

---

# How Transferable Are Self-supervised Features in Medical Image Classification Tasks?

---

**Tuan Truong**  
Bayer AG  
Leverkusen, Germany  
tuan.truong@bayer.com

**Sadegh Mohammadi**  
Bayer AG  
Leverkusen, Germany  
sadegh.mohammadi@bayer.com

**Matthias Lenga**  
Bayer AG  
Berlin, Germany  
matthias.lenga@bayer.com

## Abstract

Transfer learning has become a standard practice to mitigate the lack of labeled data in medical classification tasks. Whereas finetuning a downstream task using supervised ImageNet pretrained features is straightforward and extensively investigated in many works, there is little study on the usefulness of self-supervised pretraining. In this paper, we assess the transferability of ImageNet self-supervised pretraining by evaluating the performance of models initialized with pretrained features from three self-supervised techniques (SimCLR, SwAV, and DINO) on selected medical classification tasks. The chosen tasks cover tumor detection in sentinel axillary lymph node images, diabetic retinopathy classification in fundus images, and multiple pathological condition classification in chest X-ray images. We demonstrate that self-supervised pretrained models yield richer embeddings than their supervised counterpart, which benefits downstream tasks in view of both linear evaluation and finetuning. For example, in view of linear evaluation at a critically small subset of the data, we see an improvement up to 14.79% in Kappa score in the diabetic retinopathy classification task, 5.4% in AUC in the tumor classification task, 7.03% AUC in the pneumonia detection, and 9.4% in AUC in the detection of pathological conditions in chest X-ray. In addition, we introduce *Dynamic Visual Meta-Embedding* (DVME) as an end-to-end transfer learning approach that fuses pretrained embeddings from multiple models. We show that the collective representation obtained by DVME leads to a significant improvement in the performance of selected tasks compared to using a single pretrained model approach and can be generalized to any combination of pretrained models.

## 1 Introduction

### 1.1 Background and Motivation

The scarcity of high-quality annotated data remains a notorious challenge in medical image analysis due to the high cost of acquiring expert annotations [1]. Transfer learning from large models pretrained in a supervised fashion on natural images such as ImageNet has become a *de-facto* solution for 2D medical imaging tasks in low data regimes [2–5]. Recently, self-supervised learning shows initial success in building large-scale Deep Learning based applications by leveraging unannotated data for pretraining [6–13]. However, a bottleneck within self-supervised learning is the demanding requirement of computational resources to train compared to standard supervised learning [14–17].

For example, regarding training on ImageNet, SwAV [15] uses the batch size of 4096 distributed on 64 GPUs and SimCLR [14] uses varying batch sizes between 256 and 8192 on 32-128 TPU cores. Even when the batch size is small, the author of SimCLR notes that the training time must be extended to provide more negative examples. In pretraining medical datasets, Azizi et al. [18] observe the best performance when using the batch size of 1024 on 64 cloud TPU cores to train SimCLR on a chest X-ray dataset. While transfer learning from supervised pretraining on a large labeled dataset such as ImageNet is widely studied [19, 20], the transferability of models pretrained on ImageNet using self-supervised techniques requires further investigation.

This paper reflects on the effectiveness of transfer learning with self-supervised features. We evaluate the performance of four downstream classification tasks using ImageNet pretrained features obtained from supervised and self-supervised techniques. The four distinct tasks concern three modalities with varying data sizes and distributions. The first task is in the domain of digital pathology and aims to detect sentinel axillary lymph node metastases in hematoxylin and eosin (H&E) stained patches extracted from whole-slide images. The second task concerns the severity classification of diabetic retinopathy from colored fundus images. The last two tasks are related to reading X-ray images. One involves identifying whether a patient is suffering from pneumonia and the other involves detecting multiple findings, such as pneumothorax, nodule or mass, opacity, and fracture. In particular, we consider low data regimes ranging from approximately 1% to 10% of the original dataset size for each task (Section 5.2). We evaluate pretrained features of three self-supervised approaches, SimCLR [14], SwAV [15], and DINO [16], on aforementioned tasks by training a linear layer on top of frozen features. We find that DINO consistently outperforms other self-supervised techniques and the supervised baseline by a significant margin.

Additionally, we propose *Dynamic Visual Meta-Embeddings* (DVME) - a model-agnostic technique to combine multiple self-supervised pretrained features for downstream tasks. In natural language processing, it has been observed that different word embeddings work well for different tasks and that it is difficult to anticipate the usefulness of a given embedding technique for a certain task at hand. The usage of a meta-embedding mitigates this problem by constructing an ensemble of embedding sets to increase the lexical coverage of vocabulary which leads to improved performance on downstream tasks [21]. Similarly, in vision tasks, we propose to concatenate multiple pretrained embeddings with self-attention for transfer learning. Concatenation expands the embedding space and yields richer representation while self-attention adapts the contribution of individual embedding to a specific downstream task. We show that DVME leads to a further increase in performance across all tasks compared to the best single self-supervised pretrained model baseline.

## 1.2 Contributions

Overall, the main contributions are as following:

- Across four distinct medical image classification tasks, we assess the quality of the embeddings obtained from different models which are pretrained on ImageNet using state-of-the-art self-supervised or supervised pretraining techniques.
- We identify a single self-supervised model which consistently outperforms the other approaches on all selected downstream tasks. In particular, this effect is prominently observed in low data regimes.
- We propose Dynamic Visual Meta-Embeddings (DVME) to fully leverage the collective representations obtained from different self-supervised pretrained models. The representations obtained from the DVME model aggregation outperform all single model approaches on the selected downstream tasks.

## 2 Related work

**Self-supervised learning in medical imaging** Two main self-supervised approaches in medical imaging are in the form of *handcrafted pretext tasks* and *contrastive learning*. Early applications design tailored pretext tasks to reconstruct images from transformed or distorted inputs [6, 7, 9–12, 22]. For example, Model Genesis [6] applies in-domain transfer learning to various classification and segmentation tasks on CT and X-ray images. The proposed architecture is an autoencoder that reconstructs images from four transformations, namely non-linear, local-shuffling, out-painting, and

in-painting. The induced transformations are supposed to enable the encoder to learn features related to appearance, texture, and context. Chen et al. [7] propose context restoration as a pretext task applied in three common medical use cases: plane detection on fetal 2D ultrasound images, abdominal organ localization on CT images, and brain tumor segmentation on MRI images. The proposed method generates distorted images with different spatial contexts while maintaining the same intensity distribution by repeatedly swapping two random patches in an input image. Through reconstruction, the model learns useful semantic features transferable in subsequent target classification and segmentation tasks. Alternatively, several works [9, 10] tailor the pretext tasks as solving Jigsaw puzzles and Rubik cubes. Taleb et al. [9] create puzzles made of patches fused from different modalities, e.g., different MRI modes, of the same structure and trains a construction task to reassemble the shuffled patches. Likewise, in 3D images, Zhuang et al. [10] rearrange and rotate the CT volumes, driving the model to learn features invariant to translation and rotation. The pretrained features are transferred to solve brain hemorrhage classification and tumor segmentation tasks. The limitation of handcrafted pretext tasks is that they are highly task- and domain-specific, and thus cannot generalize well to different tasks. Lately, contrastive learning-based techniques (see Section 3) resolve this issue. Sowrirajan et al. [13] use MoCo [23] to pretrain on unlabeled CheXpert [24] dataset and finetunes with labels on external Shenzhen Hospital X-ray dataset [25] to detect pleural effusion. Azizi et al. [18] propose multi-instance contrastive learning that utilizes two crops of the same patient but with different viewpoints or lighting conditions as positive pairs so that the model learns representations invariant to different images of the same pathology.

