

# *Semixup*: In- and Out-of-Manifold Regularization for Deep Semi-Supervised Knee Osteoarthritis Severity Grading from Plain Radiographs

Huy Hoang Nguyen, Simo Saarakkala, Matthew Blaschko, and Aleksei Tiulpin

**Abstract**—Knee osteoarthritis (OA) is one of the highest disability factors in the world in humans. This musculoskeletal disorder is assessed from clinical symptoms, and typically confirmed via radiographic assessment. This visual assessment done by a radiologist requires experience, and suffers from high inter-observer variability. The recent development in the literature has shown that deep learning (DL) methods can reliably perform the OA severity assessment according to the gold standard Kellgren-Lawrence (KL) grading system. However, these methods require large amounts of labeled data, which are costly to obtain. In this study, we propose the *Semixup* algorithm, a semi-supervised learning (SSL) approach to leverage unlabeled data. *Semixup* relies on consistency regularization using in- and out-of-manifold samples, together with interpolated consistency. On an independent test set, our method significantly outperformed other state-of-the-art SSL methods in most cases, and even achieved a comparable performance to a well-tuned fully supervised learning (SL) model that required over 12 times more labeled data.

**Index Terms**—Deep Learning, knee, osteoarthritis, semi-supervised learning.

## I. INTRODUCTION

**O**STEARTHRTIS (OA) is the most common joint disorder in the world causing enormous burdens at personal and societal levels [1]. OA has an unknown etiology, and its indications at late stages are worn cartilage, bone deformity, and synovitis [2]–[4].

The most common joints affected by OA are knee and hip, and among these, the disease is more prevalent in knee [5]–[8]. At the population level, such factors as sex, body-mass index (BMI), and age are known to be associated with OA [9]–[11]. As such, it was previously shown that people with BMI over 30 have 7-fold higher risk of knee OA than ones with BMI below 25 [12], and a half of elderly people over 65 years of age have OA in at least one joint [13].

From an economic perspective, OA leads to a huge burden in terms of direct costs (e.g. hospitalization, diagnosis, and

therapy), and indirect ones (e.g. losses of working days and productivity) [14]. For example, in the United States, OA costs hundreds of billion dollars annually, and is in the top-5 of annual Europe healthcare expenditure [3], [7].

Currently, knee OA diagnosis starts with a clinical examination, and then, a radiographic confirmation takes place when necessary [5], [15]. However, such guideline enables knee OA diagnosis only at a late stage when the cartilage is already worn, and the bone deformity is present, which leads to severe pain, and even physical disability [2], [3]. Ultimately, the only remaining option for a patient in that scenario is total knee replacement (TKR) surgery.

The literature shows a large and rapidly growing number of TKR surgeries worldwide [8], [16], [17]. As such, the annual rate of TKR surgeries in the United States has doubled since 2000 for adults of 45-64 years old [6], [18]. Therefore, there is a need for prevention of global disability.

Imaging, in contrast to clinical examination, may enable the detection of early knee OA signs at the stages when behavioral interventions (e.g. exercises and weight loss programs) could slow down the disease progression [19]. Radiographic assessment is the foremost imaging tool for detecting knee OA in primary care, and Kellgren-Lawrence (KL) is one of the most common clinical scales for the assessment of OA severity from plain radiographs (Figure 1). However, visual diagnosis done by a radiologist suffers from low inter-rater agreement [20], [21], thereby introducing large inconsistencies into decision-making. One possible solution to make OA diagnosis more systematic and allow for the detection of knee OA at early stages is to leverage computer-aided methods for image analysis [22]. Deep learning (DL) has become a state-of-the-art approach in this realm, and recent studies [21], [23], [24] have demonstrated that DL-based methods allow for fully-automatic KL grading. Furthermore, these studies showed a high level of agreement between the predictions made by DL-based models and the annotations produced by a consensus of radiologists.

Despite good and promising results, all the previously published DL methods in OA domain were based on Supervised-Learning (SL) and required large amounts of labeled data, which are not currently widely available. In practical applications, such datasets as the Osteoarthritis Initiative (OAI, <https://nda.nih.gov/oai/>) and the Multicenter Osteoarthritis Study (MOST, <http://most.ucsf.edu/>) are expensive to obtain due to high costs of data collection and annotation. As such, the latter needs multiple skilled experts (e.g. radiologists or orthope-

Huy Hoang Nguyen was with Research Unit of Medical Imaging, Physics and Technology, University of Oulu, Finland. E-mail: huy.nguyen@oulu.fi.

This work was supported in part by the strategic funding of the University of Oulu, in part by KAUTE foundation, Finland and in part by Sigrid Juselius Foundation, Finland.

Simo Saarakkala was with Research Unit of Medical Imaging, Physics and Technology, University of Oulu, Finland and Department of Diagnostic Radiology, Oulu University Hospital, Finland. E-mail: simo.saarakkala@oulu.fi

Matthew Blaschko was with Center for Processing Speech & Images, KU Leuven, Belgium. E-mail: matthew.blaschko@esat.kuleuven.be.

Aleksei Tiulpin was with Research Unit of Medical Imaging, Physics and Technology, University of Oulu, Finland and Department of Diagnostic Radiology, Oulu University Hospital, Finland. E-mail: aleksei.tiulpin@oulu.fi.

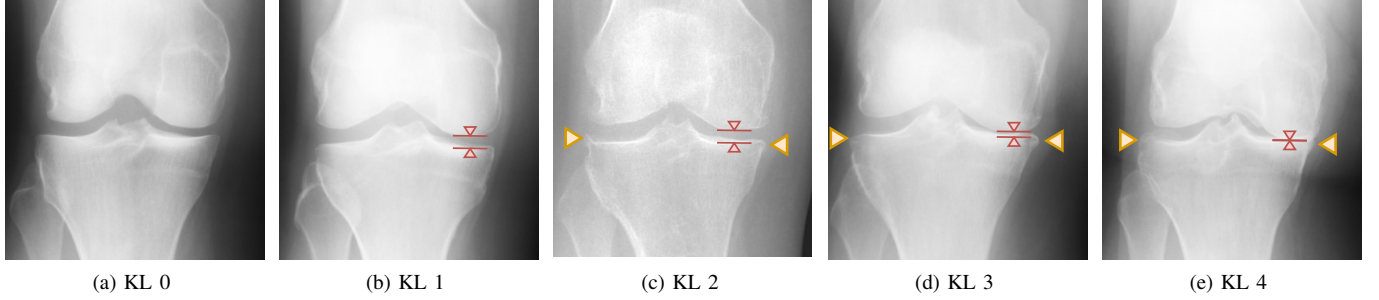


Fig. 1. Samples of knee radiographs. Joint space narrowing and osteophyte features are indicated by blue squares and red circles respectively. (a) KL 0: A healthy knee without OA, (b) KL 1 (Doubtful OA): Potential joint space narrowing, (c) KL 2 (Mild OA): Clear evidences of osteophytes, as well as slight reduction of joint space, (d) KL 3 (Moderate OA): Osteophytes grow and joint space narrowing progresses badly, and (e) KL 4 (Severe OA): Besides osteophytes, joint space is reduced so severely that the tibia and the femur are connected.

dists) making the process even more costly. Whereas labeled data are difficult to acquire, unlabeled data are available in large amounts, and can be collected from hospital imaging archives at low cost.

In the natural image recognition domain, it has been shown that leveraging small amounts of labeled, and large amounts of unlabeled data in a Semi-Supervised Learning (SSL) setting could potentially resolve the need for large labeled datasets. Recent studies [25]–[27] have developed SSL methods to utilize unlabeled data during the training processes, and have achieved competitive performances in image classification benchmarks using only a small fraction of the labeled data used in fully supervised settings.

Many SSL-based applications in the medical domain have previously been developed for automatic disease diagnosis. However, most of those are related to medical image segmentation [28]–[32], and use generative adversarial networks (GANs) [33] as a core method. To the best of our knowledge, there have been no SSL-based methods developed in the knee OA realm.

In this study, we, for the first time in the OA field, propose to leverage SSL for automatic assessment of knee OA severity from plain radiographs. Inspired by previous research in SSL [25], [26], [34], and the recently developed technique *mixup* [35], we propose a novel SSL method – *Semixup*, providing its systematic empirical comparison with the state-of-the-art approaches. Specifically, our contributions are the following:

- 1) We enhance the state-of-the-art supervised baseline [21], and propose a novel *Separable Adaptive Max-pooling* (SAM), as a drop-in replacement for the Global Average Pooling (GAP). This allows us to significantly improve over previously reported supervised results.
- 2) We introduce a novel semi-supervised DL-based method called *Semixup* for automatic KL grading of knee OA from plain radiographs. Our method yields competitive results to a well-tuned SL model trained on over 12 times more labeled data.
- 3) We systematically compare our method against several state-of-the-art SSL methods, and experimentally show that *Semixup* outperforms them in nearly all data regimes.

- 4) We follow the guidelines from [36] to conduct a realistic evaluation of our SSL approach, and provide insights into the scalability of our method with respect to the number of labeled examples, together with a tractable amount of unlabeled samples.

## II. RELATED WORK

### A. Deep Semi-Supervised Learning

There exists a wide variety of DL-based SSL methods in the literature; however, we discuss here only the ones that are close to our method and yield state-of-the-art results on generic image recognition datasets. Such approaches use two main ideas: *consistency regularization* and *pseudo-labeling*.

Consistency regularization is based on the assumption that the predictions of the model  $P(y|\mathbf{x})$  and  $P(y|T\mathbf{x})$ , where  $\mathbf{x}$  is a data point and  $T$  – class-preserving stochastic data augmentation – should not differ. The methods using the technique applied to unlabeled data include the  $\Pi$ -model [25], and Mean Teacher (MT) [37].

Label guessing (or pseudo-labeling) [38] was proposed in a deep learning setting in [39], and uses predicted labels for unlabeled samples with high confidences to update the gradients of the neural network. This technique can also be viewed as entropy regularization which favors low-density separation between classes [40].

