to values pretrained on ImageNet (left column) and random initialization (right column). Pretraining time on a source dataset is NOT counted to the runtime reported in x-axis.

On TissueMNIST (top row), SimCLR (green) and BYOL (blue) are the top two methods for both initialization types. Performance gains from pretraining are slight, BA for BYOL is around 42 with pretraining and 40 with random initialization.

On PathMNIST (bottom row), FixMatch and CoMatch are best in the pretraining case, with MixMatch and Flexmatch only a few points of balanced accuracy lower. MixMatch and CoMatch are best in the random initialization case.

Across both datasets, pretraining does not seem to impact the **top-performing methods**' ultimate accuracy by much, usually just a slight increase in BA of 0.5-3 points. One exception is FixMatch on PathMNIST, which improves by about 5 percentage points. We do not see the 10+ point gains reported by Su et al. [71] in their Table 3.

Considering more limited time budgets (e.g. after only a few hours), we do see initialization from pretraining understandably tends to improve some methods.

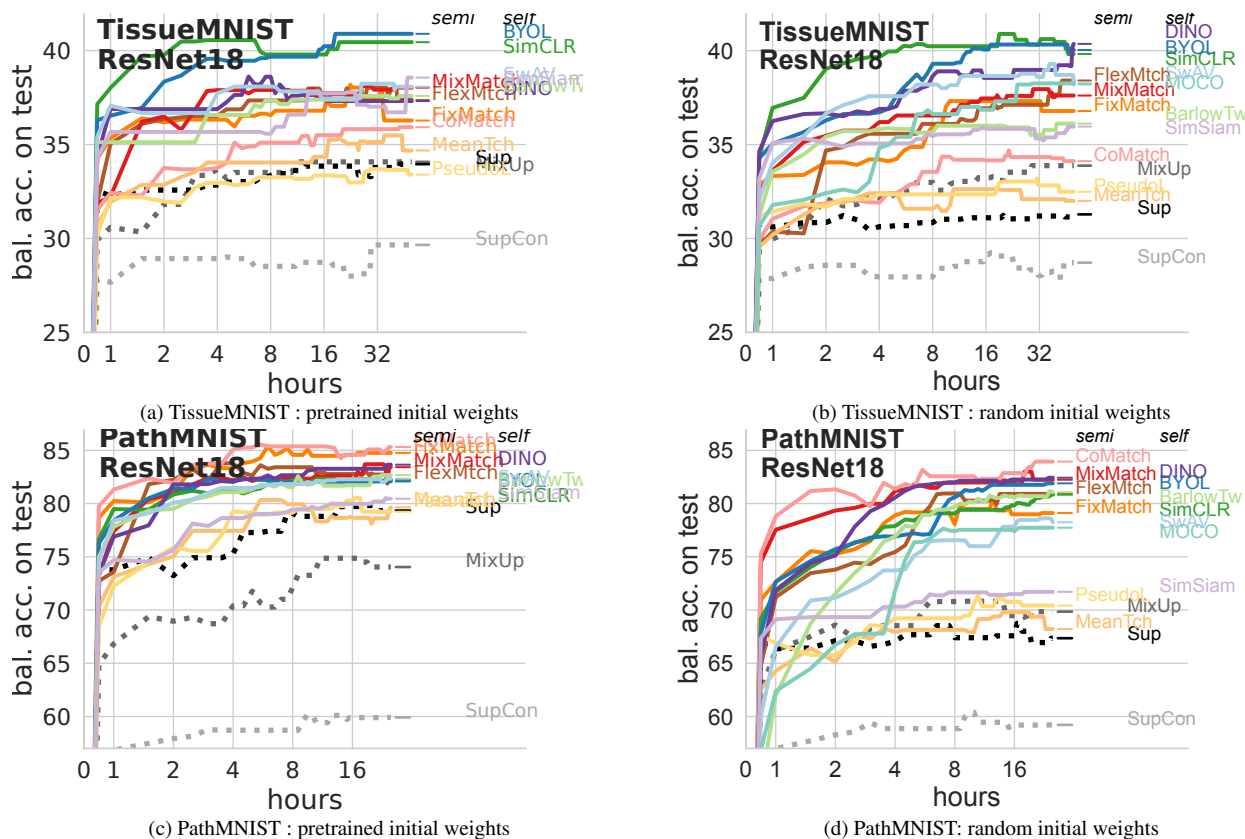


Figure C.1. Balanced accuracy on test set over time for semi- and self-supervised methods, **with (left) and without (right) initial weight pretraining on ImageNet**. Curves represent mean of each method at each time over 5 trials of Alg. D.1.

C.2. Validation-set profiles of accuracy-over-time

Fig. C.2 shows profiles of accuracy over time on the validation set, in contrast to the test set performance shown in the main paper’s Fig. 1.

All curves here by definition must be monotonically increasing, because our unified algorithm selects new checkpoints only when they improve the validation-set balanced accuracy metric. The important insight our work reveals is that the same model checkpoints selected here, based on validation-set accuracy, also tend to produce improved test-set accuracy over time (in Fig. 1). This helps provide empirical confidence in using *realistically-sized* validation sets.

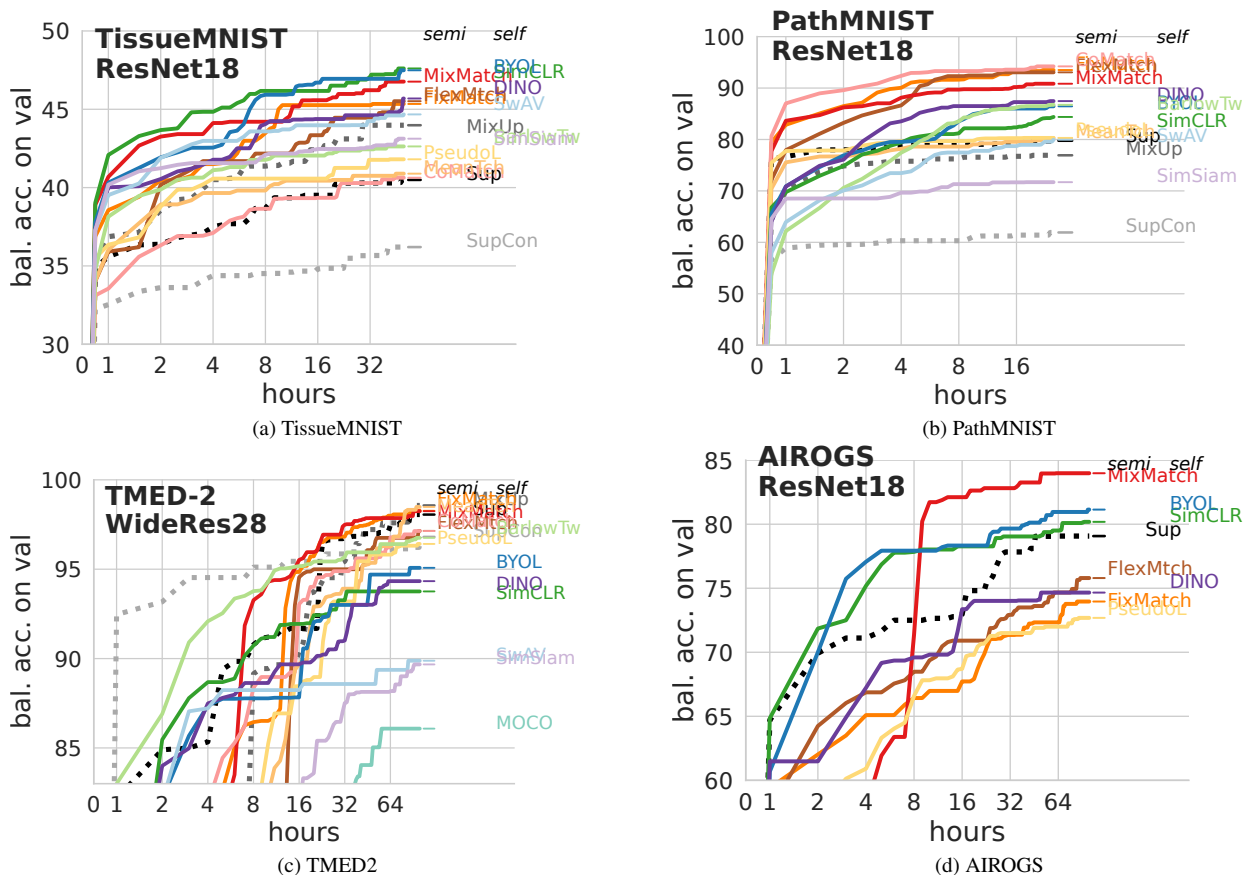


Figure C.2. **Validation-set** accuracy over time profiles of semi- and self-supervised methods on 4 datasets (panels a-d). All curves here by definition must be monotonically increasing. The increasing profiles here on the validation set translate to similar trends in test set performance in Fig. 1, indicating successful generalization.