**Transfer learning in medical imaging** Transfer learning with ImageNet pretrained features still incites debates over its actual benefits for downstream medical tasks [19, 20, 26]. In a large data regime, Raghu et al. [19] show that lightweight models with random initialization can perform on par with large architectures pretrained on ImageNet such as ResNet-50 [27] and Inception-v3 [28]. On the contrary, Ke et al. [20] argue that ImageNet pretraining can significantly boost the performance with newer architectures such as DenseNet [29] and EfficientNet [30]. In low data regimes, however, transfer learning with self-supervised approaches has been found particularly helpful in recent works [18, 31, 32]. Azizi et al. [18] perform transfer learning with SimCLR [14] on X-ray and dermatology datasets and show a significant gain compared to a supervised baseline. Chaves et al. [32] evaluate self-supervised models on multiple dermatology datasets and find the advantage of self-supervised pretraining when using low training data. Whereas prior works focus on a single self-supervised technique [18] and a unique modality, i.e., dermatology [32], our work extends the investigation by benchmarking various self-supervised approaches against the supervised baseline across a set of heterogeneous medical imaging tasks. Our primary goal is to compare the richness of feature embeddings of different self-supervised learning techniques in the scope of transfer learning on medical imaging tasks.

### 3 Contrastive Learning

Self-supervised methods can be broadly categorized into handcrafted pretext tasks [33] and contrastive learning. In Layman’s terms, contrastive learning learns the representation by comparing the similarity between images. Given the output embeddings obtained from an encoder, they are either pulled closer (similar) or pushed away (dissimilar) in the embedding space. Most of the approaches are built on the notion of *multi-instance level classification* [34], which considers each image as a unique class. Current techniques create so-called positive pairs as data augmented versions of an image and the model learns to contrast them from the rest of images in the batch. A detailed review and taxonomy of contrastive learning can be found in [35]. Among state-of-the-art contrastive techniques, our main focuses are SimCLR, SwAV, and DINO.

**SimCLR** *Simple Framework for Contrastive Learning of Visual Representation* [14] maximizes the agreement of two views from the same image. The paper proposes a set of transformations applied to input images to create positive and negative pairs. An encoder takes a transformed batch and forwards it to a projection head that maps images to an embedding space. A contrastive loss on top compares the embeddings to minimize the distance between similar (positive) embeddings. Finally, the projection head is discarded and the encoder can be transferred to downstream tasks.

**SwAV** *Swapping Assignments between multiple Views of the same image* [15] also contrasts two image views but not in a direct, sample-based fashion as SimCLR. Instead, it compares the cluster to which each view belongs. If two views come from the same image, they should fall on the same cluster assignment and vice versa. Caron et al. [15] shows that this approach has an advantage over SimCLR in avoiding the need for large batch size and improving the convergence time. In comparison to a prior clustering-based self-supervised technique in [36], the clustering assignment process is online, so gradients can be backpropagated in batch-wise manner.

**DINO** *Knowledge distillation without labels* [16] matches the output probability distributions of various image views obtained from two networks. This approach takes inspiration from Bootstrap Your Own Latent (BYOL) [17] in the perspective of self distillation task and the architecture of Vision Transformer (ViT) [37] as the backbone. Instead of passing the views into the same network, DINO passes two transformations of an image into two networks, namely the student and teacher network. The loss compares the probability outputs of both networks and the student’s parameters are updated via backpropagation while the teacher’s parameters is updated via an exponential moving average of the student ones. In addition, compared to using convolutional architectures, Caron et al.’s study [16] indicates that ViT-based DINO shows distinct properties in characterizing object boundaries and generates features that perform well using K-Nearest Neighbors without further finetuning in ImageNet classification task.

## 4 Datasets

The four datasets in our experiments are distinct in terms of modality, dataset size, and class distribution in order to partially reflect the heterogeneity of typical medical imaging tasks. We consider three common modalities in medical image analysis: digital pathology, fundus imaging, and X-ray.

**PatchCamelyon (PatchCam)** The dataset contains H&E sentinel axillary lymph node patches extracted from whole-slide image in the study at [38, 39]. All of the slides were annotated by expert pathologists. If the center of a patch contains at least one pixel of tumor tissue, it will be marked as positive. The data version we use is the curated one from Kaggle competition<sup>1</sup> that removes all the duplicated patches and comes with a default train/test split. The original train set consists of 220025 patches of size  $96 \times 96$  with binary labels indicating whether there is a tumor or not. For our training task, we randomly select a subset comprising 50000 images. The official test set comprises 57486 images for which no labels are provided. Hence, for all performance evaluations on PatchCam we therefore submit our predictions to Kaggle.

**APTOS** The dataset comprises of colored fundus images of diabetic retinopathy patients obtained from diverse clinics with different camera setups. For each image, clinicians grade the severity with a score between 0 and 4 that corresponds to no diabetic retinopathy, mild, moderate, severe, and proliferative diabetic retinopathy. The dataset is part of the challenge held on Kaggle<sup>2</sup>. The train and test set contain 3662 and 1928 images respectively. Similar to PatchCam, we submit the inference of the test set to Kaggle to obtain the scores.

**Pneumonina chest X-ray** The dataset contains chest X-ray images annotated by two expert radiologists. Each radiologist classifies each image into healthy and pneumonia. We obtain the dataset in Kaggle<sup>3</sup> with a default train/test split of 5216/624 images. Each patient can have multiple images, which is taken into consideration for patient stratification during train and validation. Further description of the dataset and acquisition can be found at [40].

**NIH chest X-ray** The dataset consists of chest X-ray images provided by NIH Clinical Center<sup>4</sup>. The images were manually reviewed by three certified radiologists from the American Board of Radiology. Each radiologist marks the presence of four medical conditions: pneumothorax, nodule

<sup>1</sup><https://www.kaggle.com/c/histopathologic-cancer-detection>

<sup>2</sup><https://www.kaggle.com/c/aptos2019-blindness-detection>

<sup>3</sup><https://www.kaggle.com/paultimothymooney/chest-xray-pneumonia>

<sup>4</sup><https://nihcc.app.box.com/v/ChestXray-NIHCC>

or mass, opacity, and fracture. We use two subsets of the original NIH Chest X-ray dataset which are referred as validation and test set in the study at [41]. We use the first subset (2414 images) for training and keep the second subset (1962 images) for evaluation. Since there can be multiple findings per image, for simplicity we exclude such cases in our experiment. In addition, we also add a class called *other* for when no mentioned conditions are found in the image.

## 5 Experimental Setup

### 5.1 Architecture

To assess the transferability of self-supervised features in medical imaging tasks, we choose ResNet-50 [27], which is the backbone for many state-of-the-art self-supervised approaches [14–16, 23]. ResNet-50 is used for the supervised baseline (ImageNet and random initialization) and two self-supervised approaches (SwAV and SimCLR). For Dino, the architecture of choice is VisionTransformer (ViT) with patch size  $8 \times 8$ . The pretrained weights for SimCLR and SwAV are obtained from VISSL<sup>5</sup> while for DINO the weights are obtained from the FAIR repository<sup>6</sup>.

### 5.2 Dataset sizes and subtasks

For each dataset we consider three different subtasks (Small, Medium, Full) each containing a different fraction of training data, see Table 1. The subtasks are generated by random sampling which is accounting for class imbalances. For each subtask, we generate a split of the training data which is used to conduct a five fold cross validation. Finally, the five models from this cross validation are evaluated on the fixed test set and the final performance is obtained by averaging the individual model scores. For PatchCam and APTOS, we submit the prediction of test set to Kaggle and obtain a final score of which the calculation is presented in Section 5.6.

Table 1: Number of samples for different subtasks

Dataset	Training			Testing
	Small (S)	Medium (M)	Full (F)	
PatchCam	500	5000	50000	57486 (*)
APTOS	50	500	3662	1928 (*)
Pneumonia Chest X-ray	50	500	5216	624
NIH Chest X-ray	20	200	2414	1962

(\*) The performance evaluation is obtained by submitting the predictions to Kaggle.

### 5.3 Linear performance and finetuning

The evaluation of *linear performance* is of the standard methods for assessing feature quality of a pretrained feature extractor. This performance measure is obtained by freezing all layers of the pretrained model and only finetuning a final linear layer, cf. [14–16, 18]. In our experiments we follow the same setup by adding and training linear layers with output dimensions adjusted to the number of classes in the respective downstream tasks. The linear layer is added after the last average pooling layer in ResNet-50 and after the concatenation of class tokens from the last four blocks in ViT, following the implementation of [16]. In addition to linear evaluation we consider *finetuning*, where all layers of the pretrained base network as well as the final linear classifier are adapted on the downstream task at hand. To avoid bias learning due to class imbalance, the number of samples per class is maintained balanced across each experiment. However, as the number of samples in the NIH Chest X-ray dataset for the underrepresented class is significantly small compared to other classes, we use oversampling during training.

