

# Deep Semi-Supervised and Self-Supervised Learning for Diabetic Retinopathy Detection

Jose Arrieta<sup>1</sup>, Oscar J. Perdomo<sup>2</sup>, and Fabio A. González<sup>1</sup>

<sup>1</sup> MindLab Research Group, Universidad Nacional de Colombia, Colombia  
 {jmarrieta<sup>1</sup>, ojperdomo<sup>2</sup>, fagonzalezo<sup>3</sup>}@unal.edu.co

<sup>2</sup> Universidad del Rosario, Colombia

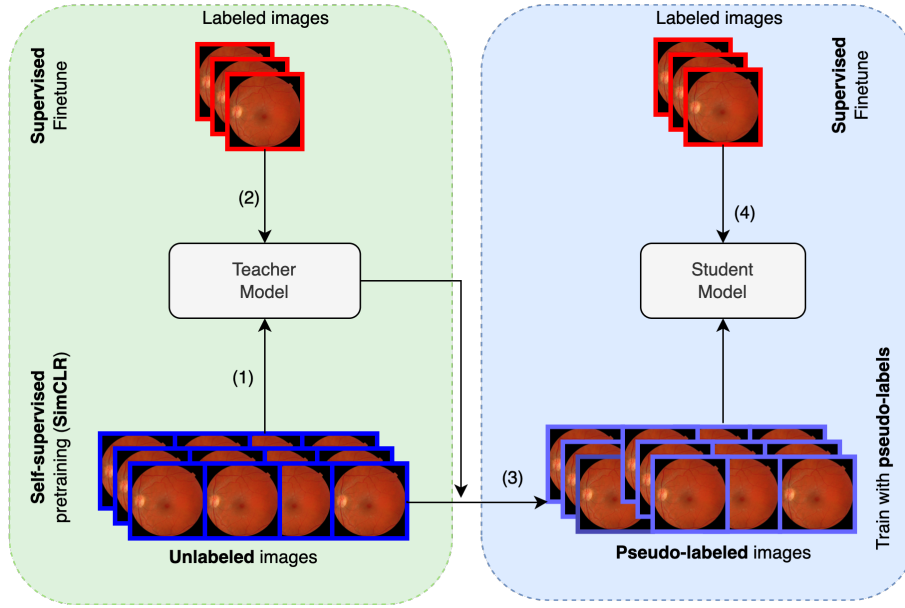
**Abstract.** Diabetic retinopathy (DR) is one of the leading causes of blindness in the working-age population of developed countries, caused by a side effect of diabetes that reduces the blood supply to the retina. Deep neural networks have been widely used in automated systems for DR classification on eye fundus images. However, these models need a large number of annotated images. In the medical domain, annotations from experts are costly, tedious, and time-consuming; as a result, a limited number of annotated images are available. This paper presents a semi-supervised method that leverages unlabeled images and labeled ones to train a model that detects diabetic retinopathy. The proposed method uses unsupervised pretraining via self-supervised learning followed by supervised fine-tuning with a small set of labeled images and knowledge distillation to increase the performance in classification task. This method was evaluated on the EyePACS test and Messidor-2 dataset achieving 0.94 and 0.89 AUC respectively using only 2% of EyePACS train labeled images.

**Keywords:** Diabetic Retinopathy · Medical Imaging · Deep Learning · Semi-supervised Learning · Self-supervised learning.

## 1 Introduction

Diabetic retinopathy (DR) is one of the leading causes of blindness in the working-age population in developed countries [15]. DR consists of a side effect of diabetes that reduces the blood supply to the retina, including lesions that appear on the surface, such as microaneurysms, exudates, hemorrhages, and cotton wool spots. Automated methods for detecting eye diseases present as an useful tool for early diagnosis, important to prevent the occurrence of blindness and lack of vision [5,3]. Recently, deep learning methods became popular for DR classification on eye fundus images because of their promising results [1,9,25,13,18,19,24]. However, the majority of these approaches use labeled images, whereas the manual labeling of medical images is expensive and time-consuming because it requires medical experts in the retina. As a result, relatively small labeled data sets are available to train deep learning models.