C.3. Additional performance metrics: Profiles over time on AIROGS

In Fig. C.3, we report the test performance over time on the AIROGS dataset across all 4 metrics of interest, including the partial AUROC and sensitivity at 95% specificity metrics recommended by the AIROGS data creators as being particularly relevant for the glaucoma detection task.

Broadly, our takeaway is that our proposed hyperparameter tuning method is viable for all these metrics, not just the BA metric covered in the main paper. Furthermore, this viability appears consistent across both ResNet-18 and ResNet-50 architectures.

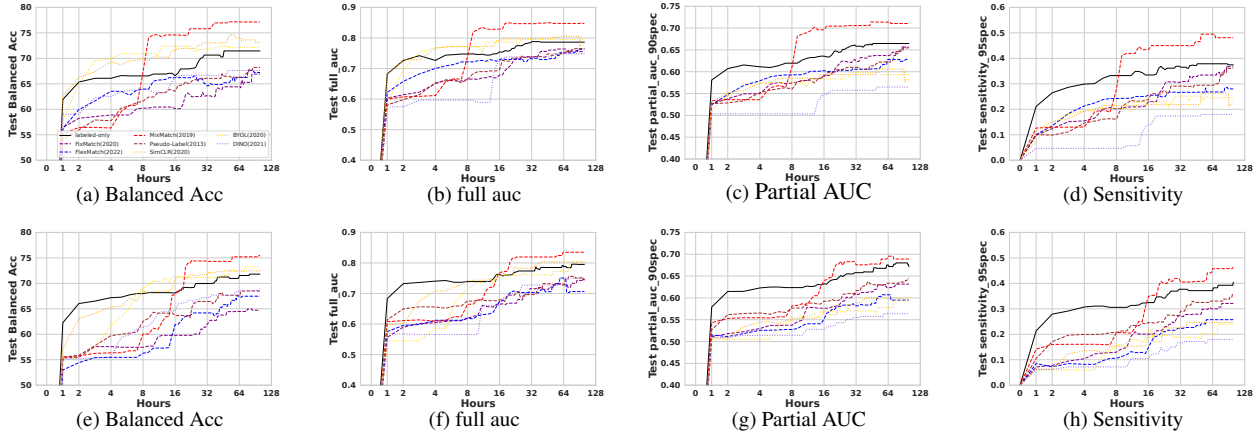


Figure C.3. **Profiles of several clinically-relevant performance metrics over time on the AIROGS test set.** *Top row:* ResNet-18. *Bottom row:* ResNet-50. *Columns, left-to-right:* Balanced Accuracy, AUROC, Partial AUROC focused on the 90% - 100% specificity regime, and sensitivity at 95% specificity. At each time, we report mean of each method over 5 trials of Alg. D.1.

C.4. Variability in Performance Across Trials

In Fig. C.4 on the next page, we explicitly visualize the variability in performance of each method across the 5 separate trials of Alg. D.1 (most other figures show the mean of these 5 trials for visual clarity).

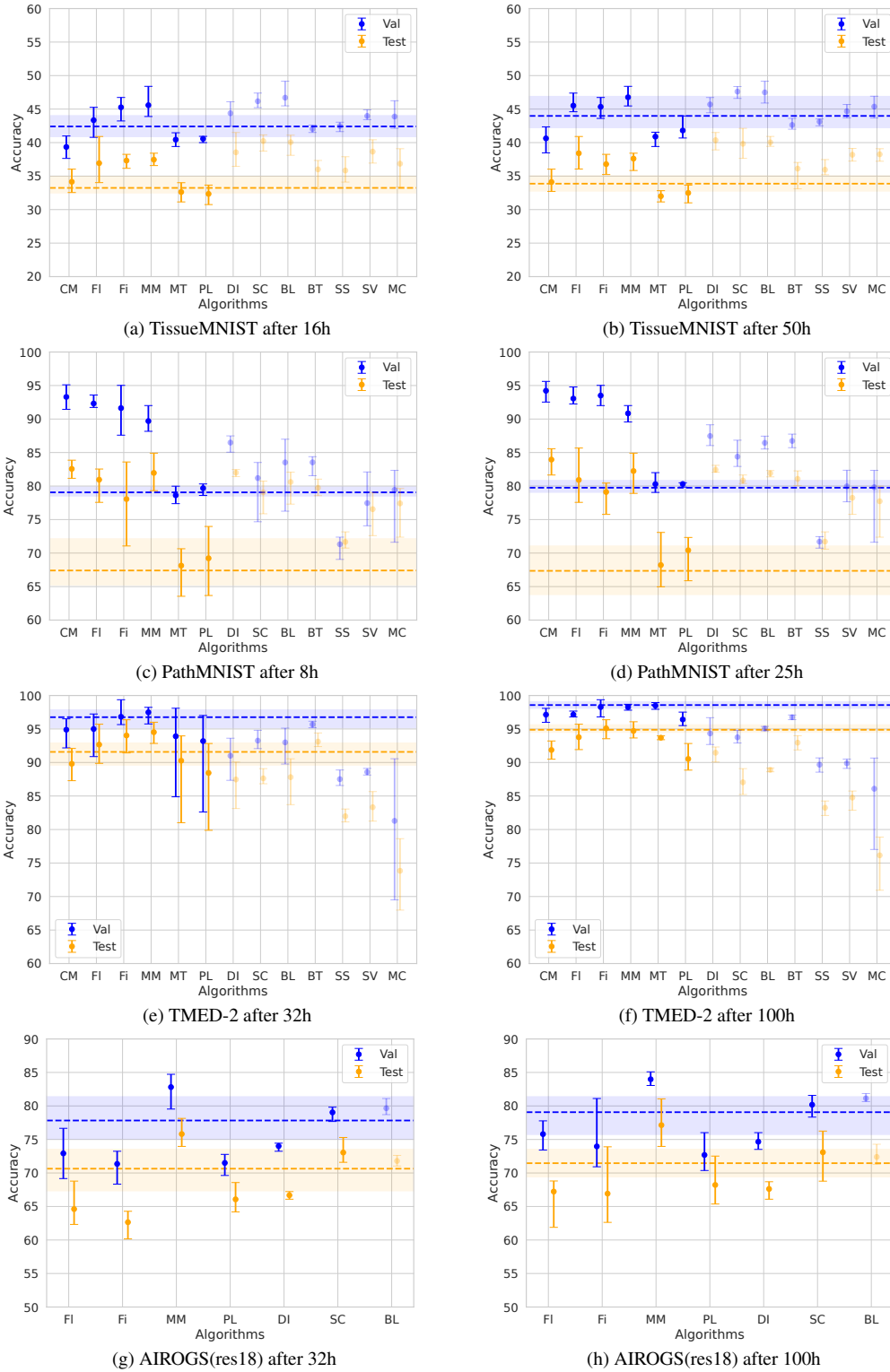


Figure C.4. Balanced accuracy of different methods across 2 time budgets (columns) and four datasets (rows). For each method, the interval indicates the low and high performance of 5 separate trials of Alg. D.1, while dot indicates the mean performance. Horizontal lines indicate the best labeled-set-only baseline at that time. Abbreviation: CM, FI, Fi, MM, MT, PL, DI, SC, BL, BT, SS, SV, MC denote CoMatch, FlexMatch, FixMatch, MixMatch, Mean Teacher, Pseudo Label, DINO, SimCLR, BYOL, Barlow Twins, SimSiam, SwAV, MOCO (v2).