<sup>5</sup><https://vissl.ai/>

<sup>6</sup><https://github.com/facebookresearch/dino>

## 5.4 Linear evaluation with DVME

Given a set of pretrained feature extractors, it may be difficult to anticipate which pretrained model to choose for a given downstream task at hand. This concern is also shared in natural language processing where there are multiple word embedding techniques trained on different domains, each having its own strengths depending on target tasks. Meta-embedding is an effective technique that takes a union over different word embeddings to tackle the out-of-vocabulary problem and fuse multi-modal information [21]. Though there is no multi-modal information in our study, we hypothesize that the pretrained features from different techniques are sufficiently independent from one another and encode certain complementary information. Therefore, we propose Dynamic Visual Meta-Embedding (DVME) for vision tasks which aggregates information by concatenating the embeddings of SimCLR, SwAV and DINO pretrained models. The newly constructed embedding space improves the separability of image features through the complementary effect of each embedding component. For SimCLR and SwAV, the embedding is obtained at the last layer before the classifier and has the dimension of 2048. For DINO, the embedding is obtained by concatenating the class tokens of the last four blocks, resulting in the length of 1536. We project each embedding to a fixed dimension of 512 and concatenate them. A self-attention module is added after the concatenated layer to learn the importance of each embedding component for a specific downstream task. This self-attention module is the same as in the Vision Transformer architecture [37] except that the attention is learned across different components of the meta-embedding instead of image patches. Figure 1 shows a sketch how self-attention is incorporated when fusing the pretrained features. Clearly, the proposed DVME approach is not limited to SimCLR, SwAV or DINO and can be used with other feature extractors. We provide a snippet of DVME implemented in PyTorch in Appendix E.

## 5.5 Hyperparameters and Augmentation

All experiments use the Adam optimizer starting with a small learning rate between  $1e-3$  and  $1e-4$  and further reducing it when the validation loss does not improve consecutively over five epochs. As our study aims to compare among different initialization and not to outperform the best performance, we do not apply intensive augmentation techniques. Images with an original size larger than  $224 \times 224$  are resized into  $256 \times 256$ , then cropped and applied further flipping or rotation depending on the nature of modality. For PatchCam dataset, we apply directly flipping without resizing or cropping as the size of the images is less than  $224 \times 224$ .

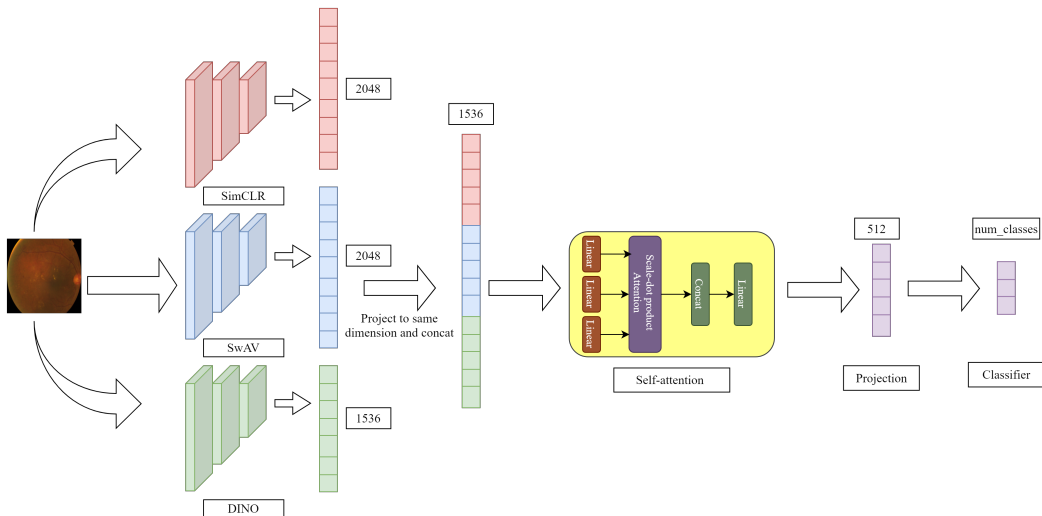


Figure 1: Dynamic Visual Meta-Embedding (DVME): The embeddings extracted from each pretrained model are projected to the same dimension and concatenated before feeding to the self-attention module.

## 5.6 Metrics

The evaluation metric for PatchCam, Pneumonia, and NIH Chest X-ray is the area under the Receiver Operating Characteristic curve (AUC) while it is the Cohen Kappa score for APTOS. We submit the predictions of APTOS and PatchCam test set to Kaggle and obtain two scores for the private and public leaderboard, which are evaluated on two different portions of the test set. We calculate the final score as the weighted average score of the private and public leaderboard. Precisely, the final score is calculated as  $s_{avg} = \alpha \times s_{private} + (1 - \alpha) \times s_{public}$  where  $\alpha$  for PatchCam is 0.51 and for APTOS is 0.85. The value of  $\alpha$  is the portion of test set that Kaggle uses to calculate the score on private leaderboard and it varies depending on the competition.

## 6 Results and Discussion

### 6.1 Evaluation of self-supervised and supervised pretrained features

For each task and each self-supervised pretrained initialization, Table 2 contains the linear evaluation performance obtained on small (S), medium (M), and full (F) sized subsets. As additional baselines we compare against the linear performance of features extracted by a ResNet-50 supervised pretrained on ImageNet and the features extracted by a Kaiming randomly initialized ResNet-50. The latter constitutes a lower bound for all other methods. In first order, the linear performance of the differently initialized feature extractors measures the separability of the dataset embeddings by a hyperplane.

Table 2 shows that SwAV and SimCLR pretrained features yield inconsistent patterns across all downstream tasks. For example, while SwAV and SimCLR initializations perform on par with each other on PatchCam and NIH Chest X-ray, they are different by approximately 10% in Kappa score and 3.7% in AUC on the S subsets of APTOS and Pneumonia Chest X-ray, respectively. Notably, DINO initialization consistently outperforms all the other initializations across all tasks by a significant margin. For example on NIH Chest X-ray S and M subtasks, DINO pretrained features yield an improvement of approximately 5-6% in AUC over SimCLR and SwAV. The single exception is the APTOS S subtask, where SwAV outperforms DINO by 3.3% in AUC. However, DINO still yields an improvement over SimCLR and ImageNet supervised initialization by 7% and 11.2% in Kappa score, respectively. We refer to Appendix B.1 for more detailed results on the performance obtained by the competing initialization methods for different dataset sizes. In comparison to ImageNet supervised pretrained features, we observe that self-supervised features are more effective in improving the performance across all downstream tasks. This suggests that the representation generated by self-supervised methods are of higher quality, leading to better performance on the test set and reducing the performance variability between folds in low data regimes, similar to the observation made in [32]. Figure 2 (a,b,d,e) visualizes the t-SNE embeddings of the features extracted by the supervised pretrained ResNet-50 and DINO on the PatchCam (binary classification) and APTOS (multi-class) downstream tasks. In Figure 2 (a,b) an improved separation in t-SNE space can be seen clearly when comparing the supervised pretrained embedding to the DINO embedding. We refer to Appendix D for the embedding visualization of other datasets.

We extend our comparison by finetuning all different model initializations separately on all downstream tasks. As our study primarily concerns assessing the quality of different pretrained features, no extensive hyperparameter tuning is conducted in the finetuning setting. Table 3 summarizes the finetuning results across all datasets and each subtasks. We refer to Appendix B.2 for more detailed results on different subset sizes as well as hyperparameters. Similar to the linear evaluation results, we consistently observe a higher performance for all self-supervised pretrained initializations compared to the supervised pretrained and randomly initialized baselines in the low data regimes (S, M), which supports the observation made by [18, 32]. DINO pretrained features outperform those from other self-supervised methods in 2/4 S subtasks and 3/4 M subtasks. When using full data for fine-tuning, all self-supervised pretrained initializations consistently outperform the baseline methods on PatchCam, APTOS and Pneumonia Chest X-ray. On full data of the NIH Chest X-ray task, only SwAV exceeds the supervised baseline performance.