Some authors have proposed semi-supervised learning approaches to leverage unlabeled images and mitigate the lack of annotated images. In specific,



**Fig. 1.** The proposed method consists of four stages as follows: (1) Self-supervised pre-training using SimCLR with unlabeled images. (2) Supervised fine-tuning on a small set of labeled images. (3) Knowledge distillation from teacher model to student model using teacher’s output as pseudo-labels to train student model. (4) Fine-tuning of student model on a small set of labeled images.

Liu et al. [14] and Xie et al. [29] propose the use of generative adversarial networks (GAN) to extend classification while performing unsupervised image reconstruction with unlabeled data in training. State-of-the-art methods in deep semi-supervised methods such as MixMatch [4] and, FixMatch [22] comprise strong and weak augmentations of unlabeled images and make class pseudo-labels match via consistency regularization. In medical imaging, Hansen et al. [10] explore MixMatch in two distinct medical imaging domains as skin lesion diagnosis and lung cancer prediction. Pooch et al. [17] evaluate the performance of different methods including pseudo-labeling, Mean Teacher [23], Unsupervised Data Augmentation (UDA) [27], MixMatch [4] and FixMatch [22] in a chest radiography classification task.

Alternatively, self-supervised learning [6] presents as a new strategy to use unlabeled data to pre-train a neural network and allows the construction of models that can learn relevant image representations from unlabeled medical images. These visual representations from unlabeled data are used to train strong semi-supervised models [7]. Azizi et al. [2] presented the effectiveness of self-supervised learning with unlabeled data for domain-specific medical images in dermatology skin and chest X-ray classification. Vu et al. [26] use self-supervised contrastive learning and propose to use patient metadata to select the pairs of

medical images requiring them to be from the same patient. Kaku et al. [12], minimize the mean squared error between the intermediate layer representations in complement to the contrastive loss using medical datasets such as NIH-Chest X-rays, breast cancer histopathology, and diabetic retinopathy.

This paper presents a method that combines self-supervised learning along with a teacher-student and knowledge distillation strategies to perform semi-supervised learning with a small subset of labeled eye fundus images. The method achieves 0.94 AUC on EyePACS test and 0.89 AUC on Messidor-2 Diabetic Retinopathy dataset using only 2% EyePACS-Kaggle train labeled images (1000 images).

## 2 Methods

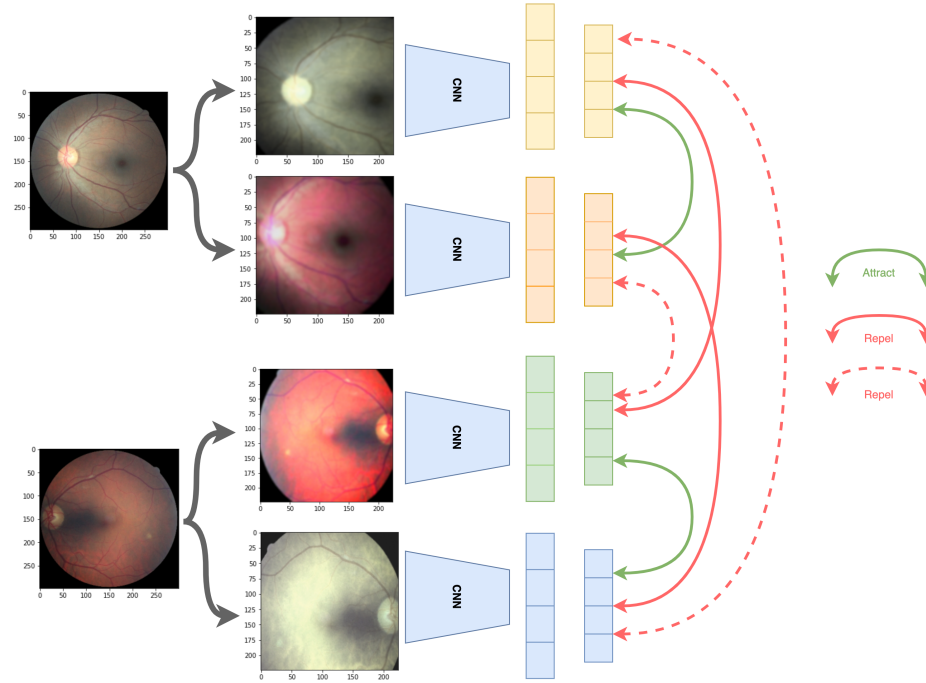
The proposed method consists of a teacher and a student network. Our method differs from previous approaches because the teacher network heavily relies on large unlabeled diabetic retinopathy (DR) images to learn useful representations in pre-training compared to the labeled ones available. Furthermore, the student network benefits from pseudo-labels created with the teacher model using unlabeled images. The complete semi-supervised workflow proposed is summarized in Figure 1 and consists of four stages.