D. Method Details

D.1. Algorithm : Unified training and hyperparameter tuning via random search on a budget

Algorithm D.1 outlines the hyperparameter tuning procedures used across all algorithms under comparison. The algorithm requires three sources of data: a labeled training set $\mathcal{L} = \{X, Y\}$, an unlabeled set for training $\mathcal{U} = X^U$, and a separate realistically-sized labeled validation set $\{X^{val}, Y^{val}\}$. We further require some budget restrictions: a common computational budget T (maximum number of hours), and a maximum training epoch per hyperparameter configuration E .

We proceed as follows: We begin by randomly sampling a hyperparameter configuration from a defined range (see Appendix F.1 for details). A model is then initialized and trained using the ADAM optimizer with the sampled hyperparameters. Each configuration is trained for a maximum of E (200) epochs or stopped early if the validation performance does not improve for 20 consecutive epochs. The model’s performance on the validation set is measured using balanced accuracy. Upon completion of training for a given hyperparameter configuration (either after reaching maximum epoch E or after early stopping), a new configuration is sampled and the process repeats until the total compute budget T is expended.

We track the best-so-far model performance every 30 minutes, and save the best-so-far model along with its validation and test performance. Semi-supervised algorithms simultaneously train the representation layers v and classifier layer w , while self-supervised algorithms train the representation layers v for each epoch and then fine-tune a linear classifier with weights w anew at the end of each epoch using an sklearn logistic regression model [64] with representation parameters v frozen.

Algorithm D.1 Unified Procedure for Training + Hyperparameter selection via random search

Input:

- Train set of features \mathbf{X} paired with labels \mathbf{Y} , with extra unlabeled features \mathbf{U}
- Validation set of features \mathbf{X}^{val} and labels \mathbf{Y}^{val}
- Runtime budget T , Max Epoch E

Output: Trained weights $\{v, w\}$, where v is the representation module, w is the classifier layer

```

1: while time elapsed < T do
2:    $\lambda \sim \text{DRAWHYPERS}$                                 ▷ Sample hyperparameters from pre-defined range (App. F.1)
3:    $\xi \leftarrow \text{CREATEOPTIM}(\lambda)$                        ▷ Initialize stateful optimizer e.g., ADAM
4:    $\{v, w\} \sim \text{INITWEIGHTS}$                              ▷ Initialize model weights
5:   for epoch  $e$  in  $1, 2, \dots, E$  do
6:     if self-supervised then
7:        $v \leftarrow \text{TRAINONEEPOCH}(\mathbf{U}, v, \lambda, \xi)$       ▷ Optimize Eq. (1) with  $\lambda^L = 0$ 
8:        $w \leftarrow \text{TRAINCLASSIFIER}(\mathbf{Y}, f_v(\mathbf{X}))$ 
9:     else if semi-supervised then
10:       $v, w \leftarrow \text{TRAINONEEPOCH}(\mathbf{X}, \mathbf{Y}, \mathbf{U}, v, w, \lambda, \xi)$   ▷ Optimize Eq. (1)
11:     else
12:       $v, w \leftarrow \text{TRAINONEEPOCH}(\mathbf{X}, \mathbf{Y}, v, w, \lambda, \xi)$   ▷ Optimize Eq. (1) with  $\lambda^U = 0$ 
13:     end if
14:      $m_e \leftarrow \text{CALCPERF}(\mathbf{X}^{val}, \mathbf{Y}^{val}, v, w)$       ▷ Record performance metric on val.
15:     if first try or  $m_e > m_*$  then
16:        $v_*, w_* \leftarrow v, w$ 
17:        $\lambda_* \leftarrow \lambda$ 
18:        $m_* \leftarrow m_e$                                 ▷ Update best config found so far
19:     end if
20:     if EARLYSTOP( $m_1, m_2, \dots, m_e$ ) or time elapsed >  $T$  then
21:       break
22:     end if
23:   end for
24: end while
25: return  $v_*, w_*, \lambda_*, m_*$ 

```

D.2. Semi-supervised method details

Semi-supervised learning trains on the labeled and unlabeled data simultaneously, usually with the total loss being a weighted sum of a labeled loss term and an unlabeled loss term. Different methods mainly differ in how unlabeled data is used to form training signals. Many approaches have been proposed and refined over the past decades. These include co-training, which involves training multiple classifiers on various views of the input data [7, 61]; graph-structure-based models [46, 89]; generative models [52, 53]; consistency regularization-based models that enforce consistent model outputs [5, 57, 73]; pseudo-label-based models that impute labels for unlabeled data [12, 58]; and hybrid models that combine several methods [69]. Comprehensive reviews can be found in Chapelle et al. [14], Van Engelen and Hoos [74], Zhu [88].

Among the deep classifier methods following Eq. (1), below we describe each method we selected and how its specific unlabeled loss is constructed.

Pseudo-Labeling uses the current model to assign class probabilities to each sample in the unlabeled batch. If, for an unlabeled sample, the maximum class probability $P(y_i)$ exceeds a certain threshold τ , this sample contributes to the calculation of the unlabeled loss for the current batch. The cross-entropy loss is computed as if the true label of this sample is class i .

Mean-Teacher constructs the unlabeled loss by enforcing consistency between the model’s output for a given sample and the output of the same sample from the Exponential Moving Average (EMA) model.

MixMatch uses the MixUp [85] technique on both labeled data (features and labels) and unlabeled data (features and guessed labels) within each batch to produce transformed labeled and unlabeled data. The labeled and unlabeled losses are then calculated using these transformed samples. Specifically, the unlabeled loss is derived from the mean squared error between the model’s output for the transformed unlabeled samples and their corresponding transformed guessed labels.

FixMatch generates two augmentations of an unlabeled sample, one with weak augmentation and the other using strong augmentations (e.g., RandAug [22]). The unlabeled loss is then formulated by enforcing the model’s output for the strongly augmented sample to closely resemble that of the weakly augmented sample using cross-entropy loss.

FlexMatch builds directly upon FixMatch by incorporating a class-specific threshold on the unlabeled samples during training.

CoMatch marks the first introduction of contrastive learning into semi-supervised learning. The model is equipped with two distinct heads: a classification head, which outputs class probabilities for a given sample, and a projection head, which maps the sample into a low-dimensional embedding. These two components interact in a unique manner. The projection head-derived embeddings inform the similarities between different samples, which are then used to refine the pseudo-labels against which the classification head is trained. Subsequently, these pseudo-labels constitute a pseudo-label graph that trains the embedding graph produced by the projection head.

D.3. Self-supervised method details

In recent years, self-supervised learning algorithms have emerged rapidly and are known as one of the most popular fields of machine learning. These include contrastive learning, which involves learning representations by maximizing agreement between differently augmented views of the same data [17, 37]; predictive models that forecast future instances in the data sequence [63]; generative models that learn to generate new data similar to the input [16]; clustering-based approaches that learn representations by grouping similar instances [9, 10]; context-based models that predict a specific part of the data from other parts [8, 25]; and hybrid models that combine various methods for more robust learning [18]. A more comprehensive review can be found in [47, 90].