The aforementioned SSL techniques have been explored separately; however, Berthelot *et al.* has recently introduced the MixMatch method [27], the idea of which was to combine label guessing and entropy regularization into a holistic framework. The authors of MixMatch made an empirical observation that applying *mixup* [35], which performed convex combinations of 2 arbitrary input data points  $x$  and  $x'$ , and their labels  $y$  and  $y'$  (i.e.  $\lambda x + (1-\lambda)x'$  and  $\lambda y + (1-\lambda)y'$ , where  $\lambda \sim \text{Beta}(\alpha, \alpha)$ , for  $\alpha \in (0, \infty)$ ), to labeled data and unlabeled data with guessed labels helps to improve the performance. In our method, we also use *mixup* for both labeled and unlabeled data; however, we avoid using label guessing due to its potential to propagate label errors that can be common in medical domains.

We finally discuss here two main limitations of all the mentioned recent studies on SSL. The first common limitation

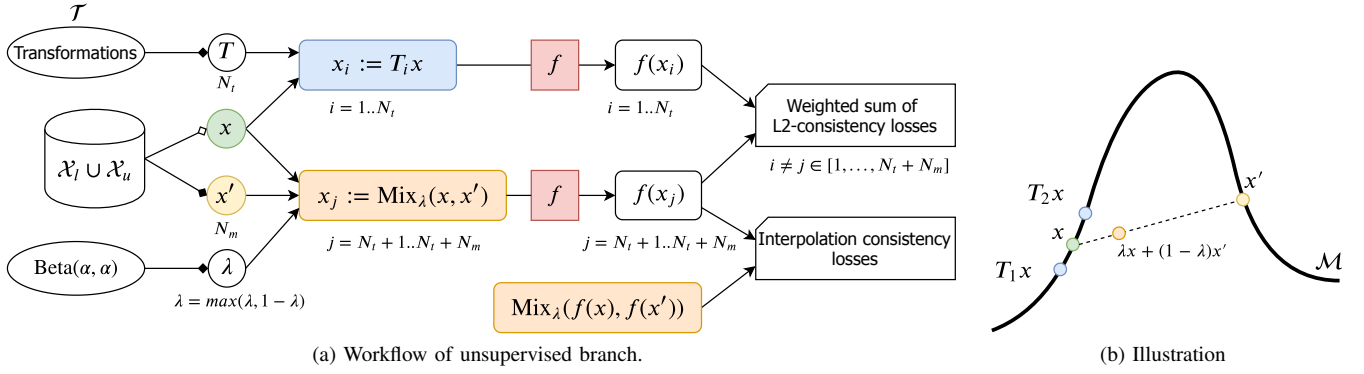


Fig. 2. The main idea of *Semixup*: (a) The workflow of its in- and out-of-manifold consistency regularization branch, and (b) its conceptual illustration. For each image  $x \in L \cup U$ , we sample  $N_l$  transformations,  $N_m$  images  $x'$ , and *mixup* coefficients  $\lambda$  (s.t.  $\lambda > 0.5$ ) to transform and blend with  $x$ . Then, we enforce consistency of predictions for every pair  $x_i, x_j$  under L2 norm. In addition, we use the interpolation consistency regularizer. Links with ■ denote sampling with replacement, ones with □ indicate sampling without replacement.

of those methods is that their evaluations were done on major image recognition benchmarks, namely CIFAR-10 and/or ImageNet without considering real-world problems, such as the ones related to medical image recognition.

The second limitation of the aforementioned studies is that only the number of labeled samples was varied in the experiments; however, none of them explored the amount of unlabeled data needed to achieve good performance. We emphasize the importance of this issue in medical imaging as from one application to another, the amounts of both labeled and unlabeled data can differ. In our work, we tackled this issue by varying these hyperparameters.

### B. Deep Learning for Knee Osteoarthritis Diagnosis

Fully supervised DL-based methods that use Convolutional Neural Networks (CNN) have recently been used to assess knee OA severity. In particular, in their pioneering work, Antony *et al.* [41] applied transfer learning and showed significant improvements compared to classifiers trained using hand-crafted features. In the follow-up work [23], Antony *et al.* proposed a CNN architecture trained from scratch to classify knee images according to the KL scale. Their new approach outperformed their previous transfer learning results.

The main common limitation of the aforementioned studies was that neither of them utilized an independent test set. In contrast, Tiulpin *et al.* [21] addressed this limitation and also proposed a novel CNN architecture that outperformed the previous methods [23], [41]. Interestingly, that model performed on-par with transfer learning baseline while having significantly less trainable parameters thanks to the inductive bias from using the relative symmetry of visual features in knee images.

The recent studies by Chen *et al.* [42], Norman *et al.* [43], and Górriz *et al.* [44] did not use any independent test set either. The only latest study where an independent test set was used for assessment of the results was by Tiulpin *et al.* [45]. In that study, the authors obtained a KL classification model as a bi-product of their main method; however, they obtained results that are similar to their previous study [21].

Despite having an independent test set and state-of-the-art results, Tiulpin *et al.* [21], [45] used large amounts of labeled

data for training. We emphasize here that none of the existing studies in which DL was applied for knee OA diagnosis from radiographs addressed the question of the quality of automatic KL grading as a function of the dataset size. Our work answers this question via a thorough experimental evaluation of both SL and SSL methods.

## III. METHOD

### A. Overview

The method proposed in this paper consists of two parts: 1) a novel extension of a previously developed Siamese network developed by Tiulpin *et al.* [21] and 2) a novel deep SSL technique. The utilized Siamese model uses shared branches of the CNN which focus their attention on the medial and the lateral sides of the analyzed knee (see Figure 3). Here, we consider pairs of lateral and medial image patches as single data points  $x \in \mathbb{R}^{2 \times H \times H}$ , where  $H$  is the size of the image patch. KL grades are the outputs of our model:  $y \in \{0, 1, 2, 3, 4\}$ .  $f_\theta$  denotes a Siamese neural network with parameters  $\theta$ . In our setting,  $p(y|x) = f_\theta(x)$ .

Our SSL method aims to perform penalization of local sharpness of the surface loss along the data manifold  $\mathcal{M}$  and also within its surroundings. The former is achieved via minimization of  $\mathbb{E}_x \|J_{\mathcal{M}}\|_F^2$ , where  $J_{\mathcal{M}}$  denotes the Jacobian along the data manifold  $\mathcal{M}$  and  $\|\cdot\|_F$  is the Frobenius norm, as well as  $\mathbb{E}_x \|J_\theta\|_F^2$ , where  $J_\theta$  denotes the Jacobian in the parameter space, using consistency regularization [25], [46]. To generate out-of-manifold samples, we use *mixup* [35]. Here, we first enforce linear behavior of the model along the *mixup* rays via interpolation consistency training (ICT) [26], and then apply the aforementioned consistency regularization to enforce consistent behavior of the model along the *mixup* rays which are in the close surroundings of  $\mathcal{M}$ . As our idea for SSL centers around *mixup*, we name our method *Semixup*. We graphically illustrate the process of sample generation for consistency regularization in Figure 2.

### B. Network Architecture

We follow the design of the previously developed Siamese model by Tiulpin *et al.* [21], and propose several modifications

essential to obtain better KL grading performance in both SL and SSL settings. The schematic illustration of our proposed network is presented in Figure 3.

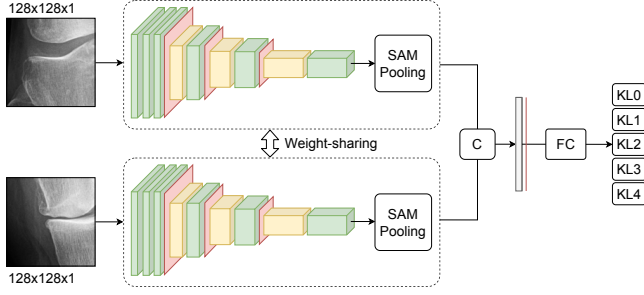


Fig. 3. The Siamese architecture of our model. Green denotes the blocks that use  $3 \times 3$  convolutions with the stride of 1, and yellow denotes  $3 \times 3$  convolutions with the stride of 2. Red indicates dropout layers. After pooling the features by Separable Adaptive Max-pooling (SAM) and concatenating (C) them into a single vector, they are passed into a fully-connected (FC) layer that predicts KL grades.

The basic building block of our model consists of a  $3 \times 3$  convolution with a zero padding  $P$  and a stride  $S$ , an instance normalization (IN), and a leaky rectified linear unit (LeakyReLU) activation with a slope of 0.2.

To enrich the representation power of our model, we first start with three consecutive  $3 \times 3$ ,  $S = 1$ ,  $P = 1$ , and one  $3 \times 3$ ,  $S = 2$ ,  $P = 1$  convolutional blocks (green blocks in Figure 3). Subsequently, we alternate  $3 \times 3$ ,  $S = 2$ ,  $P = 1$  (yellow blocks in Figure 3), and  $3 \times 3$ ,  $S = 1$ ,  $P = 1$  blocks until a feature map of size  $\frac{1}{8}H \times \frac{1}{8}H$  is obtained.

We use stridden convolutions to perform the downsampling, whereas the original model from [21] used max-pooling layers. The main motivation for our method to use stridden convolutions is that the translation invariance achieved by the use of max-pooling could potentially harm the results by removing the dependencies between the OA-related fine-grained features at higher layers of the network [47].

Similar issues could also arise in the bottleneck of the network, where Tiulpin *et al.* [21] used Global Average Pooling (GAP). We argue that averaging such large feature maps is not the most optimal pooling strategy, and we tackle this problem via our novel SAM pooling (see Section III-C).

All the blocks described above share the weights among the branches of our Siamese CNN, where each branch of the model processes an individual side of the knee image pair (lateral or medial). Each branch of this model generates  $1 \times 1$  features, which are concatenated, passed through a dropout, and subsequently fed into a fully-connected layer the OA severity stage. A detailed description of our architecture is provided in Supplementary Table S7.