**1) Self-supervised unlabeled pre-training:** First; the Teacher model is pretrained using the simple framework for contrastive learning of visual representation (SimCLR [6]). This process is unsupervised as it doesn't need the labels of the images. The whole EyePACS-Kaggle train dataset (57 146 images) without the labels and a ResNet-50 backbone were used to learn useful visual representations with domain-specific medical images. SimCLR uses strong data augmentations such as color distortions (brightness, contrast, saturation, hue), cropping and rotations as shown in Figure 2. Contrastive loss  $NT-Xent$  (the normalized temperature-scaled cross-entropy loss) is used to maximize the agreement of diverse representations of the same image in a latent space, defined as:

$$\mathbb{L}_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j) / \tau)}{\sum_{k=1}^{2N} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k) / \tau)} \quad (1)$$

where  $\mathbf{z}_i$  and  $\mathbf{z}_j$  correspond to the outputs of the same image from the non-linear projection head and  $\mathbf{z}_k$  corresponds to the remainder of representations of the other  $N$  images in the batch. Additionally,  $\tau$  is the temperature parameter  $t$ , and  $\text{sim}$  denotes a distance metric.

**2) Fine-tuning of Teacher model with labeled images:** In this phase, a small set of 1000 random samples of labeled fundus images are used to further refine the weights of the model. Fine-tuning is a common approach where a network from a previous task is used as starting point, and weights are adjusted to fit a new specific task. The proposed approach used the previously pretrained SimCLR ResNet-50 [28] model with an additional sigmoid layer at the top to



**Fig. 2.** Self-supervised pre-training with SimCLR uses strong data augmentations to maximize the agreement between different representations of the same image in the latent space.

make a binary classification from referable DR and non-referable DR.

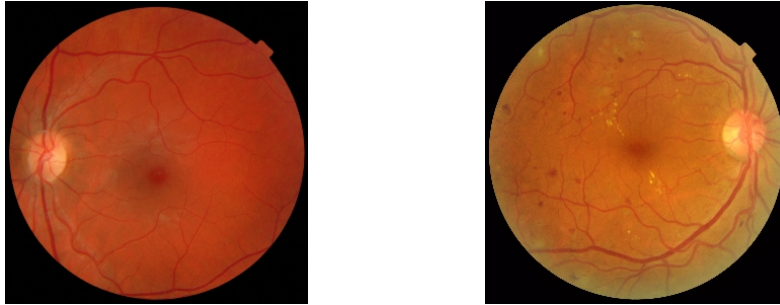
**3) Pseudo-labeling and knowledge distillation:** The third step consists of knowledge distillation from the ResNet-50 model (Teacher) to a DenseNet161 model (Student) using unlabeled data as input and teacher’s output as pseudo-labels. Knowledge distillation is the process of transferring knowledge from one model to another one, originally used to compress a large model to a smaller one, encouraging the student to match a teacher’s output [11]. However, we chose a DenseNet161, a deeper network than the teacher for Student architecture because Xie et al. [27] demonstrated that using student models that are equal to or larger than the teacher in addition to some noise to the student can improve the performance of the student beyond the teacher, making predictions with more difficult images by giving the student model enough capacity and difficult images to learn through. We transform soft-pseudo-labels to hard-pseudo-labels using a threshold of 0.5.

4) **Fine-tuning of Student model with labeled images:** Finally, the fourth task consists of a final fine-tuning of the student model with the same set of 1000 labeled eye fundus images.