Below, we provide for each selected self-supervised method a summary of its internal workings.

SimCLR generates two augmented versions of each image. Then feed these pairs of images into a base encoder network to generate image embeddings. This encoder is followed by a projection head, which is a multilayer neural network, to map these embeddings to a space where contrastive loss can be applied. Next, calculate the contrastive loss. The idea is to make the embeddings of augmented versions of the same image (positive pairs) as similar as possible and to push apart embeddings from different images (negative pairs). The loss function used is NCE loss.

MOCO V2 creates two augmented versions of each image. These pairs are processed by two encoder networks: a query encoder, and a key encoder updated by a moving average of the query encoder. The contrastive loss is computed by comparing a positive pair (the query and corresponding key) against numerous negative pairs drawn from a large queue of keys.

Note on runtime: We notice that the performance on MoCo can be increased when Shuffling BN across multiple GPUs. However, to ensure a fair comparison given our single-GPU setup, we refrained from employing any techniques to simulate multiple GPUs on one, as this would change the encoder’s structure.

SwAV begins by creating multiple augmented versions of each image. Then, these versions are input into a deep neural network to generate embeddings. Uses a clustering approach, called online stratified sampling, to predict assignments of each view’s prototypes (or cluster centers) to others, encouraging the model to match the representations of different augmentations of the same image.

Note on runtime: We’ve observed that applying multiple augmentations can enhance the effectiveness of various methods. To prevent the results from being influenced by these augmentations, we’ve standardized the number of augmentations to two in SwAV, in line with the approach taken by other methods.

BYOL starts by creating two differently augmented versions of each image. These versions are processed through two identical neural networks, known as the target and online networks, which include a backbone and a projection head. The online network is updated through backpropagation, while the target network’s weights are updated as a moving average of the online network’s weights. The unique aspect of BYOL is that it learns representations without the need for negative samples.

SimSiam creates two differently augmented versions of each image. These versions are passed through two identical networks: one predictor network and one encoder network. The encoder network contains a backbone and a projection head.

DINO utilizes two differently augmented images, processed by a student and a teacher network. The teacher’s weights evolve as a moving average of the student’s. The key idea is self-distillation, where the student’s outputs match the teacher’s for one view but differ for the other, without traditional negative samples.

Barlow Twins processes two augmented views of an image through identical networks. The aim is to have similar representations between these networks while minimizing redundancy in their components, sidestepping the need for contrasting positive and negative pairs.

E. Additional Analysis and Discussion

E.1. Effectiveness of Hyperparameter Tuning

While Oliver et al. [62] caution that extensive hyperparameter search may be futile with realistic validation set. Our experiments on the 4 dataset show that the validation set performance for each examined algorithm rise substantially over the course of hyperparameter tuning. This increase in validation set performance further translates to increased test set performance.

Given the trends we observed across 4 datasets, we think that for a chosen algorithm on a new dataset, following our hyperparameter tuning protocol (even with limited labeling budget and computation budget), we can likely expect better generalization (measured by test set performance) compared to not tuning hyperparameters at all.

E.2. Differentiating Between Methods

Oliver et al. [62] offer both empirical and theoretical analysis of how well one can distinguish if one method is truly better than another on a limited labeled dataset. Below, we revisit each analysis for our specific experiments.

E.2.1 Empirical Analysis of Differentiation

Oliver et al. [62] in their Fig 5 and 6 show that on SVHN, between 10 random samples of the validation set across several level of validation set size (1000, 500, 250, 100), the validation accuracy of the trained Pi-model, VAT, Pseudo-labeling and Mean Teacher model has substantial variability and overlap with each other. Thus, they caution that differentiating between models might be infeasible with realistic validation set size.

In our present study, we employ a relaxed notion of “realistic validation set”, by letting the validation set to be at most as large as the training set. Our experiments cover validation set size 235 (TMED), 400 (Tissue), 450 (Path), 600 (AIROGS); test set size 2019 (TMED2), 47280 (Tissue), 7180 (Path), 6000 (AIROGS). Our experiment shows that within the wide range of methods considered, differentiating between some models are possible. For example, in Fig. C.4 we can see that MixMatch is clearly better than Mean Teacher in TissueMNIST and PathMNIST, in both the validation set and test set, without overlap in the intervals. The field of semi-supervised learning has made significant advancements in recent years. It is crucial to reevaluate previous conclusions in light of the new developments.

E.2.2 Theoretical Analysis of Differentiation

Here, we show that the performance gain we observe on the test set are real. We perform the same theoretical analysis using the Hoeffding’s inequality [38] as in Oliver et al. [62].

$$\mathbf{P}(|\bar{V} - \mathbb{E}[V]| < p) > 1 - 2 \exp(-2np^2) \quad (3)$$

where \bar{V} is the empirical estimate of some model performance metric, $\mathbb{E}[V]$ is its hypothetical true value, p is the desired maximum deviation between our estimate and the true value, and n is the number of examples used.

On TissueMNIST, we have 47280 test samples, we will be more than 99.98% confident that the test accuracy is within 1% of its true value. On Path, we have 7180 test samples, we will be more than 99% confident that the test accuracy is within 2% its true value.

In Fig 1, we see that after hyperparameter tuning, the final test accuracy of each algorithms improves much more than 1% on TissueMNIST and 2% on PathMNIST showing the efficacy of hyperparameter tuning.

Similarly, we can see that the difference between top-performing algorithms (e.g., MixMatch) and worst-performing algorithm (e.g., Mean Teacher) is clearly larger than 1% on TissueMNIST, 2% on PathMNIST. Thus we can argue that differentiation between certain methods are viable. The same analysis can also be applied to TMED-2 and AIROGS.

E.3. Answers to Common Questions from Reviewers

Here we answer a few questions that were common to several reviewers of our paper.

E.3.1 For a medical image application, would a larger labeled dataset be more important than than developing semi-supervised or self-supervised methods?

Yes, in general, is preferable to collect as large of a labeled dataset as possible, at least up to the point of performance saturation. Investing in data collection likely has a larger payoff than investing in SSL. However, extensive collection of labeled examples is **not practical** for many real-world clinical tasks due to reasons like financial cost, logistics, privacy and legal issues (see Oliver et al. [62], Berthelot et al. [5], Shekoofeh et al. [68]).

For this reason, methods for overcoming limited labeled data, such as semi-SL and self-SL, are **important topics in medical imaging applications**. The clinical use case of SSL motivates several recent methodological works, such Zhang et al. [86], Azizi et al. [3], and Shekoofeh et al. [68].

E.3.2 Isn’t it already well-known that hyperparameter tuning with a realistic-sized validation set is viable?

When labeled data is abundant, as in common *supervised* learning settings, hyperparameter tuning is widely known as effective. However, our work focuses on the *semi/self-SL* setting, where labels are limited. We carefully reviewed semi/self-SL literature and argue that the viability of tuning on *realistic-sized* validation sets is **not well-known** in this setting. As our paper’s Table 1 shows, existing SSL benchmarks often use validation sets larger than the training set! Seminal work by Oliver et al. [62] cautions that “Extensive hyperparameter tuning may be somewhat futile due to an excessively small collection of held-out data ...”. Su et al. [71] use a similar claim to justify not doing *any tuning* on their Semi-Fungi dataset experiments.

E.3.3 Does MixMatch outperform Flex/Fix/CoMatch because RandAug not suitable for medical imaging?

In general, RandAug-type augmentation can be successful for medical imaging tasks [15, 86], though we agree that it might not be “optimal”. Instead, we hypothesize that MixMatch’s primary advantage is lower runtime cost per iteration compared to FixMatch and successors. In our AIROGS ResNet-18 experiments (Fig. 1), MixMatch explores at least 80% more hyperparameter combinations than its counterparts (111 vs. 59 for FixMatch).