### C. Separable Adaptive Max-pooling

As mentioned previously, we propose a replacement for GAP to deal with the potential information loss in the bottleneck of the model. The proposed SAM scheme is based on the idea of firstly applying pooling along one direction of the feature map (horizontal or vertical). Secondly, we use a  $1 \times 1$  convolutional block (with IN and LeakyReLU) to remove

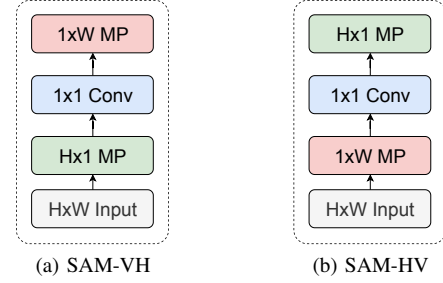


Fig. 4. Separable Adaptive Max-pooling configurations. The core idea of this approach is to inject a non-linearity between the pooling steps – vertical-horizontal and horizontal-vertical as displayed in subplots (a) and (b), respectively. Here, H and W indicate the height and the width of the input, respectively.

unnecessary correlations between the pooled features. Finally, we apply the pooling in the other direction than was done in the first phase.

As the initial pooling of the features can be done in either horizontal or vertical directions, we suggest two configurations of SAM presented in Figure 4, namely:

- 1) Max-pooling  $\frac{1}{8}H \times 1$  with stride  $\frac{1}{8}H \times 1$ ,  $1 \times 1$  convolution, IN, LeakyReLU with the slope of 0.2, and  $1 \times \frac{1}{8}H$  max-pooling with the stride of  $1 \times 1$  (SAM-VH).
- 2) Similar to the above, but the first and the last max-poolings are swapped (SAM-HV).

### D. Semi-Supervised Learning

1) *Problem setting*: Let  $X_l$  and  $X_u$  be labeled and unlabeled image sets, respectively. Let  $\mathcal{Y}$  denote the labels for  $X_l$ . In our setting, we optimize the following objective:

$$\min_{\theta} \mathcal{L}_l(\theta; X_l, \mathcal{Y}) + \mathcal{L}_u(\theta; \mathbf{w}, X_u, X_l), \quad (1)$$

where  $\mathcal{L}_l$  is a cross-entropy loss with *mixup*,  $\mathcal{L}_u$  is a combination of losses without the involvement of labels, and  $\mathbf{w}$  are hyperparameters responsible for weighing unsupervised losses. Here,  $\mathcal{L}_u$  acts as a regularizer which leverages data without labels to enhance the robustness of the model  $f_{\theta}$  via auxiliary tasks.

2) *Supervised Loss*: *mixup* proposed in [35] is a simple and effective technique to improve generalization. It can be viewed as data augmentation, and in a nutshell, it linearly mixes two samples  $x_i, x_j$  with a blending coefficient  $\lambda \sim \text{Beta}(\alpha, \alpha)$ , for  $\alpha \in (0, \infty)$ :

$$x_{mix} = \text{Mix}_{\lambda}(x_i, x_j) = \lambda x_i + (1 - \lambda)x_j. \quad (2)$$

Having the mixed sample  $x_{mix}$ , the following loss is optimized:

$$\mathcal{L}_l(\theta; x_{mix}, y_i, y_j) = \lambda \mathcal{L}_{ce}(x_{mix}, y_i) + (1 - \lambda) \mathcal{L}_{ce}(x_{mix}, y_j), \quad (3)$$

where  $\mathcal{L}_{ce}$  is a multi-class cross-entropy loss. Here and further, we call the loss in (3) "soft" cross-entropy loss.

After the introduction of *mixup*, it was shown that this technique performs out-of-manifold regularization [48], and



in our work, we exploit this property. Specifically, we view *mixup* as a data augmentation technique that generates out-of-manifold samples  $x_{mix}$  that belong to the ray between any two points  $x_i$  and  $x_j$  (Equation (2), and Figure 2b).

3) *Consistency Regularization: Penalizing Local Loss Sharpness*: The consistency regularization technique used in [25], [37] aims to minimize the following objective w.r.t  $\theta$ :

$$\mathbb{E}_{x \sim p(x)} \mathbb{E}_{T, T' \sim p(\tau)} \|f_\theta(Tx) - f_\theta(T'x)\|_2^2, \quad (4)$$

where  $p(x)$  is the distribution of training data (labeled and unlabeled), and  $p(\tau)$  is the distribution of stochastic transformations (e.g., data augmentations).

According to Athiwaratkun *et al.* [46], regularizing consistency for a model using dropout in convolutional layers implies minimization of two terms:  $\mathbb{E}_x \|J_M\|_F^2$  and  $\mathbb{E}_x \|J_\theta\|_F^2$ . That idea results in an interesting connection between the consistency-based method and the classic graph-based approaches that use Laplacian regularization for SSL [49].

The explanations provided by Athiwaratkun *et al.* [46] demonstrate that minimization of  $\mathbb{E}_x \|J_\theta\|_F^2$  leads to a broader optimum that is presumably helpful for good model generalization. Therefore, the minimization of the regularization term from (4) can lead to better performance which has been supported by experimental evidence in [25], [37], [46].

4) *Interpolation Consistency*: The *mixup* operator was introduced as an efficient data augmentation for supervised learning regularizing against out-of-manifold samples that lie close to  $\mathcal{M}$ . Recently, Verma *et al.* [26] utilized it to enforce linear behavior of the model along the *mixup* rays in the ICT method:

$$\mathbb{E}_{\lambda \sim \text{Beta}(\alpha, \alpha)} \mathbb{E}_{x_i, x_j \sim p(x)} \|\text{Mix}_\lambda(f_\theta(x_i), f_\theta(x_j)) - f_\theta(x_{mix})\|_2^2. \quad (5)$$

While Verma *et al.* [26] did not consider *mixup* as an out-of-manifold regularizer, we consider this to be an important observation [48].

### E. Semixup

1) *Motivation*: The consistency regularization and ICT aim to maximize the consistency of label assignment within and out of  $\mathcal{M}$ . However, we also note that while linear behavior of the model is achieved in ICT, it does not aim to minimize the inconsistency of label assignment for e.g.  $Tx$  and  $x_{mix}$  which can be viewed as in- and out-of-manifold augmented versions of a data point  $x$ . In this work, we address this limitation and strengthen the ability of making consistent predictions for out-of-manifold samples that are close to the data manifold  $\mathcal{M}$ . In summary, besides applying loss (3) over labeled data, *Semixup* optimizes a linear combination of the objectives shown in (4), and (5), as well as the out-of-manifold term described in the following section over both labeled and unlabeled samples. Supplementary Algorithm 1 shows a concrete implementation of our method.

2) *Out-of-Manifold Consistency Regularization*: We use the aforementioned regularizers from (4) and (5). Subsequently,

given the motivation above, we minimize the following additional objective:

$$\mathbb{E}_{\lambda \sim \text{Beta}(\alpha, \alpha)} \mathbb{E}_{x, x' \sim p(x)} \mathbb{E}_{T \sim p(\tau)} \|f_\theta(\text{Mix}_\lambda(x, x')) - f_\theta(Tx)\|_2^2. \quad (6)$$

The presented objective aims to maximize the consistent label assignment for perturbed data items  $Tx \in \mathcal{M}$ , and also the ones being out-of-manifold but are close to it. To generate the latter, when sampling  $\lambda$  in *mixup*, we enforce  $\lambda = \max(\lambda, 1 - \lambda)$ .

3) *Low-variance Sampling*: Our method uses both labeled and unlabeled data in the unsupervised regularization term of the loss (1). It is motivated by the fact that unlabeled data in medical imaging can come from different device vendors rather than labeled data, thereby the empirical distributions of labeled and unlabeled data might be misaligned.

We note that the objective in (6) uses stochastic augmentations  $T \sim p(\tau)$ . When using this loss in a combination with consistency regularization and ICT, it can be seen that there exist predictions for two versions of an image  $x$  as in (4). Therefore, when optimizing with stochastic gradient descent, we are able to obtain a lower variance stochastic estimate of the objective (6) at low marginal computational cost (see lines 10, 11 of Supplementary Algorithm 1).

## IV. EXPERIMENTS

### A. Datasets

We used knee radiographs from two large public cohorts: The OAI and the MOST. The OAI dataset was collected from 4,796 participants whose ages were from 45 to 79 years old. The cohort included a baseline, and follow-up visits after 12, 18, 24, 30, 36, 48, 60, 72, 84, 96, 108, 120, and 132 months. Radiographic imaging for bilateral fixed flexion X-ray images took place at most, but not all follow-ups. We used all available images except for those imaged during 18, 30, 60, 84, and 108-month follow-ups, where imaging was performed only for small sub-cohorts. The MOST dataset had 3,021 participants examined at its baseline, and follow-up visits after 15, 30, 60, 72, and 84 months. Except for the 72-month follow-up, each examination included radiographic imaging.

Radiographs in the OAI dataset were posterior-anterior (PA) bilateral images acquired with the protocol that called for a beam angle of 10 degrees. We utilized OAI data in training and model selection phases.

The MOST dataset included PA, and lateral images. The PA images were acquired with a beam angle of 5, 10, or 15 degrees. In this study, we used the MOST dataset as an independent test set. To ensure the reliability of the labels in MOST, we only keep 10 degree bilateral PA images that were acquired at baseline, had KL and Osteoarthritis Research Society International (OARSI) grades available, and were without implants. Eventually, our training and test data comprised 39,902 and 3,445 knee images from the OAI, and the MOST datasets, respectively. The detailed data distribution by KL grade is presented in Table I.

TABLE I  
DESCRIPTION OF THE OAI AND THE MOST DATASETS.

Dataset	Split	# Images	KL grade				
			0	1	2	3	4
OAI	Train/Val	39,902	15,954	7,636	9,617	5,228	1,467
MOST	Test	3,445	1,550	568	520	559	248

## B. Experimental Setup

1) *Data Pre-Processing*: The essential step in analyzing knee radiographs is a region of interest (ROI) localization. In this paper, we focused rather on the development of an efficient DL architecture and SSL method than on implementing a full knee X-ray analysis pipeline. Thus, we utilized a Random Forest Regression Voting Constrained Local Model approach implemented in the BoneFinder (BF) tool [50].