### 3 Experimental Evaluation

#### 3.1 Datasets

EyePACS-Kaggle [20] consists of retina images taken under a variety of imaging conditions as shown in figure 3, where the clinician rate the presence of diabetic retinopathy in each image on a scale of 0 to 4 (No DR, Mild, Moderate, Severe, Proliferative), respectively. For the binary classification task, grades 0 and 1 correspond to non-referable DR, while grades 2, 3, and 4 correspond to referable DR according to the International Clinical Diabetic Retinopathy Scale [21]. The EyePACS-Kaggle dataset is partitioned into train and tests datasets, with 35121 images and 42918 images respectively of different sizes. In this work, we used the EyePACS-Kaggle partition defined in Voets et al. [25] which consists of a train set of 57146 images and a test set of 8790 images. Training is performed on the EyePACS-Kaggle train and evaluation is performed on the EyePACS-Kaggle test and Messidor-2 datasets. Messidor-2 [8] consists of a public dataset that contains 1748 images with grades adjudicated by a panel of three retina specialists and constitutes a standard dataset used to compare performance results in Diabetic Retinopathy detection. Following Voets et al. [25] and Gulshan et al. [9] pre-processing, all fundus images are centered and resized to 299 x 299 pixels, with the fundus center in the middle of the image.



**Fig. 3.** Example of eye fundus image (EFI) non-referable DR on the left and a referable DR example on the right.

#### 3.2 Experimental setup

All models were implemented using Pytorch [16]. EyePACS-Kaggle train dataset without labels is used for self-supervised pretraining a ResNet-50 network using

SimCLR data augmentations alongside a projection head with a contrastive loss function, a learning rate of  $1 \times 10^{-5}$ , weight decay of  $5 \times 10^{-4}$  and batch of size 64 for 100 epochs. For knowledge distillation, the previously trained ResNet-50 network is used as a teacher to predict pseudo-labels of the EyePACS-Kaggle train dataset while using them to train a different architecture. DenseNet161, a deeper network is used to improve the performance of the student beyond the teacher. The data augmentation configuration consists of random color jitter and random horizontal flips with BinaryCrossEntropy loss function and stochastic gradient descent as an optimizer, a learning rate of  $1 \times 10^{-4}$  and batch size 32 for 200 epochs. A final fine-tuning of the student network is performed using the 1000 images labeled data for 100 epochs.

Three baselines were defined. First, the supervised Inception-V3 network as proposed in Voets et al. [25] and Krause et al. [13], although some other architectures were also tested with similar performance. Additionally, two state-of-the-art hybrid methods, such as MixMatch and FixMatch, were chosen to compare the performance of semi-supervised approaches.

## 4 Results

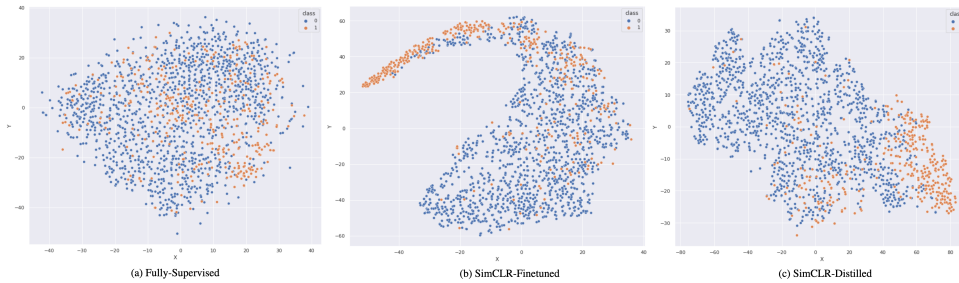
The proposed method is compared to a supervised baseline that uses the fraction of 1000 labeled images (2% of EyePACS-Kaggle) and semi-supervised learning methods such as FixMatch and MixMatch that make class pseudo-labels match via consistency regularization using additional unlabeled data. Additionally, a supervised method with the complete train dataset is used as a reference for performance. The proposed method successfully outperforms the methods trained with a labeled fraction of 2% of EyePACS train and its at similar performance compared to 100% of train labeled images method leveraging the use of unlabeled data in the pre-train stage and creating useful visual representations. Additionally, the knowledge distillation step creates a boosting under the AUC metric as reported in Table 1 completing a better generalization task of detection of diabetic retinopathy using only 1000 labeled images 2% of EyePACS-Kaggle and the rest as unlabeled images. Further, in Figure 4 the representation of visual features is presented using the t-SNE technique in the Messidor-2 dataset. The supervised representation of DR and non-DR seems not easily separable. In contrast, the SimCLR-Finetuned (Teacher) and SimCLR-Distilled (Student) seems more separable between the DR and non-DR instances, improving visual representations and generalization.