F. Hyperparameter Details

F.1. Hyperparameter Tuning Strategy: Random Search Details

Below, in a specific table for each of the 16 methods (supervised, semi-, or self-), we provide a method-specific table showing the random sampling distribution used for each hyperparameter for the random search of Alg. D.1.

Settings Common to All Methods	
Optimizer	Adam
Learning rate schedule	Cosine

F.1.1 Supervised Baselines

Labeled only	
Batch size	64
Learning rate	$3 \times 10^x, X \sim \text{Unif}(-5, -2)$
Weight decay	$4 \times 10^x, X \sim \text{Unif}(-6, -3)$
MixUp	
Batch size	64
Learning rate	$3 \times 10^x, X \sim \text{Unif}(-5, -2)$
Weight decay	$4 \times 10^x, X \sim \text{Unif}(-6, -3)$
Beta shape α	$x, X \sim \text{Unif}(0.1, 10)$

Sup Contrast	
Batch size	256
Learning rate	$3 \times 10^x, X \sim \text{Unif}(-5.5, -1.5)$
Weight decay	$4 \times 10^x, X \sim \text{Unif}(-7.5, -3.5)$
Temperature	$x, X \sim \text{Unif}(0.05, 0.15)$

F.1.2 Semi-Supervised Methods

FlexMatch	
Labeled batch size	64
Unlabeled batch size	448
Learning rate	$3 \times 10^x, X \sim \text{Unif}(-5, -2)$
Weight decay	$4 \times 10^x, X \sim \text{Unif}(-6, -3)$
Unlabeled loss coefficient	$10^x, X \sim \text{Unif}(-1, 1)$
Unlabeled loss warmup schedule	No warmup
Pseudo-label threshold	0.95
Sharpening temperature	1.0
FixMatch	
Labeled batch size	64
Unlabeled batch size	448
Learning rate	$3 \times 10^x, X \sim \text{Unif}(-5, -2)$
Weight decay	$4 \times 10^x, X \sim \text{Unif}(-6, -3)$
Unlabeled loss coefficient	$10^x, X \sim \text{Unif}(-1, 1)$
Unlabeled loss warmup schedule	No warmup
Pseudo-label threshold	0.95
Sharpening temperature	1.0
CoMatch	
Labeled batch size	64
Unlabeled batch size	448
Learning rate	$3 \times 10^x, X \sim \text{Unif}(-5, -2)$
Weight decay	$4 \times 10^x, X \sim \text{Unif}(-6, -3)$
Unlabeled loss coefficient	$10^x, X \sim \text{Unif}(-1, 1)$
Unlabeled loss warmup schedule	No warmup
Contrastive loss coefficient	$5 \times 10^x, X \sim \text{Unif}(-1, 1)$
Pseudo-label threshold	0.95
Sharpening temperature	0.2

MixMatch	
Labeled batch size	64
Unlabeled batch size	64
Learning rate	$3 \times 10^x, X \sim \text{Unif}(-5, -2)$
Weight decay	$4 \times 10^x, X \sim \text{Unif}(-6, -3)$
Beta shape α	$x, X \sim \text{Unif}(0.1, 1)$
Unlabeled loss coefficient	$7.5 \times 10^x, X \sim \text{Unif}(0, 2)$
Unlabeled loss warmup schedule	linear
Sharpening temperature	0.5
Mean Teacher	
Labeled batch size	64
Unlabeled batch size	64
Learning rate	$3 \times 10^x, X \sim \text{Unif}(-5, -2)$
Weight decay	$4 \times 10^x, X \sim \text{Unif}(-6, -3)$
Unlabeled loss coefficient	$8 \times 10^x, X \sim \text{Unif}(-1, 1)$
Unlabeled loss warmup schedule	linear
Pseudo-label	
Labeled batch size	64
Unlabeled batch size	64
Learning rate	$3 \times 10^x, X \sim \text{Unif}(-5, -2)$
Weight decay	$4 \times 10^x, X \sim \text{Unif}(-6, -3)$
Unlabeled loss coefficient	$10^x, X \sim \text{Unif}(-1, 1)$
Unlabeled loss warmup schedule	Linear
Pseudo-label threshold	0.95

F.1.3 Self-supervised Methods

SwAV	
Batch size	256
Learning rate	$1 \times 10^x, X \sim \text{Unif}(-4.5, -1.5)$
Weight decay	$1 \times 10^x, X \sim \text{Unif}(-6.5, -3.5)$
Temperature	$x, X \sim \text{Unif}(0.07, 0.12)$
Num. prototypes	$1 \times 10^x, X \sim \text{Unif}(1, 3)$
MoCo	
Batch size	256
Learning rate	$1 \times 10^x, X \sim \text{Unif}(-4.5, -1.5)$
Weight decay	$1 \times 10^x, X \sim \text{Unif}(-6.5, -3.5)$
Temperature	$x, X \sim \text{Unif}(0.07, 0.12)$
Momentum	$x, X \sim \text{Unif}(0.99, 0.9999)$
SimCLR	
Batch size	256
Learning rate	$1 \times 10^x, X \sim \text{Unif}(-4.5, -1.5)$
Weight decay	$1 \times 10^x, X \sim \text{Unif}(-6.5, -3.5)$
Temperature	$x, X \sim \text{Unif}(0.07, 0.12)$
SimSiam	
Batch size	256
Learning rate	$1 \times 10^x, X \sim \text{Unif}(-4.5, -1.5)$
Weight decay	$1 \times 10^x, X \sim \text{Unif}(-6.5, -3.5)$

BYOL	
Batch size	256
Learning rate	$1 \times 10^x, X \sim \text{Unif}(-4.5, -1.5)$
Weight decay	$1 \times 10^x, X \sim \text{Unif}(-6.5, -3.5)$
Temperature	$x, X \sim \text{Unif}(0.07, 0.12)$
Momentum	$x, X \sim \text{Unif}(0.99, 0.9999)$
DINO	
Batch size	256
Learning rate	$1 \times 10^x, X \sim \text{Unif}(-4.5, -1.5)$
Weight decay	$1 \times 10^x, X \sim \text{Unif}(-6.5, -3.5)$
Temperature	$x, X \sim \text{Unif}(0.07, 0.12)$
Momentum	$x, X \sim \text{Unif}(0.99, 0.9999)$
Barlow Twins	
Batch size	256
Learning rate	$1 \times 10^x, X \sim \text{Unif}(-4.5, -1.5)$
Weight decay	$1 \times 10^x, X \sim \text{Unif}(-6.5, -3.5)$
Temperature	$x, X \sim \text{Unif}(0.07, 0.12)$
Momentum	$x, X \sim \text{Unif}(0.99, 0.9999)$

F.2. Hyperparameter transfer strategy

To make the most of limited labeled data, one potential strategy recommended by Su et al. [71] is to use the entire labeled set for training, reserving no validation set at all. This relies on pre-established hyperparameters from other dataset/experiments. In this study, we experiment with two scenarios: using pre-determined hyperparameters tuned for CIFAR-10, or using hyperparameters tuned for TissueMNIST.

The CIFAR-10 hyperparameters are sourced from repositories published by each method’s original authors, as this is a common benchmark in the SSL literature. We ensure that each hyperparameter choice, when applied using the re-implemented code for each method in our codebase, matches previously reported results on CIFAR-10.

The TissueMNIST hyperparameters originate from our experiments as depicted in Figure C.2 (a). For exact values, see App. F.2.1.

For each method using the transfer strategy, we perform training on the combined train+validation set, setting the maximum number of epochs to 100 for PathMNIST and AIROGS (80 epochs for TMED2). Training is terminated early if the train loss does not improve over 20 consecutive epochs. Empirically, we observe that all models which did not trigger early stopping reached a plateau in training loss.