Moreover, we note that finetuning can result in slightly lower performance than linear evaluation in some of the initializations when the number of training samples falls less than a hundred due to overfitting effects. This can be observed for the S subtasks of APTOS, Pneumonia Chest X-ray, and NIH Chest X-ray dataset where there are 50 samples or less, see Table 1. For example, while

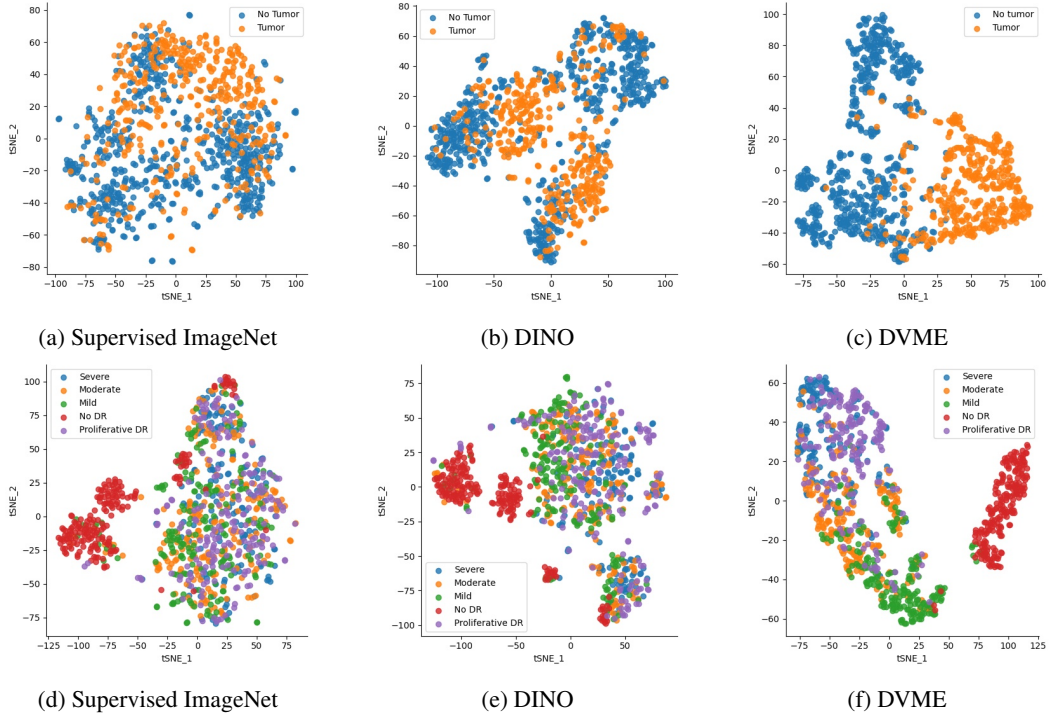


Figure 2: t-SNE visualization of embeddings obtained using different pretrained feature extractors (supervised ImageNet, DINO, proposed method DVME). Top row: **PatchCam** dataset, bottom row: **APTOS** dataset

DINO achieves the highest AUC of 0.6323 on the S subtask of NIH Chest X-ray in linear evaluation, the best performance for finetuning is obtained by SwAV and reaches only 0.5903 AUC.

## 6.2 Evaluation of DVME performance

Table 4 summarizes the results obtained from fusing SwAV, SimCLR and DINO with the DVME approach on each dataset and for each subtask. We evaluate DVME similar to the linear performance evaluation from Section 6.1 by training only the meta-embedding on top of the three frozen feature extractors, cf. Section 5.4. As a performance benchmark, we selected for each dataset and each fraction of data the self-supervised initialization that leads to the best linear evaluation performance, cf. Table 2. DVME outperforms this benchmark in 4/4 of the S subtasks, in 3/4 of the M subtasks and in 2/4 F subtasks. For the subtasks where DVME is not exceeding the benchmark performance, the performance difference lies within one standard deviation of the DVME linear evaluation score. The improvement of DVME over the benchmark is particularly pronounced for the APTOS and NIH Chest X-ray tasks. For example, DVME helps gain roughly 6% in Kappa score over the best individual baseline for the S and M subtask of the APTOS dataset.

The t-SNE visualizations of the DVME embeddings in Figure 2 (c,f) qualitatively indicate that the clusters are better separated, particularly in the case of multiclass classification. When analyzing the attention matrix, we find that SwAV and SimCLR pay little attention to each other but firmly into DINO, suggesting the representation from SwAV and SimCLR could be more similar and thus not so informative compared to DINO.

In order to better understand the effect of self-attention on the model fusion, we conduct an ablation study on DVME in Appendix C. In the setting without self-attention, the meta-embedding is directly connected to the linear classifier. Without self-attention, the feature fusion still yields a significant improvement over the baseline, which supports our hypothesis in Section 5.4 that each embedding contains complementary information. However, self-attention is particularly beneficial to certain tasks. For example, on APTOS the Kappa scores increase by 5.6%, 4.4% and 4.8% for the S, M and F subtask, respectively.



Table 2: Linear evaluation performance of different self-supervised initializations, supervised pre-training and random initialization on different scales (small, medium, full) of the data.

Dataset	Method	Small (S)	Medium (M)	Full (F)
PatchCam	Random	0.6594 ± 0.0319	0.6994 ± 0.0079	0.7990 ± 0.0021
	Supervised ImageNet	0.7517 ± 0.0136	0.7863 ± 0.0063	0.7975 ± 0.0032
	SwAV	0.7834 ± 0.0112	0.8043 ± 0.0072	0.8088 ± 0.0025
	SimCLR	0.7895 ± 0.0091	0.8053 ± 0.0069	0.8084 ± 0.0026
	DINO	<b>0.8058 ± 0.0100</b>	<b>0.8359 ± 0.0053</b>	<b>0.8487 ± 0.0014</b>
APTOS (*)	Random	0.0324 ± 0.0602	0.0624 ± 0.0459	0.1550 ± 0.116
	Supervised ImageNet	0.4851 ± 0.0811	0.6822 ± 0.0257	0.7331 ± 0.0124
	SwAV	<b>0.6330 ± 0.0204</b>	0.7274 ± 0.0095	0.7617 ± 0.0128
	SimCLR	0.5305 ± 0.0539	0.6500 ± 0.0138	0.6989 ± 0.0084
	DINO	0.6003 ± 0.0691	<b>0.7372 ± 0.0167</b>	<b>0.7790 ± 0.0083</b>
Pneumonia Chest X-ray	Random	0.6899 ± 0.0339	0.8258 ± 0.0237	0.8907 ± 0.0144
	Supervised ImageNet	0.8789 ± 0.0234	0.8954 ± 0.0151	0.9397 ± 0.0033
	SwAV	0.8808 ± 0.0222	0.9215 ± 0.0252	0.9709 ± 0.0047
	SimCLR	0.9168 ± 0.0006	0.9346 ± 0.0072	0.9665 ± 0.0027
	DINO	<b>0.9492 ± 0.0170</b>	<b>0.9718 ± 0.0055</b>	<b>0.9868 ± 0.0008</b>
NIH Chest X-ray	Random	0.5212 ± 0.0344	0.5317 ± 0.0176	0.5392 ± 0.0346
	Supervised ImageNet	0.5383 ± 0.0392	0.6688 ± 0.0148	0.7109 ± 0.0084
	SwAV	0.5785 ± 0.0258	0.6889 ± 0.0089	0.7225 ± 0.0139
	SimCLR	0.5792 ± 0.0435	0.6645 ± 0.0067	0.6983 ± 0.0231
	DINO	<b>0.6323 ± 0.0131</b>	<b>0.7373 ± 0.0112</b>	<b>0.7438 ± 0.0228</b>

(\*) The evaluation metric for APTOS is Cohen-Kappa score while for others is AUC score.