## 5 Discussion and Conclusion

Despite the deep learning model’s need for a large amount of labeled data for their training process, this paper presents a method that combines self-supervised learning along with a teacher-student and knowledge distillation strategies to perform semi-supervised learning with a small subset of labeled eye fundus images. The method achieves 0.94 AUC on the EyePACS test and

**Table 1.** AUC performance measures of baselines and proposed method on test datasets. Supervised baseline is trained using 2% of EyePACS train labeled images. Semi-supervised methods use 2% of EyePACS train labeled images and the rest as unlabeled images.

Label fraction	Method	Architecture	Kaggle Test	Messidor2
100%	Supervised	InceptionV3	0.96	0.88
2% <i>Train</i>	Supervised	InceptionV3	0.88	0.66
	MixMatch [4]	ResNet50	0.85	0.64
	FixMatch [22]	ResNet50	0.83	0.79
	SimCLR-Finetuned (Teacher)	ResNet50	0.92	0.85
	SimCLR-Distilled (Student)	DenseNet161	<b>0.94</b>	<b>0.89</b>



**Fig. 4.** t-SNE visual representation in Messidor-2: (a) Supervised, (b) SimCLR-Finetuned (Teacher), and (c) SimCLR-Distilled methods (Student).

0.89 AUC on Messidor-2 Diabetic Retinopathy dataset using only 2% EyePACS train labeled images improving the performance from the baseline supervised approach, and with greater performance than the other two state-of-the-art methods, alongside a better generalization compared to the teacher model. Self-supervised learning resulted in a great pre-training strategy to use unlabeled medical images and create richer and better visual representations improving downstream tasks such as diabetic retinopathy detection in this case. Although one disadvantage of this approach consists of computing and memory intensive, as it creates several representations from each image in the training batch and could require a lengthy training process. Also, knowledge distillation from a teacher to a student model improved performance in the presented results. Furthermore, the student network benefits from pseudo-labels created with the teacher model using unlabeled images, improving the generalization of outcomes and enhancing the performance beyond the teacher. In future works, the validation of our proposed model in a clinical environment, combining images from different sources, equipment, and audiences, especially in the unsupervised steps. Additionally, it was hypothesized that this method could be applied to other medical image domains where labeled images are scarce and constitute a vital research topic for deep learning models.