F.2.1 Best Hyperparameters on TissueMNIST for Semi-Supervised Methods

Below we report the chosen hyperparameters on TissueMNIST for each semi-supervised method, as used in the hyperparameter transfer experiments.

	FlexMatch				
	seed0	seed1	seed2	seed3	seed4
Learning rate	0.00036	0.00016	0.00016	0.00068	0.00006
Weight decay	0.00259	0.00001	0.00371	0.00023	0.002103
Unlabeled loss coefficient	2.22	0.82	5.00	1.94	6.09

	FixMatch				
	seed0	seed1	seed2	seed3	seed4
Learning rate	0.00074	0.00034	0.00392	0.00102	0.00037
Weight decay	0.00045	0.00315	0.00001	0.00005	0.00058
Unlabeled loss coefficient	3.08	6.70	1.85	1.46	0.47

CoMatch					
	seed0	seed1	seed2	seed3	seed4
Learning rate	0.00124	0.00145	0.00061	0.00026	0.00113
Weight decay	0.00042	0.00009	0.00005	0.00009	0.00017
Unlabeled loss coefficient	0.30	1.71	1.26	2.74	0.46
Contrastive loss coefficient	1.26	2.21	3.71	0.56	1.37

MixMatch					
	seed0	seed1	seed2	seed3	seed4
Learning rate	0.00028	0.00003	0.00018	0.00009	0.00005
Weight decay	0.000005	0.00195	0.00005	0.00085	0.00082
Beta shape α	0.2	0.9	0.9	0.8	0.7
Unlabeled loss coefficient	9.13	37.96	8.06	25.16	11.17

Mean Teacher					
	seed0	seed1	seed2	seed3	seed4
Learning rate	0.00062	0.00022	0.00005	0.00128	0.00125
Weight decay	0.00189	0.00001	0.00008	0.00001	0.00001
Unlabeled loss coefficient	67.67	0.87	1.25	7.60	13.56

Pseudo-label					
	seed0	seed1	seed2	seed3	seed4
Learning rate	0.00007	0.00021	0.00005	0.00063	0.00060
Weight decay	0.00033	0.00093	0.00383	0.00005	0.00087
Unlabeled loss coefficient	0.19	0.16	8.73	0.82	0.25

F.2.2 Best Hyperparameters on TissueMNIST for Self-Supervised Methods

Below we report the chosen hyperparameters on TissueMNIST for each self-supervised method, as used in the hyperparameter transfer experiments.

SwAV					
	seed0	seed1	seed2	seed3	seed4
Learning rate	0.00065	0.00325	0.00012	0.00086	0.00196
Weight decay	0.0001497	0.0000056	0.0000006	0.0000021	0.0000003
Num. prototypes	845	131	36	201	59

MoCo					
	seed0	seed1	seed2	seed3	seed4
Learning rate	0.00288	0.00023	0.00043	0.00005	0.02629
Weight decay	0.000002	0.0000008	0.0000003	0.0000005	0.0000004
temperature	0.09331	0.07097	0.10987	0.07414	0.07080
Momentum	0.99242	0.99672	0.99267	0.99950	0.99538

SimCLR					
	seed0	seed1	seed2	seed3	seed4
Learning rate	0.00217	0.00131	0.000640	0.00380	0.00136
Weight decay	0.00002	0.00001	0.00001	0.00001	0.00001
temperature	0.11719	0.10426	0.08652	0.07784	0.11478

SimSiam					
	seed0	seed1	seed2	seed3	seed4
Learning rate	0.0002	0.00056	0.00013	0.00338	0.00098
Weight decay	0.000066	0.000046	0.000023	0.000001	0.000001

BYOL					
	seed0	seed1	seed2	seed3	seed4
Learning rate	0.000245	0.001308	0.000371	0.001653	0.001959
Weight decay	0.0000007	0.0000057	0.0000004	0.000003	0.000001
Momentum	0.9928618	0.996167	0.9988484	0.9940063	0.9934791

DINO					
	seed0	seed1	seed2	seed3	seed4
Learning rate	0.000245	0.001308	0.000371	0.001653	0.001959
Weight decay	0.0000007	0.0000057	0.0000004	0.000003	0.000001
Momentum	0.9928618	0.996167	0.9988484	0.9940063	0.9934791

Barlow Twins					
	seed0	seed1	seed2	seed3	seed4
Learning rate	0.000245	0.001308	0.000371	0.001653	0.001959
Weight decay	0.0000007	0.0000057	0.0000004	0.000003	0.000001
Momentum	0.9928618	0.996167	0.9988484	0.9940063	0.9934791

References

- [1] N Ahmadi, MY Tsang, AN Gu, TSM Tsang, and P Abolmaesumi. Transformer-based spatio-temporal analysis for classification of aortic stenosis severity from echocardiography cine series. *IEEE Transactions on Medical Imaging*, 2023. 8
- [2] Shekoofeh Azizi, Basil Mustafa, Fiona Ryan, Zachary Beaver, Jan Freyberg, Jonathan Deaton, Aaron Loh, Alan Karthikesalingam, Simon Kornblith, Ting Chen, Vivek Natarajan, and Mohammad Norouzi. Big Self-Supervised Models Advance Medical Image Classification. In *International Conference on Computer Vision (ICCV)*. arXiv, 2021. 4
- [3] Shekoofeh Azizi, Laura Culp, Jan Freyberg, Basil Mustafa, Sebastien Baur, Simon Kornblith, Ting Chen, Nenad Tomasev, Jovana Mitrović, Patricia Strachan, S. Sara Mahdavi, Ellery Wulczyn, Boris Babenko, Megan Walker, Aaron Loh, Po-Hsuan Cameron Chen, Yuan Liu, Pinal Bavishi, Scott Mayer McKinney, Jim Winkens, Abhijit Guha Roy, Zach Beaver, Fiona Ryan, Justin Krogue, Mozziyar Etemadi, Umesh Telang, Yun Liu, Lily Peng, Greg S. Corrado, Dale R. Webster, David Fleet, Geoffrey Hinton, Neil Houlsby, Alan Karthikesalingam, Mohammad Norouzi, and Vivek Natarajan. Robust and data-efficient generalization of self-supervised machine learning for diagnostic imaging. *Nature Biomedical Engineering*, 7(6):756–779, 2023. 4, 23
- [4] James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *Journal of machine learning research*, 13(2), 2012. 6
- [5] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. *Advances in neural information processing systems*, 32, 2019. 1, 3, 21, 23
- [6] Benjamin Billot, Colin Magdamo, Steven E Arnold, Sudeshna Das, and Juan Eugenio Iglesias. Robust segmentation of brain mri in the wild with hierarchical cnns and no retraining. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 538–548. Springer, 2022. 5
- [7] Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In *Proceedings of the eleventh annual conference on Computational learning theory*, pages 92–100, 1998. 21
- [8] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 21
- [9] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *Proceedings of the European conference on computer vision (ECCV)*, pages 132–149, 2018. 21
- [10] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems*, 33:9912–9924, 2020. 1, 3, 21
- [11] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 3
- [12] Paola Cascante-Bonilla, Fuwen Tan, Yanjun Qi, and Vicente Ordonez. Curriculum labeling: Revisiting pseudo-labeling for semi-supervised learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021. 21
- [13] Leo Anthony Celi, Jacqueline Cellini, Marie-Laure Charpignon, Edward Christopher Dee, Franck Dernoncourt, Rene Eber, William Greig Mitchell, Lama Moukheiber, Julian Schirmer, et al. Sources of bias in artificial intelligence that perpetuate healthcare disparities—A global review. *PLOS Digital Health*, 1(3), 2022. 8
- [14] Olivier Chapelle, Bernhard Scholkopf, and Alexander Zien. Semi-supervised learning (chapelle, o. et al., eds.; 2006)[book reviews]. *IEEE Transactions on Neural Networks*, 20(3):542–542, 2009. 21
- [15] Chen Chen, Chen Qin, Huaqi Qiu, Cheng Ouyang, Shuo Wang, Liang Chen, Giacomo Tarroni, Wenjia Bai, and Daniel Rueckert. Realistic adversarial data augmentation for mr image segmentation. In *Medical Image Computing and Computer Assisted Intervention MICCAI*, 2020. 23
- [16] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pre-training from pixels. In *International conference on machine learning*, pages 1691–1703. PMLR, 2020. 21
- [17] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 1, 2, 3, 21
- [18] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E Hinton. Big self-supervised models are strong semi-supervised learners. *Advances in neural information processing systems*, 33:22243–22255, 2020. 21
- [19] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15750–15758, 2021. 1, 3
- [20] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 1, 3
- [21] Yanbei Chen, Massimiliano Mancini, Xiatian Zhu, and Zeynep Akata. Semi-supervised and unsupervised deep visual learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 1, 2, 4
- [22] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 702–703, 2020. 6, 21
- [23] Victor Guilherme Turrissi Da Costa, Enrico Fini, Moin Nabi, Nicu Sebe, and Elisa Ricci. solo-learn: A library of self-supervised methods for visual representation learning. *J. Mach. Learn. Res.*, 23(56):1–6, 2022. 4