In each bilateral X-ray, we localized the anatomical landmarks using BF and cropped the ROIs of  $140\text{mm} \times 140\text{mm}$  centered at each knee joint (2 ROIs per image at most). Subsequently, based on the anatomical landmarks, we performed the standardization of each of the ROIs by a horizontal alignment of the tibial plateau.

To standardize the intensity of the histograms, we performed intensity clipping to the 5<sup>th</sup> and 99<sup>th</sup> percentiles as well as global contrast normalization prior to converting 16-bit raw DICOM images into an 8-bit intensity range. We then center-cropped the obtained 8-bit images to  $110\text{mm} \times 110\text{mm}$  and resized them to  $300 \times 300$  pixels (0.37mm pixel spacing).

At the next step of the pre-processing, similar to [21], we flipped the left images to look like the right ones, and cropped the medial and the lateral patches with a size of  $128 \times 128$  pixels ( $H = 128$ ). These patches were cropped with the common top anchor at one third of the image height. Both patches were cropped at the lateral and medial image sides. After cropping, we flipped the medial patch horizontally (Figure 3). Finally, we normalized the intensities of each item in the obtained pair into the intensity range of  $[-1, 1]$ . Whereas the previous studies [21], [41] required the statistics of their training sets to normalize input data, we utilized the mean of 0.5 and the standard deviation of 0.5.

### 2) Training and Evaluation Protocol:

a) *Data Split*: First, we split the OAI dataset into 2 parts so that 25% and 75% were for labeled and unlabeled data, respectively. Here, we stratified the splits by KL grade, and ensured that patient IDs do not belong to both of these parts. We then applied another stratification as above to divide each of the aforementioned splits into 5 folds. In each fold, we generated 24 data settings, having 4 labeled (50, 100, 500, and 1000 samples per KL grade), and 6 unlabeled data configurations. Here, the amount of unlabeled data was  $N$ ,  $2N$ ,  $3N$ ,  $4N$ ,  $5N$ , or  $6N$  samples, where  $N$  is the corresponding amount of labeled data.

In addition, we also used the whole OAI dataset to assess the performance of our SL baseline. The training and the validation sets of the setting had 31,922 and 7,980 samples, respectively. Therefore, the 4 aforementioned labeled data settings varied from 0.8% to 16% of the whole training set.

b) *Architecture Selection and Training Setup*: Following the protocol proposed by Oliver *et al.* [36], we firstly tuned the SL setting before comparing it to SSL methods. As such, we considered 6 feasible architectures, each of which was the combination of either the architecture of the SL baseline [21], or ours with each type of pooling layer (GAP, SAM-VH, or SAM-HV). In our experiments, we selected the top-3 architectures for further comparisons based on cross-validation. The best model among those was utilized as the base model of *Semixup* and SSL baselines.

c) *SL and SSL Comparisons*: We investigated effects of SSL methods such as the  $\Pi$ -model [25], MixMatch [27], and *Semixup* without the use of unlabeled data. In such scenario, MixMatch [27] is equivalent to *mixup* [35]. Ultimately, we compared *Semixup* to 3 SSL baselines such as the  $\Pi$  model [25], Interpolation Consistency Training (ICT) [26], and MixMatch [27]. Each method was trained on the 24 data settings of the first fold previously generated from the OAI, and finally evaluated on the independent test set from the MOST. Detailed implementations of *Semixup* and all the SSL baselines are provided in Supplementary Section S5-A and Section S5-B.

TABLE II  
ABLATION STUDY OF THE SAM POOLING (SL SETTING). THE VALUES INDICATE BALANCED ACCURACIES (%). THE RESULTS WERE ESTIMATED OUT-OF-FOLD USING A 5-FOLD CROSS-VALIDATION.

Base model	Pooling	# data / KL grade				
		50	100	500	1000	Average
Tiulpin <i>et al.</i> [21]	GAP	43.7	52.5	62.1	66.8	56.3
	SAM-VH	<b>54.0</b>	50.7	62.5	<u>69.1</u>	<u>59.1</u>
	SAM-HV	39.0	51.2	63.7	62.7	54.1
Ours (Sec. III-B)	GAP	46.7	49.6	<b>67.0</b>	68.6	58.0
	SAM-VH	47.7	<u>54.6</u>	65.8	65.5	58.4
	SAM-HV	<u>48.3</u>	<b>56.3</b>	<u>66.7</u>	<b>69.7</b>	<b>60.3</b>

d) *Metrics*: In the training phase of all the experiments, the best models were selected based on both balanced multi-class accuracy (BA) [51] and Cohen’s quadratic kappa coefficient (KC) [52]. In the final model evaluation and comparison to the baseline methods, we used BA. To be in line with the metrics used in the previous studies, such as [21], [24], [41], we also used confusion matrix, KC, and mean square error (MSE). To assess the performance of detecting radiographic OA ( $\text{KL} \geq 2$ ), we used receiver operating characteristic (ROC) curves, area under the ROC curve (AUC), precision-recall (PR) curves, and average precision (AP).

e) *Statistical Analyses*: In our initial experiments, we noticed that several factors such as data acquisition center, knee side (left or right), and subject ID may affect the validity of the statistical analyses. To verify this, we used a generalized linear mixed effects model [53], and noticed that neither of these factors has an impact on the results. Therefore, for simplicity of evaluation, we used standard error and one-sided Wilcoxon signed-rank test [54]. Here, we split the test into 20 equal-sized chunks (no overlapping patients), and calculated

the BA on each of them. Finally, these values were used for the statistical analyses.

### C. Supervised Learning

1) *Ablation Study*: Table II shows that, on average, the models with our base architecture performed 4.1% better compared to the models with [21] in the case of 500 samples per KL grade. With respect to the base architecture of [21], SAM-VH was the most suitable pooling operator, especially in the cases of 50 and 1000 samples per KL grade. On the other hand, our base architecture with SAM-HV yielded the highest average BA.

Based on the average BAs in Table II, we selected the three best architectures for further comparisons. The best one among those (ours with SAM-HV) was chosen as the base model of all the SSL methods.

2) *Performance on the Test Set*: In the first group of results, Table III shows that our architecture with SAM outperformed the SL baseline model [21] in all the data settings. Specifically, our architecture with SAM-HV yielded BAs 9% better than the baseline [21] in the data settings of 500 and 1000 images per KL grade. Notably, in the case of 500 samples per KL grade, our model with SAM-HV surpassed by 6% the baseline model that would require a double amount of data. Finally, our model with SAM-HV achieved a BA of 67.5% when trained on only 16% of the full OAI set.

TABLE III

SL AND SSL METHODS EVALUATION ON THE TEST SET. WE REPORTED THE TOP-3 SL MODELS WITH GAP<sup>\*</sup>, SAM-VH<sup>†</sup>, OR SAM-HV<sup>‡</sup>. OUR MODEL WITH SAM-HV<sup>‡</sup> IS THE COMMON ARCHITECTURE FOR ALL SSL METHODS. BOLD HIGHLIGHTS MODELS PERFORMING SUBSTANTIALLY BETTER THAN THE SECOND BEST MODEL IN EACH CATEGORY. ◊ INDICATES NO SUBSTANTIAL DIFFERENCE AMONG THE BEST MODELS.

Method	# labels / KL grade			
	50	100	500	1000
Fully SL				
Tiulpin <i>et al.</i> <sup>*</sup> [21]	40.5±0.8	49.7±0.9	55.1±0.8	58.5±0.8
Tiulpin <i>et al.</i> <sup>†</sup> [21]	45.2±0.8◊	48.6±0.9	57.6±0.8	58.5±0.8
Our SL <sup>†</sup> (Sec. III-B)	45.6±0.8◊	53.2±0.9◊	61.5±0.8	63.5±0.8
Our SL <sup>‡</sup> (Sec. III-B)	41.5±0.8	52.9±0.9◊	<b>64.0±0.8</b>	<b>67.5±0.8</b>
SL and SSL methods - Without unlabeled data				
<i>mixup</i> [35]	39.6±0.8	53.9±0.8	64.5±0.8	67.1±0.8◊
Π model [25]	<b>42.3±0.8</b>	56.3±0.8◊	64.3±0.8	67.9±0.8◊
<i>Semixup</i> (Ours)	31.7±0.8	56.1±0.8◊	<b>66.6±0.8</b>	68.0±0.8◊
SSL methods - Best performing models <sup>*</sup> comparison				
ICT [26]	47.7±0.9◊	53.5±0.8	64.1±0.8	67.8±0.8
Π model [25]	45.9±0.8◊	56.2±0.8	67.1±0.8	69.0±0.8
MixMatch [27]	45.1±0.8◊	56.0±0.8	67.6±0.8	68.4±0.8
<i>Semixup</i> (Ours)	46.9±0.9◊	<b>58.8±0.8</b>	<b>69.7±0.8</b>	<b>71.0±0.8</b>
Fully SL - Trained on the full OAI (31,922 samples)				
Our SL <sup>‡</sup> (Sec. III-B)	70.9±0.8			

### D. Semi-Supervised Learning

We summarized the results of SL and SSL models trained on 4 labeled data settings, and the SL model trained on full OAI data in Table III. Detailed performance evaluations of all the SSL methods are described in Supplementary Table S8.

1) *SSL Methods without Unlabeled Data*: The performance of *Semixup* scaled well with the amount of training data, and outperformed its corresponding fully SL model when we had at least 100 samples per KL grade. In the comparison with other SSL baselines, *Semixup* achieved comparable results.

2) *Impact of Consistency Regularization Terms*: Because our in- and out-of-manifold regularizer comprises several individual losses, it is essential to understand the impact of each of them onto the method's performance. To assess the contributions of each used regularizer – (4), (5), and (6), we consecutively removed each of them from our loss function. Table IV shows this ablation study. Here, we report the best validation performance among 6 different unlabeled data settings. The results show that *Semixup* performs better when all the regularizers are used jointly.

TABLE IV  
ABLATION STUDY OF REGULARIZATION TERMS USED IN *Semixup*.  
THE VALUES ARE BAs (%).