## References

1. Norah Asiri, Muhammad Hussain, Fadwa Al, and Nazih Alzaiddi. Deep learning based computer-aided diagnosis systems for diabetic retinopathy : A survey. *Artificial Intelligence In Medicine*, 99(December 2018):101701, 2019.
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. pages 1–19.
3. Emma Beede, Elizabeth Baylor, Fred Hersch, Anna Iurchenko, Lauren Wilcox, Dr. Paisan Raumviboonsuk, and Laura Vardoulakis. A Human-Centered Evaluation of a Deep Learning System Deployed in Clinics for the Detection of Diabetic Retinopathy. *CHI 2020 Paper*, pages 1–12, 2020.
4. David Berthelot, Avital Oliver, Nicholas Carlini, Ian Goodfellow, Colin Raffel, and Nicolas Papernot. MixMatch : A Holistic Approach to Semi-Supervised Learning. (NeurIPS):1–11, 2019.
5. Malavika Bhaskaranand, Chaithanya Ramachandra, Sandeep Bhat, Jorge Cuadros, Muneeswar G. Nittala, Srinivas R. Sadda, and Kaushal Solanki. The value of automated diabetic retinopathy screening with the EyeArt system: A study of more than 100,000 consecutive encounters from people with diabetes. *Diabetes Technology and Therapeutics*, 21(11):635–643, 2019.
6. Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A Simple Framework for Contrastive Learning of Visual Representations. 2019.
7. Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big Self-Supervised Models are Strong Semi-Supervised Learners. (NeurIPS):1–18, 2020.
8. Laboratoire de Traitement de l’Information Médicale. Messidor-2 dataset. <https://www.adcis.net/en/third-party/messidor2/>, 2011.
9. Varun Gulshan, Lily Peng, Marc Coram, Martin C Stumpe, Derek Wu, Arunachalam Narayanaswamy, Subhashini Venugopalan, Kasumi Widner, Tom Madams, Jorge Cuadros, Ramasamy Kim, Rajiv Raman, Philip C Nelson, Jessica L Mega, and Dale R Webster. Development and Validation of a Deep Learning Algorithm for Detection of Diabetic Retinopathy in Retinal Fundus Photographs. 94043:1–9, 2016.
10. Colin B. Hansen, Vishwesh Nath, Riqiang Gao, Camilo Bermudez, Yuankai Huo, Kim L. Sandler, Pierre P. Massion, Jeffrey D. Blume, Thomas A. Lasko, and Bennett A. Landman. Semi-supervised Machine Learning with MixMatch and Equivalence Classes. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 12446 LNCS:112–121, 2020.
11. Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the Knowledge in a Neural Network. pages 1–9, 2015.
12. Aakash Kaku, Sahana Upadhyay, and Narges Razavian. Intermediate layers matter in momentum contrastive self supervised learning. 2021.
13. Jonathan Krause, Varun Gulshan, Ehsan Rahimy, Peter Karth, Kasumi Widner, Greg S. Corrado, Lily Peng, and Dale R. Webster. Grader Variability and the Importance of Reference Standards for Evaluating Machine Learning Models for Diabetic Retinopathy. *Ophthalmology*, 125(8):1264–1272, 2018.
14. Sijie Liu, Jingmin Xin B, Jiayi Wu, and Peiwen Shi. Semi-supervised Adversarial Learning for Diabetic Retinopathy Screening. pages 60–68, 2019.



15. World Health Organisation. *World report on vision. 2019*, volume 214. 2019.
16. Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
17. Eduardo H.P. Pooch, Pedro Ballester, and Rodrigo C. Barros. Semi-supervised Classification of Chest Radiographs. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 12446 LNCS:172–179, 2020.
18. Gwenolé Quéllec, Katia Charrière, Yassine Boudi, and Béatrice Cochener. Deep image mining for diabetic retinopathy screening. 39:178–193, 2017.
19. Alexander Rakhlin. Diabetic Retinopathy detection through integration of Deep Learning classification framework. pages 1–11, 2018.
20. Diabetic retinopathy detection. Eyepacs. [www.kaggle.com/competitions/diabetic-retinopathy-detection](http://www.kaggle.com/competitions/diabetic-retinopathy-detection), 2015.
21. International Clinical Diabetic Retinopathy Disease Severity Scale. International Clinical Diabetic Retinopathy. (October):8500, 2002.
22. Kihyuk Sohn, David Berthelot, Chun-liang Li Zizhao, Zhang Nicholas, Ekin D Cubuk, Alex Kurakin, Han Zhang, and Colin Raffel. FixMatch: Simplifying Semi-Supervised Learning with Consistency and Confidenc.
23. 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*, 2017-December:1196–1205, 2017.
24. Santiago Toledo-cort, Melissa De La Pava, and Oscar Perd. Diabetic Retinopathy Diagnosis and Uncertainty.
25. Mike Voets, Mollersen Kajsa, and Lars Bongo. Reproduction study using public data of : Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. pages 1–11, 2019.
26. Yen Nhi Truong Vu, Richard Wang, Niranjana Balachandar, Can Liu, Andrew Y. Ng, and Pranav Rajpurkar. MedAug: Contrastive learning leveraging patient meta-data improves representations for chest X-ray interpretation. pages 1–14, 2021.
27. Qizhe Xie, Zihang Dai, Eduard Hovy, Minh-thang Luong, and Quoc V Le. Unsupervised Data Augmentation for Consistency Training. (NeurIPS):1–20, 2020.
28. Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, 2017-Janua:5987–5995, 2017.
29. Yingpeng Xie, Qiwei Wan, Guozhen Chen, Yanwu Xu, and Baiying Lei B. *Retinopathy Diagnosis Using Generative Adversarial Network*, volume 1. Springer International Publishing.