Table 3: Finetuning performance of different self-supervised initializations, supervised pretraining and random initialization on different scales (small, medium, full) of the data.

Dataset	Method	Small (S)	Medium (M)	Full (F)
PatchCam	Random	0.7355 ± 0.0282	0.7660 ± 0.0223	0.8515 ± 0.0023
	Supervised ImageNet	0.7897 ± 0.0162	0.8274 ± 0.0051	0.8483 ± 0.0097
	SwAV	0.7895 ± 0.0336	0.8399 ± 0.0142	<b>0.8619 ± 0.0090</b>
	SimCLR	0.8021 ± 0.0138	0.8329 ± 0.0085	0.8553 ± 0.0110
	DINO	<b>0.8366 ± 0.0092</b>	<b>0.8440 ± 0.0172</b>	0.8517 ± 0.0158
APTOS (*)	Random	0.0177 ± 0.0954	0.3233 ± 0.0822	0.5927 ± 0.0545
	Supervised ImageNet	0.4817 ± 0.0991	0.7369 ± 0.0310	0.8057 ± 0.0149
	SwAV	0.4928 ± 0.0378	0.7594 ± 0.0246	0.8293 ± 0.0133
	SimCLR	0.5916 ± 0.0570	0.7603 ± 0.0249	0.8264 ± 0.0103
	DINO	<b>0.6601 ± 0.0447</b>	<b>0.7945 ± 0.0079</b>	<b>0.8365 ± 0.0213</b>
Pneumonia Chest X-ray	Random	0.6895 ± 0.0512	0.9183 ± 0.0186	0.9820 ± 0.0043
	Supervised ImageNet	0.8649 ± 0.0442	0.9698 ± 0.0066	0.9910 ± 0.0015
	SwAV	<b>0.9289 ± 0.0291</b>	0.9814 ± 0.0087	0.9927 ± 0.0016
	SimCLR	0.9197 ± 0.0168	0.9781 ± 0.0085	<b>0.9950 ± 0.0013</b>
	DINO	0.9256 ± 0.0235	<b>0.9867 ± 0.0051</b>	0.9948 ± 0.0010
NIH Chest X-ray	Random	0.5015 ± 0.0253	0.6404 ± 0.0165	0.6616 ± 0.0345
	Supervised ImageNet	0.5251 ± 0.0238	0.6816 ± 0.0429	0.7618 ± 0.0116
	SwAV	<b>0.5903 ± 0.0384</b>	0.6973 ± 0.0227	<b>0.7737 ± 0.0212</b>
	SimCLR	0.5570 ± 0.0450	<b>0.7228 ± 0.0287</b>	0.7358 ± 0.0295
	DINO	0.5552 ± 0.0546	0.6652 ± 0.0114	0.7404 ± 0.0240

(\*) The evaluation metric for APTOS is Cohen-Kappa score while for others is AUC score.

Table 4: Linear evaluation performance of Dynamic Visual Meta-Embedding (DVME) in comparison with the best score obtained using a single pretrained model on different downstream tasks on different scales (small, medium, full) of the data.

Dataset	Method	Small (S)	Medium (M)	Full (F)
PatchCam	DVEM	<b>0.8227 ± 0.0148</b>	<b>0.8399 ± 0.0059</b>	0.8467 ± 0.0094
	Best single baseline	0.8058 ± 0.0100	0.8359 ± 0.0100	<b>0.8487 ± 0.0014</b>
APTOS (*)	DVME	<b>0.6913 ± 0.0575</b>	<b>0.7925 ± 0.0265</b>	<b>0.8242 ± 0.0279</b>
	Best single baseline	0.6330 ± 0.0204	0.7372 ± 0.0167	0.7790 ± 0.0083
Pneumonia Chest X-ray	DVME	<b>0.9539 ± 0.0025</b>	0.9696 ± 0.0101	0.9842 ± 0.0029
	Best single baseline	0.9492 ± 0.0170	<b>0.9718 ± 0.0055</b>	<b>0.9868 ± 0.0008</b>
NIH Chest X-ray	DVME	<b>0.6566 ± 0.0564</b>	<b>0.7601 ± 0.0146</b>	<b>0.7538 ± 0.0234</b>
	Best single baseline	0.6323 ± 0.0131	0.7373 ± 0.0112	0.7438 ± 0.0228

(\*) The evaluation metric for APTOS is Cohen-Kappa score while for others is AUC score.

## 7 Conclusion

In this study, we assess the quality of ImageNet self-supervised pretrained features in four selected medical image classification tasks. We demonstrate that feature extractors which were pretrained using SwAV, SimCLR or DINO consistently yield richer embeddings on the downstream tasks compared to a supervised pretrained baseline model. Among all self-supervised techniques, DINO outperforms the other methods on the majority of datasets and subtasks. Furthermore, we show that the representations from each individual pretrained model encode complementary information which can be fused to yield even more meaningful features. To that end we propose Dynamic Visual Meta-Embedding (DVME), a model-agnostic meta-embedding approach. Our experiments indicate that DVME outperforms the best single model baseline on numerous tasks. As a model-agnostic approach, DVME is not limited to SwAV, SimCLR or DINO. With slight modifications other models can be combined using DVME to generate enriched representations.

## References

- [1] Daniel C Castro, Ian Walker, and Ben Glocker. Causality matters in medical imaging. *Nature Communications*, 11(1):1–10, 2020.
- [2] Carson Lam, Darwin Yi, Margaret Guo, and Tony Lindsey. Automated detection of diabetic retinopathy using deep learning. *AMIA summits on translational science proceedings*, 2018: 147, 2018.
- [3] Neslihan Bayramoglu and Janne Heikkilä. Transfer learning for cell nuclei classification in histopathology images. In *European Conference on Computer Vision*, pages 532–539. Springer, 2016.
- [4] Bens Pardamean, Tjeng Wawan Cenggoro, Reza Rahutomo, Arif Budiarto, and Ettikan Kandasamy Karupiah. Transfer learning from chest x-ray pre-trained convolutional neural network for learning mammogram data. *Procedia Computer Science*, 135:400–407, 2018.
- [5] Yang Yang, Lin-Feng Yan, Xin Zhang, Yu Han, Hai-Yan Nan, Yu-Chuan Hu, Bo Hu, Song-Lin Yan, Jin Zhang, Dong-Liang Cheng, et al. Glioma grading on conventional mr images: a deep learning study with transfer learning. *Frontiers in neuroscience*, 12:804, 2018.
- [6] Zongwei Zhou, Vatsal Sodha, Md Mahfuzur Rahman Siddiquee, Ruibin Feng, Nima Tajbakhsh, Michael B. Gotway, and Jianming Liang. Models Genesis: Generic Autodidactic Models for 3D Medical Image Analysis. *arXiv:1908.06912 [cs, eess]*, August 2019. URL <http://arxiv.org/abs/1908.06912>. arXiv: 1908.06912.
- [7] Liang Chen, Paul Bentley, Kensaku Mori, Kazunari Misawa, Michitaka Fujiwara, and Daniel Rueckert. Self-supervised learning for medical image analysis using image context restoration. *Medical Image Analysis*, 58:101539, December 2019. ISSN 13618415. doi:

- 10.1016/j.media.2019.101539. URL <https://linkinghub.elsevier.com/retrieve/pii/S1361841518304699>.
- [8] Aiham Taleb, Winfried Loetzsch, Noel Danz, Julius Severin, Thomas Gaertner, Benjamin Bergner, and Christoph Lippert. 3D Self-Supervised Methods for Medical Imaging. *arXiv:2006.03829 [cs, eess]*, November 2020. URL <http://arxiv.org/abs/2006.03829>. arXiv: 2006.03829.
- [9] Aiham Taleb, Christoph Lippert, Tassilo Klein, and Moin Nabi. Multimodal Self-Supervised Learning for Medical Image Analysis. *arXiv:1912.05396 [cs]*, October 2020. URL <http://arxiv.org/abs/1912.05396>. arXiv: 1912.05396.
- [10] Xinrui Zhuang, Yuexiang Li, Yifan Hu, Kai Ma, Yujiu Yang, and Yefeng Zheng. Self-supervised Feature Learning for 3D Medical Images by Playing a Rubik’s Cube. In Dinggang Shen, Tianming Liu, Terry M. Peters, Lawrence H. Staib, Caroline Essert, Sean Zhou, Pew-Thian Yap, and Ali Khan, editors, *Medical Image Computing and Computer Assisted Intervention – MICCAI 2019*, volume 11767, pages 420–428. Springer International Publishing, Cham, 2019. ISBN 978-3-030-32250-2 978-3-030-32251-9. doi: 10.1007/978-3-030-32251-9\_46. URL [http://link.springer.com/10.1007/978-3-030-32251-9\\_46](http://link.springer.com/10.1007/978-3-030-32251-9_46). Series Title: Lecture Notes in Computer Science.
- [11] Wenjia Bai, Chen Chen, Giacomo Tarroni, Jinming Duan, Florian Guitton, Steffen E. Petersen, Yike Guo, Paul M. Matthews, and Daniel Rueckert. Self-Supervised Learning for Cardiac MR Image Segmentation by Anatomical Position Prediction. In Dinggang Shen, Tianming Liu, Terry M. Peters, Lawrence H. Staib, Caroline Essert, Sean Zhou, Pew-Thian Yap, and Ali Khan, editors, *Medical Image Computing and Computer Assisted Intervention – MICCAI 2019*, volume 11765, pages 541–549. Springer International Publishing, Cham, 2019. ISBN 978-3-030-32244-1 978-3-030-32245-8. doi: 10.1007/978-3-030-32245-8\_60. URL [http://link.springer.com/10.1007/978-3-030-32245-8\\_60](http://link.springer.com/10.1007/978-3-030-32245-8_60). Series Title: Lecture Notes in Computer Science.
- [12] Christian Abbet, Inti Zlobec, Behzad Bozorgtabar, and Jean-Philippe Thiran. Divide-and-Rule: Self-Supervised Learning for Survival Analysis in Colorectal Cancer. In Anne L. Martel, Purang Abolmaesumi, Danail Stoyanov, Diana Mateus, Maria A. Zuluaga, S. Kevin Zhou, Daniel Racoceanu, and Leo Joskowicz, editors, *Medical Image Computing and Computer Assisted Intervention – MICCAI 2020*, volume 12265, pages 480–489. Springer International Publishing, Cham, 2020. ISBN 978-3-030-59721-4 978-3-030-59722-1. doi: 10.1007/978-3-030-59722-1\_46. URL [http://link.springer.com/10.1007/978-3-030-59722-1\\_46](http://link.springer.com/10.1007/978-3-030-59722-1_46). Series Title: Lecture Notes in Computer Science.
- [13] Hari Sowrirajan, Jingbo Yang, Andrew Y. Ng, and Pranav Rajpurkar. MoCo-CXR: MoCo Pre-training Improves Representation and Transferability of Chest X-ray Models. *arXiv:2010.05352 [cs]*, February 2021. URL <http://arxiv.org/abs/2010.05352>. arXiv: 2010.05352.
- [14] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A Simple Framework for Contrastive Learning of Visual Representations. *arXiv:2002.05709 [cs, stat]*, June 2020. URL <http://arxiv.org/abs/2002.05709>. arXiv: 2002.05709.
- [15] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised Learning of Visual Features by Contrasting Cluster Assignments. *arXiv:2006.09882 [cs]*, January 2021. URL <http://arxiv.org/abs/2006.09882>. arXiv: 2006.09882.
- [16] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging Properties in Self-Supervised Vision Transformers. *arXiv:2104.14294 [cs]*, April 2021. URL <http://arxiv.org/abs/2104.14294>. arXiv: 2104.14294.
- [17] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised Learning. *arXiv:2006.07733 [cs, stat]*, September 2020. URL <http://arxiv.org/abs/2006.07733>. arXiv: 2006.07733.
- [18] 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.

- arXiv:2101.05224 [cs, eess]*, January 2021. URL <http://arxiv.org/abs/2101.05224>. arXiv: 2101.05224.
- [19] Maithra Raghu, Chiyuan Zhang, Jon Kleinberg, and Samy Bengio. Transfusion: Understanding Transfer Learning for Medical Imaging. *arXiv:1902.07208 [cs, stat]*, October 2019. URL <http://arxiv.org/abs/1902.07208>. arXiv: 1902.07208.
- [20] Alexander Ke, William Ellsworth, Oishi Banerjee, Andrew Y. Ng, and Pranav Rajpurkar. CheXtransfer: Performance and Parameter Efficiency of ImageNet Models for Chest X-Ray Interpretation. *Proceedings of the Conference on Health, Inference, and Learning*, pages 116–124, April 2021. doi: 10.1145/3450439.3451867. URL <http://arxiv.org/abs/2101.06871>. arXiv: 2101.06871.
- [21] Douwe Kiela, Changhan Wang, and Kyunghyun Cho. Dynamic meta-embeddings for improved sentence representations, 2018.
- [22] Antoine Rivail, Ursula Schmidt-Erfurth, Wolf-Dieter Vogl, Sebastian M. Waldstein, Sophie Riedl, Christoph Grechenig, Zhichao Wu, and Hrvoje Bogunovic. Modeling Disease Progression in Retinal OCTs with Longitudinal Self-supervised Learning. In Islem Rekik, Ehsan Adeli, and Sang Hyun Park, editors, *Predictive Intelligence in Medicine*, volume 11843, pages 44–52. Springer International Publishing, Cham, 2019. ISBN 978-3-030-32280-9 978-3-030-32281-6. doi: 10.1007/978-3-030-32281-6\_5. URL [http://link.springer.com/10.1007/978-3-030-32281-6\\_5](http://link.springer.com/10.1007/978-3-030-32281-6_5). Series Title: Lecture Notes in Computer Science.
- [23] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum Contrast for Unsupervised Visual Representation Learning. *arXiv:1911.05722 [cs]*, March 2020. URL <http://arxiv.org/abs/1911.05722>. arXiv: 1911.05722.
- [24] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, Jayne Seekins, David A. Mong, Safwan S. Halabi, Jesse K. Sandberg, Ricky Jones, David B. Larson, Curtis P. Langlotz, Bhavik N. Patel, Matthew P. Lungren, and Andrew Y. Ng. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison, 2019.
- [25] Stefan Jaeger, Sema Candemir, Sameer Antani, Yi-Xiáng J. Wáng, Pu-Xuan Lu, and George Thoma. Two public chest x-ray datasets for computer-aided screening of pulmonary diseases. *Quantitative Imaging in Medicine and Surgery*, 4(6), 2014. ISSN 2223-4306. URL <https://qims.amegroups.com/article/view/5132>.
- [26] Kaiming He, Ross Girshick, and Piotr Dollár. Rethinking ImageNet Pre-training. *arXiv:1811.08883 [cs]*, November 2018. URL <http://arxiv.org/abs/1811.08883>. arXiv: 1811.08883.
- [27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.
- [28] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–9, 2015. doi: 10.1109/CVPR.2015.7298594.
- [29] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks, 2018.
- [30] Mingxing Tan and Quoc Le. EfficientNet: Rethinking model scaling for convolutional neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 6105–6114. PMLR, 09–15 Jun 2019. URL <http://proceedings.mlr.press/v97/tan19a.html>.
- [31] Alejandro Newell and Jia Deng. How Useful is Self-Supervised Pretraining for Visual Tasks? *arXiv:2003.14323 [cs]*, March 2020. URL <http://arxiv.org/abs/2003.14323>. arXiv: 2003.14323.
- [32] Levy Chaves, Alceu Bissoto, Eduardo Valle, and Sandra Avila. An Evaluation of Self-Supervised Pre-Training for Skin-Lesion Analysis. *arXiv:2106.09229 [cs]*, June 2021. URL <http://arxiv.org/abs/2106.09229>. arXiv: 2106.09229.