Method	# data / KL grade	
	100	500
<i>Semixup</i> w/o in-manifold terms	60.8	73.1
<i>Semixup</i> w/o out-of-manifold term	60.8	72.0
<i>Semixup</i> w/o interpolation consistency	60.8	73.4
<i>Semixup</i>	<b>65.6</b>	<b>74.1</b>

3) *Semixup Compared to SL and SSL Baselines*: After evaluating 24 trained models of each SSL method, we derived maximum BAs grouped by labeled data setting as in the third part of Table III. With 500 labeled samples per KL grade, *Semixup* achieved a BA of 69.7%, which was 5.7% higher than the result of the SL baseline, and surpassed the SL baseline by 2.2% even when trained on 2 times less labeled data than the baseline was. With 1000 labeled data per KL grade, our SSL method reached the BA of 71% that was comparable to the result of the SL baseline trained on the full OAI setting (over 6 times more data used).

TABLE V  
COMPARISON OF OUR BEST MODELS TRAINED WITH *Semixup* AGAINST  
OUR WELL-TUNED SL MODEL WITH SAM-HV.

Method	# labels	KC	MSE	AUC (KL ≥ 2)	AP (KL ≥ 2)
<i>Semixup</i>	250	0.708	1.080	0.880	0.861
		[0.692, 0.724]	[1.000, 1.022]	[0.869, 0.891]	[0.849, 0.872]
	500	0.789	0.810	0.933	0.916
		[0.776, 0.801]	[0.764, 0.860]	[0.926, 0.940]	[0.906, 0.925]
	2,500	0.877	0.442	0.967	0.956
		[0.870, 0.884]	[0.416, 0.466]	[0.962, 0.972]	[0.950, 0.962]
SL (full OAI)	31,922	0.878	0.458	0.972	0.959
		[0.870, 0.885]	[0.430, 0.487]	[0.968, 0.976]	[0.953, 0.965]
SL (full OAI)	31,922	0.881	0.440	0.974	0.963
		[0.873, 0.889]	[0.414, 0.471]	[0.969, 0.978]	[0.958, 0.969]

Besides comparing to the SL settings, we also compared *Semixup* to the state-of-the-art SSL baselines. On the test set,

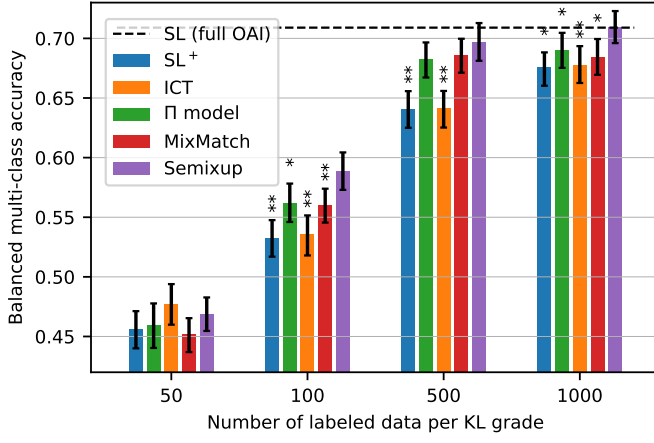


Fig. 5. Graphical comparison of *Semixup* and the baseline methods (MOST dataset). The bars indicate the 95% confidence intervals,  $SL^+$  indicates our fully SL models equipped with either SAM-HV or SAM-VH, and \* and \*\* indicate the statistically significant difference ( $p < 0.05$  and  $p < 0.001$ , respectively). The dash line indicates the BA of the fully SL model with SAM-HV trained on the full OAI dataset.

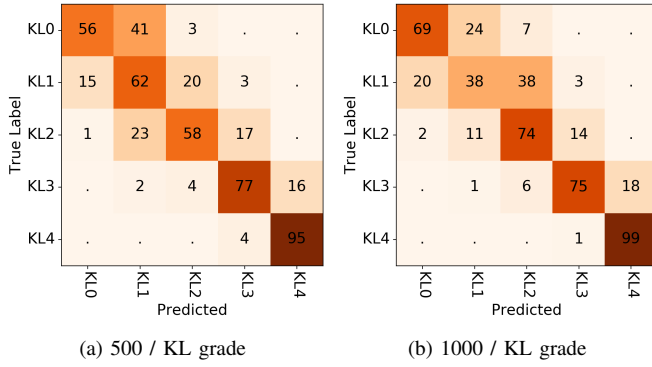


Fig. 6. Confusion matrices showing performance (%) of our models trained with *Semixup* on the test set (MOST dataset) in two labeled data settings (a) and (b).

our method outperformed the others in the data settings of 100, 500, and 1000 labels per KL grade. In the data setting of 2500 labeled images (7.8% of the full OAI size), *Semixup* yielded a BA 2.1% higher than the best SSL baseline, MixMatch, with the same labeled data amount, and 0.7% higher than the best SSL baseline, the  $\Pi$  model, with twice the amount of labeled data. In addition, *Semixup* achieved the highest average BAs in the data settings having at least 100 labeled images per KL grade (Supplementary Table S8).

The comparisons between *Semixup* and each of the baseline models across different labeled data settings are graphically illustrated in Figure 5. The best models of *Semixup* in the cases of 100 and 1000 labels per KL grade significantly outperformed all the baselines ( $p < 0.05$ ). In the case of 500 labeled samples per KL grade, *Semixup* was significantly better than the SL baseline and ICT with the same amount of labels. Furthermore, our method also surpassed the SL baseline that was trained with twice more labeled data (1000 samples per KL grade). The statistical analyses indicate that SL model

trained on the full OAI dataset did not differ significantly from the *Semixup* models that were trained with 500 and 1000 labeled samples per KL grade ( $p$ -values were of 0.054 and 0.368, respectively). The details of the statistical testing are presented in Supplementary Table S9.

4) *Detection of Early Radiographic Osteoarthritis*: Figure 6 presents the confusion matrices of our 2 best models. Here, we evaluate the accuracy of our method to detect early OA (i.e., KL = 2). With 500 and 1000 labels per KL grade, *Semixup* was able to detect early OA with 58% and 74% accuracies. Notably, without the doubtful OA (KL = 1), we achieved a substantially high BA of 79.25%.

With regard to detection of radiographic OA (KL  $\geq 2$ ), Figure 7 and Table V show how our SSL model with 500 labels per KL grade was comparable to the well-tuned SL model trained on the large training set (more than 12 times labeled samples). The detailed comparisons with respect to different amounts of unlabeled data are presented in Supplementary Figure S8 and Figure S9.

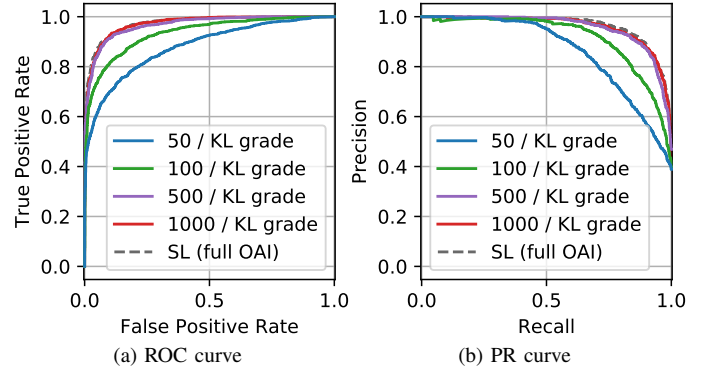


Fig. 7. Comparison of the best models of trained with *Semixup* on radiographic OA detection task (KL  $\geq 2$ ). The models were trained with 50, 100, 500, or 1000 labeled examples per KL grade.

## V. DISCUSSION AND CONCLUSIONS

In this study, we presented a novel SSL method – *Semixup*. This method leverages in- and out-of-manifold consistency regularization, and we demonstrated its application in the task of automatic grading of knee OA severity. Furthermore, we also proposed a novel state-of-the-art architecture for this task.

The core novelty of this work lies in using *mixup* for generating out-of-manifold samples in close proximity of data manifold, and then ensuring consistent predictions made by a neural network on such samples and the ones drawn from the data manifold itself. We experimentally showed that the proposed technique helps to train more robust models in limited data regimes even when unlabeled data is not available. Another novelty of this work is the proposed pooling approach that allows to efficiently process large feature maps, thereby leveraging fine-grained information required for knee OA grading. Our work demonstrated not only the new method for SSL, but also an efficient baseline for supervised setting when the amount of training samples is limited. To our knowledge, this is the first work that systematically studied the problem



of learning robust classifiers for knee OA severity assessment in the limited data regime. Our implementation is publicly available at <https://github.com/MIPT-Oulu/semixup>.

Besides all the aforementioned benefits, this study has several limitations. Specifically, we did not use EMA in our methods as it is costly to train and the methods using EMA would require significantly more computational power than their competitors. Another limitation of this study is that we did not leverage the power of SAM-VH and SAM-HV in our model. We foresee potential improvements of our results if both of these pooling schemes are used, e.g. by ensembling the results produced by models trained with SAM-VH and SAM-HV, respectively. From a clinical point of view, using KL grades as a reference can also be considered a limitation, thus we suggest future studies to focus on OARS grading of the knee joints [45], [55]. Finally, another clinical limitation of this work is that OA is typically graded at a late stage, when it is already present, and we think that the next steps should rather focus on developing models for predicting OA progression [56] using SSL and partially labeled data. These data are available in hospital archives and can be leveraged at a low cost.

To conclude, we would like to highlight the clinical implications of our work. Firstly, we demonstrated that highly accurate KL grading can be done with only small amounts of labeled data, which allows small research teams and medical device vendors to build generalizable models suitable for clinical use. We highlight that this work demonstrates results drastically outperforming all the previous studies on automatic KL grading, and also human-level agreement (KC of 0.56 – 0.85 [57]), while the proposed method requires drastically less data than has been previously used. Secondly the proposed method can significantly reduce routine work done by radiologists, which on the societal level can lead to cost savings while at the same time improving the quality of health care. Thirdly, we expect the proposed method and the network architecture to generalize to other domains, such as hips and hands. Finally, we think that this study is an important step towards data efficient medical image recognition, which is currently lacking thoroughly validated methodologies.