- [24] Coen de Vente, Koenraad A. Vermeer, Nicolas Jaccard, He Wang, Hongyi Sun, Firas Khader, Daniel Truhn, Temirgali Aimyshev, Yerkebulan Zhanibekuly, Tien-Dung Le, Adrian Galdran, Miguel Ángel González Ballester, Gustavo Carneiro, Devika R. G, Hrishikesh P. S, Densen Puthussery, Hong Liu, Zekang Yang, Satoshi Kondo, Satoshi Kasai, Edward Wang, Ashritha Durvasula, Jónathan Heras, Miguel Ángel Zapata, Teresa Araújo, Guilherme Aresta, Hrvoje Bogunović, Mustafa Arikian, Yeong Chan Lee, Hyun Bin Cho, Yoon Ho Choi, Abdul Qayyum, Imran Razzak, Bram van Ginneken, Hans G. Lemij, and Clara I. Sánchez. AIROGS: Artificial Intelligence for ROBust Glaucoma Screening Challenge, 2023. [5](#), [15](#)
- [25] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. [21](#)
- [26] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. [8](#)
- [27] Linus Ericsson, Henry Gouk, and Timothy M Hospedales. How well do self-supervised models transfer? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5414–5423, 2021. [1](#), [3](#), [4](#)
- [28] Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. *nature*, 542(7639):115–118, 2017. [5](#)
- [29] Loveleen Gaur, Ujwal Bhatia, NZ Jhanjhi, Ghulam Muhammad, and Mehedi Masud. Medical image-based detection of covid-19 using deep convolution neural networks. *Multimedia systems*, 29(3):1729–1738, 2023.
- [30] Hemant Ghayvat, Muhammad Awais, AK Bashir, Sharnil Pandya, Mohd Zuhair, Mamoon Rashid, and Jamel Nebhen. Ai-enabled radiologist in the loop: novel ai-based framework to augment radiologist performance for covid-19 chest ct medical image annotation and classification from pneumonia. *Neural Computing and Applications*, 35(20):14591–14609, 2023.
- [31] Amirata Ghorbani, David Ouyang, Abubakar Abid, Bryan He, Jonathan H Chen, Robert A Harrington, David H Liang, Euan A Ashley, and James Y Zou. Deep learning interpretation of echocardiograms. *NPJ digital medicine*, 3(1):10, 2020. [5](#)
- [32] Priya Goyal, Dhruv Mahajan, Abhinav Gupta, and Ishan Misra. Scaling and benchmarking self-supervised visual representation learning. In *Proceedings of the IEEE/CVF International Conference on computer vision*, pages 6391–6400, 2019. [4](#)
- [33] Margherita Grandini, Enrico Bagli, and Giorgio Visani. Metrics for multi-class classification: an overview. *arXiv preprint arXiv:2008.05756*, 2020. [5](#)
- [34] Jean-Bastien Grill, Florian Strub, Florent Althé, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020. [3](#)
- [35] Isabelle Guyon, Kristin Bennett, Gavin Cawley, Hugo Jair Escalante, Sergio Escalera, Tin Kam Ho, Núria Macià, Bisakha Ray, Mehreen Saeed, Alexander Statnikov, et al. Design of the 2015 ChaLearn AutoML challenge. In *2015 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2015. [5](#)
- [36] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [5](#)
- [37] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. [1](#), [3](#), [21](#)
- [38] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *The collected works of Wassily Hoeffding*, pages 409–426, 1994. [23](#)
- [39] Zhe Huang, Gary Long, Benjamin Wessler, and Michael C Hughes. A new semi-supervised learning benchmark for classifying view and diagnosing aortic stenosis from echocardiograms. In *Proceedings of the Machine Learning for Healthcare Conference*. PMLR, 2021. [5](#), [15](#)
- [40] Zhe Huang, Gary Long, Benjamin S Wessler, and Michael C Hughes. TMED 2: A dataset for semi-supervised classification of echocardiograms. In *DataPerf: Benchmarking Data for Data-Centric AI Workshop*, 2022. [5](#), [15](#)
- [41] Zhe Huang, Mary-Joy Sidhom, Benjamin S Wessler, and Michael C Hughes. Fix-a-step: Semi-supervised learning from uncurated unlabeled data. In *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2023. [5](#), [8](#)
- [42] Zhe Huang, Benjamin S Wessler, and Michael C Hughes. Detecting heart disease from multi-view ultrasound images via supervised attention multiple instance learning. In *Machine Learning for Healthcare Conference*, pages 285–307. PMLR, 2023. [5](#), [8](#)
- [43] Zhe Huang, Xiaowei Yu, Benjamin S Wessler, and Michael C Hughes. Semi-supervised multimodal multi-instance learning for aortic stenosis diagnosis. *arXiv preprint arXiv:2403.06024*, 2024. [8](#)
- [44] Zhe Huang, Xiaowei Yu, Dajiang Zhu, and Michael C Hughes. Interlude: Interactions between labeled and unlabeled data to enhance semi-supervised learning. *arXiv preprint arXiv:2403.10658*, 2024. [8](#)
- [45] Tongtong Huo, Yi Xie, Ying Fang, Ziyi Wang, Pengran Liu, Yuyu Duan, Jiayao Zhang, Honglin Wang, Mingdi Xue, Songxiang Liu, et al. Deep learning-based algorithm improves radiologists’ performance in lung cancer bone metastases detection on computed tomography. *Frontiers in Oncology*, 13:1125637, 2023. [5](#)
- [46] Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. Label propagation for deep semi-supervised learning.