- [33] Longlong Jing and Yingli Tian. Self-supervised Visual Feature Learning with Deep Neural Networks: A Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2020. ISSN 0162-8828, 2160-9292, 1939-3539. doi: 10.1109/TPAMI.2020.2992393. URL <https://ieeexplore.ieee.org/document/9086055/>.
- [34] Alexey Dosovitskiy, Philipp Fischer, Jost Tobias Springenberg, Martin Riedmiller, and Thomas Brox. Discriminative unsupervised feature learning with exemplar convolutional neural networks, 2015.
- [35] Phuc H. Le-Khac, Graham Healy, and Alan F. Smeaton. Contrastive Representation Learning: A Framework and Review. *IEEE Access*, 8:193907–193934, 2020. ISSN 2169-3536. doi: 10.1109/ACCESS.2020.3031549. URL <https://ieeexplore.ieee.org/document/9226466/>.
- [36] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep Clustering for Unsupervised Learning of Visual Features. *arXiv:1807.05520 [cs]*, March 2019. URL <http://arxiv.org/abs/1807.05520>. arXiv: 1807.05520.
- [37] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *arXiv:2010.11929 [cs]*, June 2021. URL <http://arxiv.org/abs/2010.11929>. arXiv: 2010.11929.
- [38] Babak Ehteshami Bejnordi, Mitko Veta, Paul Johannes van Diest, Bram van Ginneken, Nico Karssemeijer, Geert Litjens, Jeroen A. W. M. van der Laak, , and the CAMELYON16 Consortium. Diagnostic Assessment of Deep Learning Algorithms for Detection of Lymph Node Metastases in Women With Breast Cancer. *JAMA*, 318(22):2199–2210, 12 2017. ISSN 0098-7484. doi: 10.1001/jama.2017.14585. URL <https://doi.org/10.1001/jama.2017.14585>.
- [39] Bastiaan S. Veeling, Jasper Linmans, Jim Winkens, Taco Cohen, and Max Welling. Rotation equivariant cnns for digital pathology, 2018.
- [40] Daniel S. Kermany, Michael Goldbaum, Wenjia Cai, Carolina C.S. Valentim, Huiying Liang, Sally L. Baxter, Alex McKeown, Ge Yang, Xiaokang Wu, Fangbing Yan, Justin Dong, Made K. Prasadha, Jacqueline Pei, Magdalene Y.L. Ting, Jie Zhu, Christina Li, Sierra Hewett, Jason Dong, Ian Ziyar, Alexander Shi, Runze Zhang, Lianghong Zheng, Rui Hou, William Shi, Xin Fu, Yaou Duan, Viet A.N. Huu, Cindy Wen, Edward D. Zhang, Charlotte L. Zhang, Oulan Li, Xiaobo Wang, Michael A. Singer, Xiaodong Sun, Jie Xu, Ali Tafreshi, M. Anthony Lewis, Huimin Xia, and Kang Zhang. Identifying medical diagnoses and treatable diseases by image-based deep learning. *Cell*, 172(5):1122–1131.e9, 2018. ISSN 0092-8674. doi: <https://doi.org/10.1016/j.cell.2018.02.010>. URL <https://www.sciencedirect.com/science/article/pii/S0092867418301545>.
- [41] Anna Majkowska, Sid Mittal, David F. Steiner, Joshua J. Reicher, Scott Mayer McKinney, Gavin E. Duggan, Krish Eswaran, Po-Hsuan Cameron Chen, Yun Liu, Sreenivasa Raju Kalidindi, Alexander Ding, Greg S. Corrado, Daniel Tse, and Shravya Shetty. Chest radiograph interpretation with deep learning models: Assessment with radiologist-adjudicated reference standards and population-adjusted evaluation. *Radiology*, 294(2):421–431, 2020. doi: 10.1148/radiol.2019191293. URL <https://doi.org/10.1148/radiol.2019191293>. PMID: 31793848.

## A Dataset splits

For each experiment on a single dataset, we split the dataset into 5 training and validation folds. Each training fold contains the same number of samples per each class. An exception is the NIH dataset since the number of samples across classes is highly imbalanced. In this case, we continue sampling to maximize the number of samples per each class as much as possible and use oversampling during the training process to compensate for class imbalance. Regarding each sample size, we use the absolute number of samples by rounding after multiplying with the corresponding percentage. Hence in section B, we report the sample size in absolute values rather than percentage.

## B Detailed results

### B.1 Linear Evaluation

Table B1: Linear evaluation on the PatchCam dataset with various initializations. The mean AUC is obtained across 5 folds.

No. of training samples	50	100	200	500 (S)	1000	2000	5000 (M)	10000	20000	Full
Mean AUC										
Random	0.5041	0.5001	0.5629	0.6594	0.6886	0.6955	0.6994	0.7110	0.7210	0.7990
Supervised ImageNet	0.7193	0.7287	0.7183	0.7517	0.7667	0.7709	0.7863	0.7894	0.7970	0.7975
SwAV	0.7543	0.7756	0.7510	0.7834	0.7686	0.8022	0.8043	0.8338	0.8310	0.8088
SimCLR	0.7502	0.7714	0.7675	0.7895	0.7981	0.7996	0.8053	0.8051	0.8110	0.8084
DINO	0.7721	0.8040	0.8010	0.8058	0.8204	0.8214	0.8359	0.8399	0.8446	0.8487

Table B2: Linear evaluation on the APTOS dataset with various initializations. The mean Kappa score is obtained across 5 folds.

No. of training samples	50 (S)	100	200	500 (M)	Full
Mean Kappa Score					
Random	0.0324	-0.0272	0.0083	0.0624	0.1550
Supervised ImageNet	0.4851	0.5758	0.6752	0.6822	0.7331
SwAV	0.6330	0.6559	0.7100	0.7274	0.7617
SimCLR	0.5305	0.5657	0.6369	0.6500	0.6989
DINO	0.6003	0.6889	0.7339	0.7372	0.7790

Table B3: Linear evaluation on the Pneumonia Chest X-Ray dataset with various initializations. The mean AUC is obtained across 5 folds.

No. of training samples	50 (S)	100	200	500 (M)	Full
Mean AUC					
Random	0.6899	0.7323	0.7720	0.8258	0.8907
Supervised ImageNet	0.8789	0.8788	0.8789	0.8954	0.9397
SwAV	0.8808	0.8731	0.8753	0.9215	0.9709
SimCLR	0.9168	0.9010	0.9176	0.9346	0.9665
DINO	0.9492	0.9466	0.9553	0.9718	0.9868

Table B4: Linear evaluation on the NIH Chest X-ray dataset with various initializations. The mean AUC is obtained across 5 folds.

No. of training samples	20 (S)	50	100	150	200 (M)	Full
Mean AUC						
Random	0.5212	0.5127	0.5031	0.5044	0.5317	0.5392
Supervised ImageNet	0.5383	0.5897	0.6388	0.6432	0.6688	0.7109
SwAV	0.5785	0.6469	0.6563	0.6673	0.6889	0.7225
SimCLR	0.5792	0.6273	0.6359	0.6686	0.6645	0.6983
DINO	0.6323	0.6831	0.7108	0.7385	0.7373	0.7438

## B.2 Finetuning

For all pretrained models except DINO, we update the gradients across all layers with a learning rate between 1e-3 and 1e-4 and a batch size of 64. For DINO, we use a smaller learning between 5e-4 and 1e-5 and a batch size of 16.

Table B5: Finetuning on the PatchCam dataset with various initializations. The mean AUC is obtained across 5 folds.

No. of training samples	50	100	200	500 (S)	1000	2000	5000 (M)	10000	20000	Full
Mean AUC										
Random	0.5181	0.5725	0.7016	0.7355	0.7642	0.7674	0.7660	0.7846	0.8114	0.8515
Supervised ImageNet	0.6359	0.6844	0.7656	0.7897	0.7870	0.7950	0.8274	0.8338	0.8491	0.8483
SwAV	0.6237	0.6164	0.6539	0.7895	0.8174	0.8221	0.8399	0.8338	0.8587	0.8619
SimCLR	0.6631	0.6116	0.6305	0.8021	0.8160	0.7944	0.8329	0.8402	0.8492	0.8553
DINO	0.7901	0.8115	0.8028	0.8366	0.8454	0.8438	0.8440	0.8379	0.8745	0.8517

Table B6: Finetuning on the APTOS dataset with various initializations. The mean Kappa score is obtained across 5 folds.