#### ACKNOWLEDGMENTS

The OAI is a public-private partnership comprised of five contracts (N01-AR-2-2258; N01-AR-2-2259; N01-AR-2-2260; N01-AR-2-2261; N01-AR-2-2262) funded by the National Institutes of Health, a branch of the Department of Health and Human Services, and conducted by the OAI Study Investigators. Private funding partners include Merck Research Laboratories; Novartis Pharmaceuticals Corporation, GlaxoSmithKline; and Pfizer, Inc. Private sector funding for the OAI is managed by the Foundation for the National Institutes of Health.

The MOST is comprised of four cooperative grants (Felson - AG18820; Torner - AG18832; Lewis - AG18947; and Nevitt - AG19069) funded by the National Institutes of Health, a branch of the Department of Health and Human Services, and conducted by MOST study investigators. This manuscript was

prepared using MOST data and does not necessarily reflect the opinions or views of MOST investigators.

We would like to acknowledge the strategic funding of the University of Oulu, KAUTE foundation and Sigrid Juselius Foundation, Finland.

Dr. Claudia Lindner is acknowledged for providing BoneFinder. Iaroslav Melekhov is acknowledged for proposing the ablation study of the pooling method. Phuoc Dat Nguyen is acknowledged for discussions about *mixup*.

#### REFERENCES

- [1] D. J. Hunter and S. Bierma-Zeinstra, "Osteoarthritis," *The Lancet*, vol. 393, no. 10182, pp. 1745–1759, Apr. 2019.
- [2] P. A. Dieppe and L. S. Lohmander, "Pathogenesis and management of pain in osteoarthritis," *The Lancet*, vol. 365, no. 9463, pp. 965–973, 2005.
- [3] A. Mobasheri and M. Batt, "An update on the pathophysiology of osteoarthritis," *Annals of physical and rehabilitation medicine*, vol. 59, no. 5–6, pp. 333–339, 2016.
- [4] A. Mathiessen and P. G. Conaghan, "Synovitis in osteoarthritis: current understanding with therapeutic implications," *Arthritis research & therapy*, vol. 19, no. 1, p. 18, 2017.
- [5] C. Palazzo, C. Nguyen, M.-M. Lefevre-Colau, F. Rannou, and S. Poiradeau, "Risk factors and burden of osteoarthritis," *Annals of physical and rehabilitation medicine*, vol. 59, no. 3, pp. 134–138, 2016.
- [6] B. S. Ferket, Z. Feldman, J. Zhou, E. H. Oei, S. M. Bierma-Zeinstra, and M. Mazumdar, "Impact of total knee replacement practice: cost effectiveness analysis of data from the osteoarthritis initiative," *bmj*, vol. 356, p. j1131, 2017.
- [7] M. Cross, E. Smith, D. Hoy, S. Nolte, I. Ackerman, M. Fransen, L. Bridgett, S. Williams, F. Guillemin, C. L. Hill *et al.*, "The global burden of hip and knee osteoarthritis: estimates from the global burden of disease 2010 study," *Annals of the rheumatic diseases*, vol. 73, no. 7, pp. 1323–1330, 2014.
- [8] S. Salih and A. Hamer, "Hip and knee replacement," *Surgery (Oxford)*, vol. 31, no. 9, pp. 482 – 487, 2013, orthopaedics: Lower limb.
- [9] A. A. Guccione, D. T. Felson, J. J. Anderson, J. M. Anthony, Y. Zhang, P. Wilson, M. Kelly-Hayes, P. A. Wolf, B. E. Kreger, and W. B. Kannel, "The effects of specific medical conditions on the functional limitations of elders in the framingham study," *American journal of public health*, vol. 84, no. 3, pp. 351–358, 1994.
- [10] K. N. Murtagh and H. B. Hubert, "Gender differences in physical disability among an elderly cohort," *American journal of public health*, vol. 94, no. 8, pp. 1406–1411, 2004.
- [11] S. Glyn-Jones, A. Palmer, R. Agricola, A. Price, T. Vincent, H. Weinans, and A. Carr, "Osteoarthritis," *The Lancet*, vol. 386, no. 9991, pp. 376–387, 2015.
- [12] A. T. Toivanen, M. Heliövaara, O. Impivaara, J. P. Arokoski, P. Knekt, H. Lauren, and H. Kröger, "Obesity, physically demanding work and traumatic knee injury are major risk factors for knee osteoarthritis: a population-based study with a follow-up of 22 years," *Rheumatology*, vol. 49, no. 2, pp. 308–314, 2009.
- [13] M. E. Miller, W. J. Rejeski, S. P. Messier, and R. F. Loeser, "Modifiers of change in physical functioning in older adults with knee pain: the observational arthritis study in seniors (oasis)," *Arthritis Care & Research: Official Journal of the American College of Rheumatology*, vol. 45, no. 4, pp. 331–339, 2001.
- [14] G. Leardini, F. Salaffi, R. Caporali, B. Canesi, L. Rovati, R. Montanelli, I. G. for Study of the Costs of Arthritis *et al.*, "Direct and indirect costs of osteoarthritis of the knee," *Clin Exp Rheumatol*, vol. 22, no. 6, pp. 699–706, 2004.
- [15] D. J. Hunter and S. Bierma-Zeinstra, "Osteoarthritis," *The Lancet*, vol. 393, no. 10182, pp. 1745 – 1759, 2019.
- [16] A. J. Carr, O. Robertsson, S. Graves, A. J. Price, N. K. Arden, A. Judge, and D. J. Beard, "Knee replacement," *The Lancet*, vol. 379, no. 9823, pp. 1331 – 1340, 2012.
- [17] A. J. Price, A. Alvand, A. Troelsen, J. N. Katz, G. Hooper, A. Gray, A. Carr, and D. Beard, "Knee replacement," *The Lancet*, vol. 392, no. 10158, pp. 1672 – 1682, 2018.
- [18] P. Baker, K. Muthumayandi, C. Gerrand, B. Kleim, K. Bettinson, and D. Deehan, "Influence of body mass index (bmi) on functional improvements at 3 years following total knee replacement: a retrospective cohort study," *PloS one*, vol. 8, no. 3, p. e59079, 2013.