- In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5070–5079, 2019. 21
- [47] Longlong Jing and Yingli Tian. Self-supervised visual feature learning with deep neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):4037–4058, 2020. 21
- [48] Jakob Nikolas Kather, Johannes Krisam, Pornpimol Charoentong, Tom Luedde, Esther Herpel, Cleo-Aron Weis, Timo Gaiser, Alexander Marx, Nektarios A Valous, Dyke Ferber, et al. Predicting survival from colorectal cancer histology slides using deep learning: A retrospective multicenter study. *PLoS medicine*, 16(1):e1002730, 2019. 15
- [49] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in neural information processing systems*, 33:18661–18673, 2020. 3
- [50] Byoungjip Kim, Jinho Choo, Yeong-Dae Kwon, Seongho Joe, Seungjai Min, and Youngjune Gwon. Selfmatch: Combining contrastive self-supervision and consistency for semi-supervised learning. *arXiv preprint arXiv:2101.06480*, 2021. 4
- [51] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [52] Durk P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised learning with deep generative models. *Advances in neural information processing systems*, 27, 2014. 21
- [53] Abhishek Kumar, Prasanna Sattigeri, and Tom Fletcher. Semi-supervised learning with gans: Manifold invariance with improved inference. *Advances in neural information processing systems*, 30, 2017. 21
- [54] Devidas T Kushnure, Shweta Tyagi, and Sanjay N Talbar. Lim-net: Lightweight multi-level multiscale network with deep residual learning for automatic liver segmentation in ct images. *Biomedical Signal Processing and Control*, 80: 104305, 2023. 5
- [55] Zhengfeng Lai, Chao Wang, Luca Cerny Oliveira, Brittany N Dugger, Sen-Ching Cheung, and Chen-Nee Chuah. Joint semi-supervised and active learning for segmentation of gigapixel pathology images with cost-effective labeling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 591–600, 2021.
- [56] Zhengfeng Lai, Chao Wang, Henry Gunawan, Sen-Ching S Cheung, and Chen-Nee Chuah. Smoothed adaptive weighting for imbalanced semi-supervised learning: Improve reliability against unknown distribution data. In *International Conference on Machine Learning*, pages 11828–11843. PMLR, 2022. 5
- [57] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. *arXiv preprint arXiv:1610.02242*, 2016. 21
- [58] Dong-Hyun Lee. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning at ICML*, 2013. 2, 3, 21
- [59] Junnan Li, Caiming Xiong, and Steven CH Hoi. Comatch: Semi-supervised learning with contrastive graph regularization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9475–9484, 2021. 3, 4
- [60] Ali Madani, Ramy Arnaout, Mohammad Mofrad, and Rima Arnaout. Fast and accurate view classification of echocardiograms using deep learning. *NPJ digital medicine*, 1(1):6, 2018. 5
- [61] Shaobo Min, Xuejin Chen, Hongtao Xie, Zheng-Jun Zha, and Yongdong Zhang. A mutually attentive co-training framework for semi-supervised recognition. *IEEE Transactions on Multimedia*, 23:899–910, 2020. 21
- [62] Avital Oliver, Augustus Odena, Colin A Raffel, Ekin Dogus Cubuk, and Ian Goodfellow. Realistic evaluation of deep semi-supervised learning algorithms. *Advances in neural information processing systems*, 31, 2018. 1, 2, 3, 4, 6, 8, 13, 22, 23
- [63] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 21
- [64] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011. 20
- [65] Guo-Jun Qi and Jiebo Luo. Small data challenges in big data era: A survey of recent progress on unsupervised and semi-supervised methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(4):2168–2187, 2020. 1, 4
- [66] Maithra Raghu, Chiyuan Zhang, Jon Kleinberg, and Samy Bengio. Transfusion: Understanding transfer learning for medical imaging. *Advances in neural information processing systems*, 32, 2019. 7
- [67] Kuniaki Saito, Donghyun Kim, and Kate Saenko. Openmatch: Open-set consistency regularization for semi-supervised learning with outliers. *arXiv preprint arXiv:2105.14148*, 2021. 8
- [68] Azizi Shekoofeh, Mustafa Basil, Ryan Fiona, Beaver Zachary, Freyberg Jan, Deaton Jonathan, Loh Aaron, Karthikesalingam Alan, Kornblith Simon, Chen Ting, et al. Big self-supervised models advance medical image classification. *arXiv preprint arXiv:2101.05224*, 2021. 23
- [69] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in neural information processing systems*, 33:596–608, 2020. 1, 3, 5, 21
- [70] Ewout W Steyerberg and Yvonne Vergouwe. Towards better clinical prediction models: seven steps for development and an abcd for validation. *European Heart Journal*, 35(29), 2014. 8
- [71] Jong-Chyi Su, Zezhou Cheng, and Subhransu Maji. A realistic evaluation of semi-supervised learning for fine-grained classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12966–12975, 2021. 1, 3, 4, 5, 7, 8, 13, 16, 23, 25

- [72] Tepei Suzuki. Consistency regularization for semi-supervised learning with pytorch. <https://github.com/perrying/pytorch-consistency-regularization>, 2020. 13
- [73] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 30, 2017. 3, 21
- [74] Jesper E Van Engelen and Holger H Hoos. A survey on semi-supervised learning. *Machine learning*, 109(2):373–440, 2020. 1, 21
- [75] Diane Wagner, Fabio Ferreira, Danny Stoll, Robin Tibor Schirrmeyer, Samuel Müller, and Frank Hutter. On the importance of hyperparameters and data augmentation for self-supervised learning. *arXiv preprint arXiv:2207.07875*, 2022. 1, 5
- [76] Yidong Wang, Hao Chen, Yue Fan, Wang Sun, Ran Tao, Wenxin Hou, Renjie Wang, Linyi Yang, Zhi Zhou, Lan-Zhe Guo, et al. Usb: A unified semi-supervised learning benchmark for classification. *Advances in Neural Information Processing Systems*, 35:3938–3961, 2022. 2, 3, 4, 6
- [77] Benjamin S Wessler, Zhe Huang, Gary M Long Jr, Stefano Pacifici, Nishant Prashar, Samuel Karmiy, Roman A Sandler, Joseph Z Sokol, Daniel B Sokol, Monica M Dehn, et al. Automated detection of aortic stenosis using machine learning. *Journal of the American Society of Echocardiography*, 36(4): 411–420, 2023. 5
- [78] Nan Wu, Jason Phang, Jungkyu Park, Yiqiu Shen, Zhe Huang, Masha Zorin, Stanisław Jastrzębski, Thibault Févry, Joe Katsnelson, Eric Kim, et al. Deep neural networks improve radiologists’ performance in breast cancer screening. *IEEE transactions on medical imaging*, 39(4):1184–1194, 2019. 5, 8
- [79] Jiancheng Yang, Rui Shi, and Bingbing Ni. Medmnist classification decathlon: A lightweight automl benchmark for medical image analysis. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pages 191–195. IEEE, 2021. 14
- [80] Jiancheng Yang, Rui Shi, Donglai Wei, Zequan Liu, Lin Zhao, Bilian Ke, Hanspeter Pfister, and Bingbing Ni. Medmnist v2-a large-scale lightweight benchmark for 2d and 3d biomedical image classification. *Scientific Data*, 10(1):41, 2023. 4, 14, 15
- [81] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2016. 5
- [82] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *International Conference on Machine Learning*, pages 12310–12320. PMLR, 2021. 3
- [83] Xiaohua Zhai, Avital Oliver, Alexander Kolesnikov, and Lucas Beyer. S4l: Self-supervised semi-supervised learning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1476–1485, 2019. 4
- [84] Bowen Zhang, Yidong Wang, Wenxin Hou, Hao Wu, Jindong Wang, Manabu Okumura, and Takahiro Shinozaki. Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling. *Advances in Neural Information Processing Systems*, 34:18408–18419, 2021. 3
- [85] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017. 3, 6, 21
- [86] Wenqiao Zhang, Lei Zhu, James Hallinan, Shengyu Zhang, Andrew Makmur, Qingpeng Cai, and Beng Chin Ooi. BoostMIS: Boosting medical image semi-supervised learning with adaptive pseudo labeling and informative active annotation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 23
- [87] Mingkai Zheng, Shan You, Lang Huang, Fei Wang, Chen Qian, and Chang Xu. Simmatch: Semi-supervised learning with similarity matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14471–14481, 2022. 4
- [88] Xiaojin Zhu. Semi-Supervised Learning Literature Survey. Technical Report 1530, Department of Computer Science, University of Wisconsin Madison., 2005. 1, 21
- [89] Xiaojin Zhu, Zoubin Ghahramani, and John D Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *Proceedings of the 20th International conference on Machine learning (ICML-03)*, pages 912–919, 2003. 21
- [90] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1):43–76, 2020. 21