- [19] K. Baker and T. McAlindon, "Exercise for knee osteoarthritis," *Current Opinion in Rheumatology*, vol. 12, no. 5, pp. 456–463, 2000.
- [20] S. Reichenbach, A. Guermazi, J. Niu, T. Neogi, D. J. Hunter, F. W. Roemer, C. E. McLennan, G. Hernandez-Molina, and D. T. Felson, "Prevalence of bone attrition on knee radiographs and mri in a community-based cohort," *Osteoarthritis and cartilage*, vol. 16, no. 9, pp. 1005–1010, 2008.
- [21] A. Tiulpin, J. Thevenot, E. Rahtu, P. Lehenkari, and S. Saarakkala, "Automatic knee osteoarthritis diagnosis from plain radiographs: A deep learning-based approach," *Scientific reports*, vol. 8, no. 1, p. 1727, 2018.
- [22] J. Buckland-Wright, I. Carmichael, and S. Walker, "Quantitative microfocal radiography accurately detects joint changes in rheumatoid arthritis," *Annals of the rheumatic diseases*, vol. 45, no. 5, pp. 379–383, 1986.
- [23] J. Antony, K. McGuinness, K. Moran, and N. E. O'Connor, "Automatic detection of knee joints and quantification of knee osteoarthritis severity using convolutional neural networks," in *International conference on machine learning and data mining in pattern recognition*. Springer, 2017, pp. 376–390.
- [24] A. J. Antony, "Automatic quantification of radiographic knee osteoarthritis severity and associated diagnostic features using deep convolutional neural networks," Ph.D. dissertation, Dublin City University, 2018.
- [25] S. Laine and T. Aila, "Temporal ensembling for semi-supervised learning," *arXiv preprint arXiv:1610.02242*, 2016.
- [26] V. Verma, A. Lamb, J. Kannala, Y. Bengio, and D. Lopez-Paz, "Interpolation consistency training for semi-supervised learning," *arXiv preprint arXiv:1903.03825*, 2019.
- [27] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. Raffel, "Mixmatch: A holistic approach to semi-supervised learning," *arXiv preprint arXiv:1905.02249*, 2019.
- [28] H. Su, Z. Yin, S. Huh, T. Kanade, and J. Zhu, "Interactive cell segmentation based on active and semi-supervised learning," *IEEE Transactions on Medical Imaging*, vol. 35, no. 3, pp. 762–777, March 2016.
- [29] B. Wang, K. W. Liu, K. M. Prastawa, A. Irima, P. M. Vespa, J. D. van Horn, P. T. Fletcher, and G. Gerig, "4d active cut: An interactive tool for pathological anatomy modeling," in *2014 IEEE 11th International Symposium on Biomedical Imaging (ISBI)*, April 2014, pp. 529–532.
- [30] J. E. Iglesias, C.-Y. Liu, P. Thompson, and Z. Tu, "Agreement-based semi-supervised learning for skull stripping," in *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2010*, T. Jiang, N. Navab, J. P. W. Pluijm, and M. A. Viergever, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 147–154.
- [31] N. Kumar, P. Uppala, K. Duddu, H. Sreedhar, V. Varma, G. Guzman, M. Walsh, and A. Sethi, "Hyperspectral tissue image segmentation using semi-supervised nmf and hierarchical clustering," *IEEE Transactions on Medical Imaging*, vol. 38, no. 5, pp. 1304–1313, May 2019.
- [32] X. Li, L. Yu, H. Chen, C.-W. Fu, and P.-A. Heng, "Semi-supervised Skin Lesion Segmentation via Transformation Consistent Self-ensembling Model," *arXiv e-prints*, p. arXiv:1808.03887, Aug 2018.
- [33] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems 27*, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2014, pp. 2672–2680.
- [34] M. Belkin, P. Niyogi, and V. Sindhwani, "Manifold regularization: A geometric framework for learning from labeled and unlabeled examples," *Journal of machine learning research*, vol. 7, no. Nov, pp. 2399–2434, 2006.
- [35] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," *arXiv preprint arXiv:1710.09412*, 2017.
- [36] A. Oliver, A. Odena, C. A. Raffel, E. D. Cubuk, and I. Goodfellow, "Realistic evaluation of deep semi-supervised learning algorithms," in *Advances in Neural Information Processing Systems*, 2018, pp. 3239–3250.
- [37] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," in *Advances in neural information processing systems*, 2017, pp. 1195–1204.
- [38] X. Zhu and Z. Ghahramani, "Learning from labeled and unlabeled data with label propagation," Carnegie Mellon University, Tech. Rep., 2002.
- [39] D.-H. Lee, "Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks," in *Workshop on Challenges in Representation Learning, ICML*, vol. 3, 2013, p. 2.
- [40] O. Chapelle and A. Zien, "Semi-supervised classification by low density separation," in *AISTATS*, vol. 2005. Citeseer, 2005, pp. 57–64.
- [41] J. Antony, K. McGuinness, N. E. O'Connor, and K. Moran, "Quantifying radiographic knee osteoarthritis severity using deep convolutional neural networks," in *2016 23rd International Conference on Pattern Recognition (ICPR)*. IEEE, 2016, pp. 1195–1200.
- [42] P. Chen, L. Gao, X. Shi, K. Allen, and L. Yang, "Fully automatic knee osteoarthritis severity grading using deep neural networks with a novel ordinal loss," *Computerized Medical Imaging and Graphics*, 2019.
- [43] B. Norman, V. Pedoia, N. Noworolski, T. M. Link, and S. Majumdar, "Applying densely connected convolutional neural networks for staging osteoarthritis severity from plain radiographs," *Journal of digital imaging*, vol. 32, no. 3, pp. 471–477, 2019.
- [44] M. Górriz, J. Antony, K. McGuinness, X. Giró-i Nieto, and N. E. O'Connor, "Assessing knee oa severity with cnn attention-based end-to-end architectures," *arXiv preprint arXiv:1908.08856*, 2019.
- [45] A. Tiulpin and S. Saarakkala, "Automatic grading of individual knee osteoarthritis features in plain radiographs using deep convolutional neural networks," *arXiv preprint arXiv:1907.08020*, 2019.
- [46] B. Athiwaratun, M. Finzi, P. Izmailov, and A. G. Wilson, "There are many consistent explanations of unlabeled data: Why you should average," *arXiv preprint arXiv:1806.05594*, 2018.
- [47] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, "Striving for simplicity: The all convolutional net," *arXiv preprint arXiv:1412.6806*, 2014.
- [48] H. Guo, Y. Mao, and R. Zhang, "Mixup as locally linear out-of-manifold regularization," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 3714–3722.
- [49] O. Chapelle, B. Scholkopf, and A. Zien, "Semi-supervised learning (chapelle, o. et al., eds.; 2006)[book reviews]," *IEEE Transactions on Neural Networks*, vol. 20, no. 3, pp. 542–542, 2009.
- [50] C. Lindner, S. Thiagarajah, J. M. Wilkinson, G. A. Wallis, T. F. Cootes, arcOGEN Consortium *et al.*, "Fully automatic segmentation of the proximal femur using random forest regression voting," *IEEE transactions on medical imaging*, vol. 32, no. 8, pp. 1462–1472, 2013.
- [51] R. J. Urbanowicz and J. H. Moore, "Exstracs 2.0: description and evaluation of a scalable learning classifier system," *Evolutionary intelligence*, vol. 8, no. 2-3, pp. 89–116, 2015.
- [52] J. Cohen, "A coefficient of agreement for nominal scales," *Educational and psychological measurement*, vol. 20, no. 1, pp. 37–46, 1960.
- [53] N. E. Breslow and D. G. Clayton, "Approximate inference in generalized linear mixed models," *Journal of the American statistical Association*, vol. 88, no. 421, pp. 9–25, 1993.
- [54] F. Wilcoxon, "Individual comparisons by ranking methods," in *Breakthroughs in statistics*. Springer, 1992, pp. 196–202.
- [55] R. D. Altman and G. Gold, "Atlas of individual radiographic features in osteoarthritis, revised," *Osteoarthritis and cartilage*, vol. 15, pp. A1–A56, 2007.
- [56] A. Tiulpin, S. Klein, S. Bierma-Zeinstra, J. Thevenot, E. Rahtu, J. van Meurs, E. H. Oei, and S. Saarakkala, "Multimodal machine learning-based knee osteoarthritis progression prediction from plain radiographs and clinical data," *arXiv preprint arXiv:1904.06236*, 2019.
- [57] K. Klara, J. E. Collins, E. Gurary, S. A. Elman, D. S. Stenquist, E. Losina, and J. N. Katz, "Reliability and accuracy of cross-sectional radiographic assessment of severe knee osteoarthritis: role of training and experience," *The Journal of rheumatology*, vol. 43, no. 7, pp. 1421–1426, 2016.
- [58] H. H. Nguyen, E. Panfilov, and A. Tiulpin, "Collagen: Deep learning framework for reproducible experiments," 2019. [Online]. Available: <https://github.com/MIPT-Oulu/Collagen>
- [59] A. Tiulpin, "Solt: Streaming over lightweight transformations," 2019. [Online]. Available: <https://github.com/MIPT-Oulu/solt>
- [60] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

## SUPPLEMENT

## A. Experimental Details

All our experiments were conducted on V100 NVidia GPUs. We used PyTorch, Collagen [58], and SOLT [59] libraries in our codebase. To train all the models, we utilized the Adam optimizer [60] with a learning rate of  $1e-4$  and without weight decay regularization. Except for the SL baseline model [21], which had a dropout of 0.2 in the bottleneck of the model, we set the dropout rate to 0.35, and used it in multiple blocks our model. Note that we did not use Exponential Mean Averaging (EMA) [37] in any methods. For *Semixup*, we sampled 1 pair of random augmentations (i.e.,  $N_t = 2$ ) and 1 mixing operator (i.e.,  $N_m = 1$ ) for each knee image. We set the parameter  $\alpha$  to 0.75, and the unsupervised weight vector  $\mathbf{w}$  to  $[2, 2, 4]^T$  (Supplementary Algorithm 1).

In the pre-processing step, we transformed images by random Gaussian noise, rotation, cropping, and Gamma correction augmentations, whose parameters are described in Supplementary Table S6. A description of our network architecture is described in Supplementary Table S7. Detailed implementations of related methods are expressed in Supplementary Section S5-B.

## B. Baseline Methods

We reimplemented the previous SSL methods (<https://github.com/MIPT-Oulu/semixup>) in a common codebase as suggested by [36]. All the baseline methods used the same architecture and pre-processing operators as our approach. We trained each method with a batch size of 40 for 500 epochs, each of which is a full pass through labeled data.

The weight of consistency regularization terms,  $\mathbf{w}$ , varied among the methods. As such, for  $\Pi$ -model, we searched it through the set  $\{1, 10, 50\}$ , and found that it worked well with  $w = 1$ . In ICT, we varied  $w$  from 0 to 100 using a sigmoid ramp-up schedule within the first 80 epochs as in [26]. Following Verma *et al.*, we did not use *mixup* for labeled data in ICT. For MixMatch, we used 2 augmentations, a sharpening hyperparameter of 0.5 ( $T$  in the original paper), and a Beta( $\alpha, \alpha$ ) distribution with  $\alpha$  of 0.75. We run a search for its unsupervised coefficient from  $\{1, 10, 50\}$  based on [27], and found the coefficient of 10 to be the best choice. In addition, we trained *mixup* using a Beta( $\alpha, \alpha$ ) distribution with  $\alpha$  of 0.75.

TABLE S6  
ORDERED LIST OF AUGMENTATIONS

Order	Augmentation	Probability	Parameter
1	Gaussian noise	0.5	0.3
2	Rotation	1	[-10, 10]
3	Padding	1	5%
4	Cropping	1	128x128
5	Gamma correction	0.5	[0.5, 1.5]

## Algorithm 1: Semixup

---

**Data:**  $\mathcal{X}_{ul}$ : combined labeled and unlabeled examples  
**Data:**  $(\mathcal{X}_l, \mathcal{Y})$ : labeled examples  
**Input:**  $f_\theta$ : a neural network with parameters  $\theta$ .  
**Input:**  $\alpha$ : parameter of Beta distribution  
**Input:**  $\mathbf{w}$ : weights of consistency terms  
**Input:**  $N_T$ : the number of iterations  
**Input:**  $N_c$ : the number of classes  
**Input:**  $N_b$ : batch size

---

```

1 for  $t = 1 \dots N_T$  do
2    $B_l \leftarrow \{(x_i, y_i)\}_{i=1}^{N_b} \subset (\mathcal{X}_l, \mathcal{Y})$ ,  $B_{ul} \leftarrow \{x_i\}_{i=1}^{N_b} \subset \mathcal{X}_{ul}$ 
3    $\mathcal{L}_l, \mathcal{L}_u \leftarrow 0, 0$ 
4   for  $x_i \in B_{ul}$  do
5     Sample  $x_j \in B_{ul}$ ,  $\lambda \sim \text{Beta}(\alpha, \alpha)$ ,  $T$  and  $T' \sim p(\tau)$ 
6      $\lambda \leftarrow \max(\lambda, 1 - \lambda)$ 
7      $x_{mix} \leftarrow \text{Mix}_\lambda(Tx_i, x_j)$ 
8      $p_{mix} \leftarrow \text{Mix}_\lambda(f_\theta(Tx_i), f_\theta(x_j))$ 
9      $\mathcal{L}_u \leftarrow \mathcal{L}_u + \mathbf{w}^{(0)} \|f_\theta(Tx_i) - f_\theta(T'x_i)\|_2^2$ 
10     $\mathcal{L}_u \leftarrow \mathcal{L}_u + \mathbf{w}^{(1)} \|f_\theta(x_{mix}) - f_\theta(Tx_i)\|_2^2$ 
11     $\mathcal{L}_u \leftarrow \mathcal{L}_u + \mathbf{w}^{(1)} \|f_\theta(x_{mix}) - f_\theta(T'x_i)\|_2^2$ 
12     $\mathcal{L}_u \leftarrow \mathcal{L}_u + \mathbf{w}^{(2)} \|f_\theta(x_{mix}) - p_{mix}\|_2^2$ 
13  end
14  for  $(x_i, y_i) \in B_l$  do
15    Sample  $(x_j, y_j) \in B_l$ ,  $\lambda \sim \text{Beta}(\alpha, \alpha)$ 
16     $x_{mix} \leftarrow \text{Mix}_\lambda(x_i, x_j)$ 
17     $\mathcal{L}_l \leftarrow \lambda \mathcal{L}_{ce}(x_{mix}, y_i) + (1 - \lambda) \mathcal{L}_{ce}(x_{mix}, y_j)$ 
18  end
19   $\mathcal{L} \leftarrow \frac{1}{N_b} \mathcal{L}_l + \frac{1}{N_b N_c} \mathcal{L}_u$ 
20  UpdateStep( $\theta$ ,  $\nabla_\theta \mathcal{L}$ )
21 end
```

---

TABLE S7  
DETAILED DESCRIPTION OF OUR ARCHITECTURE.

Layer name	Output shape	Layer
Input	128x128	
Conv1_1	128x128	ConvBlock 3x3, 32, S = 1
Conv1_2	128x128	ConvBlock 3x3, 32, S = 1
Conv1_3	128x128	ConvBlock 3x3, 32, S = 1
Dropout_1	128x128	Dropout
Conv2_1	64x64	ConvBlock 3x3, 64, S = 2
Conv2_2	64x64	ConvBlock 3x3, 64, S = 1
Dropout_2	64x64	Dropout
Conv3_1	32x32	ConvBlock 3x3, 128, S = 2
Conv3_2	32x32	ConvBlock 3x3, 128, S = 1
Dropout_3	32x32	Dropout
Conv4_1	16x16	ConvBlock 3x3, 256, S = 2
Conv4_2	16x16	ConvBlock 3x3, 256, S = 1
Dropout_4	16x16	Dropout
SAM	1x1	SAM block
Merge	512	Concat
Dropout_5	512	Dropout
Output	5	Linear, Softmax

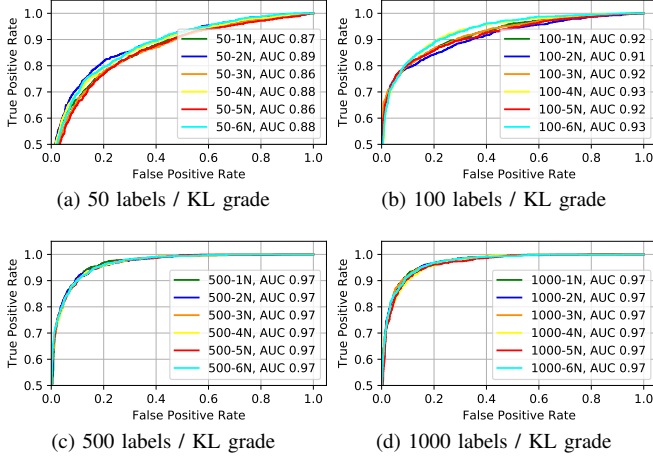


Fig. S8. ROC curves and AUC of models trained by *Semixup* using  $N$  labeled samples per KL grade. Each subplot shows the results of 6 models trained on 6 different amounts of unlabeled data. (a)  $N = 5 \times 50$ . (b)  $N = 5 \times 100$ . (c)  $N = 5 \times 500$  (d)  $N = 5 \times 1000$ .

TABLE S8

RESULTS (BA, %) ON AN INDEPENDENT TEST SET DERIVED FROM THE MOST DATA FOR ALL THE SL AND SSL MODELS. HERE, WE VARIED THE AMOUNT OF LABELS PER KL GRADE AS WELL AS THE AMOUNT OF UNLABELED DATA IN EACH SETTING (24 SETTINGS PER EACH SSL METHOD). THE RESULTS IN BOLD ARE THE BEST ONES IN EACH SETTING. THE UNDERLINE HIGHLIGHTS THE SECOND BEST MODELS. THE BOTTOM PART OF THE TABLE SHOWS THE AVERAGE PERFORMANCE ACROSS ALL UNLABELED DATA CONFIGURATIONS.

Method	# unlabeled	# labels / KL grade ( $N/5$ )			
	data	50	100	500	1000
ICT [26]	$N$	42.8	42.1	60.6	67.8
$\Pi$ model [25]		42.6	<u>52.9</u>	61.3	65.7
MixMatch [27]		41.9	<b>56.0</b>	65.0	64.6
<i>Semixup</i> (Ours)		<b>45.5</b>	52.4	<u>66.7</u>	<b>71.0</b>
ICT [26]	$2N$	40.7	50.5	65.4	66.7
$\Pi$ model [25]		<b>43.2</b>	<b>56.2</b>	64.7	63.9
MixMatch [27]		37.7	51.7	<u>66.1</u>	61.2
<i>Semixup</i> (Ours)		<u>43.0</u>	<u>53.4</u>	<b>66.2</b>	<b>66.8</b>
ICT [26]	$3N$	<b>45.3</b>	48.4	66.5	65.8
$\Pi$ model [25]		<u>43.8</u>	<u>52.0</u>	<b>66.8</b>	<u>69.0</u>
MixMatch [27]		43.4	47.6	65.7	68.2
<i>Semixup</i> (Ours)		41.3	<b>54.5</b>	<b>69.7</b>	<b>70.1</b>
ICT [26]	$4N$	46.8	52.3	63.1	66.9
$\Pi$ model [25]		45.9	<u>55.8</u>	65.5	<b>67.7</b>
MixMatch [27]		45.1	<u>55.4</u>	<b>67.6</b>	64.6
<i>Semixup</i> (Ours)		<b>46.9</b>	<b>58.8</b>	<u>66.1</u>	<u>66.9</u>
ICT [26]	$5N$	<b>47.7</b>	<u>52.7</u>	64.8	<u>64.6</u>
$\Pi$ model [25]		44.9	51.3	<b>67.1</b>	63.9
MixMatch [27]		40.1	51.0	62.6	<b>64.9</b>
<i>Semixup</i> (Ours)		44.7	<b>55.8</b>	<u>65.9</u>	61.7
ICT [26]	$6N$	<b>46.2</b>	53.5	62.6	66.1
$\Pi$ model [25]		44.2	<b>54.1</b>	<b>66.9</b>	68.5
MixMatch [27]		42.2	53.7	62.7	<u>68.4</u>
<i>Semixup</i> (Ours)		<u>45.4</u>	<u>53.8</u>	<u>65.4</u>	<b>69.5</b>
Average					
ICT [26]	$N \dots 6N$	<b>44.9</b>	49.9	63.8	66.3
$\Pi$ model [25]		44.1	<u>53.7</u>	<u>65.4</u>	<u>66.5</u>
MixMatch [27]		41.7	52.6	64.9	65.3
<i>Semixup</i> (Ours)		44.5	<b>54.8</b>	<b>66.7</b>	<b>67.7</b>

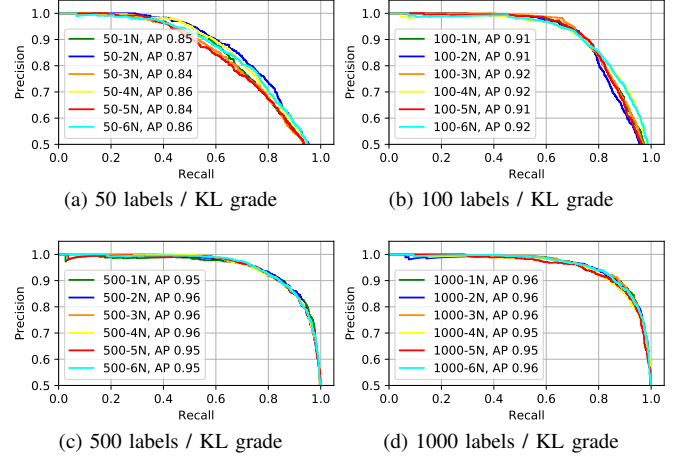


Fig. S9. PR curves and AP of the models trained by *Semixup* using  $N$  labeled samples per KL grade. Each subplot shows the results of 6 models trained on 6 different amounts of unlabeled data. (a)  $N = 5 \times 50$ ; (b)  $N = 5 \times 100$ ; (c)  $N = 5 \times 500$ ; (d)  $N = 5 \times 1000$ .

TABLE S9

STATISTICAL COMPARISONS USING ONE-SIDED WILCOXON SIGNED-RANK TEST.  $SL^+$  INDICATES OUR FULLY SL MODELS EQUIPPED WITH EITHER SAM-HV OR SAM-VH.  $SL^\ddagger$  INDICATES OUR MODEL WITH SAM-HV. \* AND \*\* SIGNS CORRESPOND TO  $p < 0.05$  AND  $p < 0.001$ , RESPECTIVELY.

Main method		Compared method		p-value
Name	# labels	Name	# labels	
<i>Semixup</i>	250	$SL^+$	250	5.4e-2
		$\Pi$ model [25]	250	1.5e-1
		ICT [26]	250	6.3e-1
		MixMatch [27]	250	5.4e-2
	500	$SL^+$	500	5.2e-5 **
		$\Pi$ model [25]	500	1.6e-3 *
		ICT [26]	500	1.1e-4 **
		MixMatch [27]	500	5.2e-5 **
	2,500	$SL^+$	2,500	7.5e-4 **
		$\Pi$ model [25]	2,500	8.4e-2
		ICT [26]	2,500	2.6e-4 **
		MixMatch [27]	2,500	1.2e-1
	5,000	$SL^+$	5,000	3.7e-2 *
		$\Pi$ model [25]	5,000	2.5e-1
		ICT [26]	5,000	6.8e-2
		MixMatch [27]	5,000	1.6e-1
$SL^\ddagger$	35,000	<i>Semixup</i> (Ours)	2,500	5.4e-2
			5,000	3.7e-1