No. of training samples	50 (S)	100	200	500 (M)	Full
Mean Kappa Score					
Random	0.0177	0.0808	0.1914	0.3233	0.5927
Supervised ImageNet	0.4817	0.5289	0.6634	0.7369	0.8057
SwAV	0.4928	0.6015	0.7354	0.7594	0.8293
SimCLR	0.5916	0.6085	0.6860	0.7603	0.8264
DINO	0.6601	0.7144	0.7754	0.7945	0.8365

Table B7: Finetuning on the Pneumonia Chest X-ray dataset with various initializations. The mean AUC is obtained across 5 folds.

No. of training samples	50 (S)	100	200	500 (M)	Full
Mean AUC					
Random	0.6895	0.8210	0.9004	0.9183	0.9820
Supervised ImageNet	0.8649	0.9032	0.9157	0.9698	0.9910
SwAV	0.9289	0.9248	0.9593	0.9814	0.9927
SimCLR	0.9197	0.9199	0.9436	0.9781	0.9950
DINO	0.9256	0.9363	0.9687	0.9867	0.9948

Table B8: Finetuning on the NIH Chest X-ray with various initializations. The mean AUC is obtained across 5 folds.

No. of training samples	20 (S)	50	100	150	200 (M)	Full
Mean AUC						
Random	0.5015	0.5492	0.5961	0.6114	0.6404	0.6616
Supervised ImageNet	0.5251	0.6105	0.6467	0.6639	0.6816	0.7618
SwAV	0.5903	0.6172	0.6616	0.7037	0.6973	0.7737
SimCLR	0.5570	0.6227	0.6768	0.6795	0.7228	0.7358
DINO	0.5552	0.6348	0.6689	0.6551	0.6652	0.7404

## C Detailed results of DVME

Embedding from each self-supervised pretrained model is projected into a dimension of 512. Embeddings from SimCLR, SwAV, and DINO add up to a dimension of 1536. The self-attention module is implemented based on the timm library<sup>7</sup>. The output from the self-attention layer is further projected to a 512-dimensional layer followed by a ReLU layer, and the final linear layer. Table C1-C4 show the detailed result of linear evaluation using DVME with and without self-attention.

Table C1: Linear evaluation using DVME on the PatchCam dataset. The mean AUC is obtained across 5 folds.

No. of training samples	50	100	200	500 (S)	1000	2000	5000 (M)	10000	20000	Full
Mean AUC										
DVME w/o self-attention	0.7376	0.7906	0.8076	0.8196	0.8200	0.8242	0.8442	0.8417	0.8525	0.8478
DVME	0.7456	0.7864	0.8026	0.8227	0.8316	0.8243	0.8399	0.8404	0.8444	0.8467

Table C2: Linear evaluation using DVME on the APTOS dataset. The mean Kappa score is obtained across 5 folds.

No. of training samples	50 (S)	100	200	500 (M)	Full
Mean Kappa Score					
DVME w/o self-attention	0.6354	0.7018	0.7351	0.7681	0.7759
DVME	0.6913	0.6992	0.7787	0.7925	0.8242

Table C3: Linear evaluation using DVME on the Pneumonia Chest X-ray dataset. The mean AUC is obtained across 5 folds.

No. of training samples	50 (S)	100	200	500 (M)	Full
Mean AUC					
DVME w/o self-attention	0.9543	0.9528	0.9532	0.9725	0.9865
DVME	0.9539	0.9469	0.9569	0.9696	0.9842

Table C4: Linear evaluation using DVME on the NIH Chest X-ray dataset. The mean AUC is obtained across 5 folds.

No. of training samples	20 (S)	50	100	150	200 (M)	Full
Mean AUC						
DVME w/o self-attention	0.6525	0.7051	0.7260	0.7209	0.7232	0.7575
DVME	0.6566	0.6871	0.7091	0.7437	0.7601	0.7538

<sup>7</sup>[https://github.com/rwightman/pytorch-image-models/blob/master/timm/models/vision\\_transformer.py](https://github.com/rwightman/pytorch-image-models/blob/master/timm/models/vision_transformer.py)



## D Embedding visualization

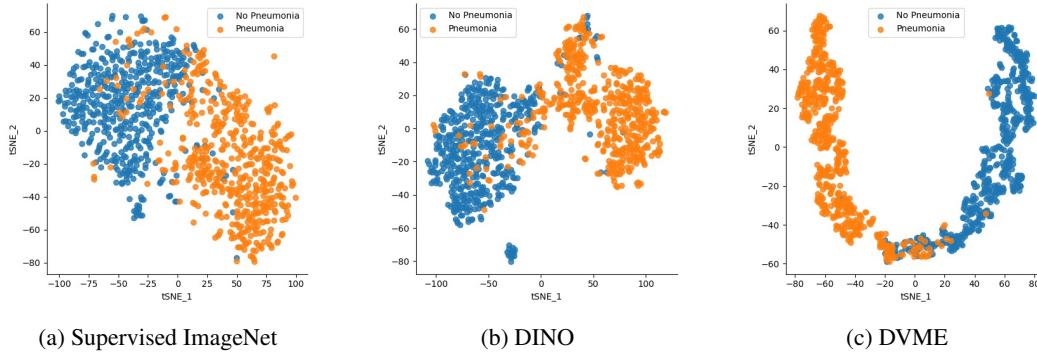


Figure D1: t-SNE visualization of the pre-trained embeddings from supervised ImageNet, DINO, and our proposed method DVME on **Pneumony Chest X-ray dataset**

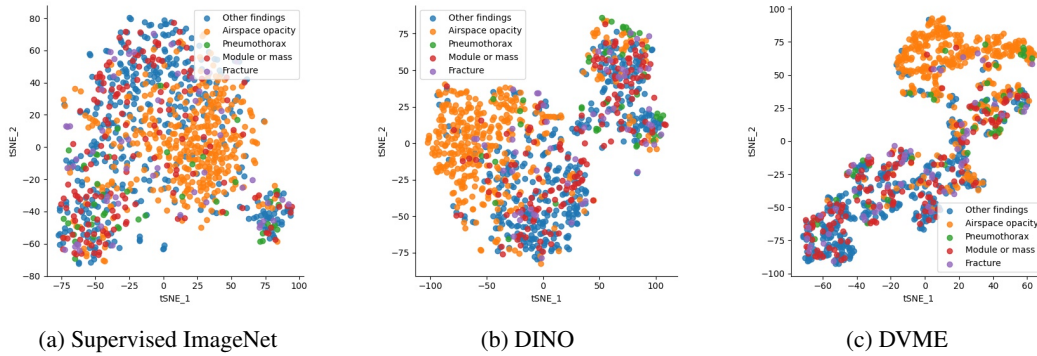


Figure D2: t-SNE visualization of the pre-trained embeddings from supervised ImageNet, DINO, and our proposed method DVME on **NIH Chest X-ray dataset**

## E Dynamic Visual Meta-embeddings

```

1 import torch.nn as nn
2 import torch
3
4 class DVME(nn.Module):
5
6     def __init__(self, proj_dim, num_cls, attn):
7         # proj_dim: dimension of projection (default is 512)
8         # num_cls: number of classes
9         # attn: self-attention module
10
11         super(DVME, self).__init__()
12         self.simclr_head = nn.Linear(2048, proj_dim)
13         self.swav_head = nn.Linear(2048, proj_dim)
14         self.dino_head = nn.Linear(1536, proj_dim)
15         self.attn = attn
16         self.normlayer = nn.LayerNorm(proj_dim*3)
17         self.proj_head = nn.Linear(proj_dim*3, proj_dim)
18         self.classifier = nn.Linear(proj_dim, num_cls)
19         self.dropout = nn.Dropout(0.2)
20
21
22

```

```

23 def forward(self, x):
24     # x: dictionary containing extracted embeddings from
25     # pretrained models SimCLR, SwAV, DINO
26
27     simclr_out = self.simclr_head(x['simclr'])
28     swav_out = self.swav_head(x['swav'])
29     dino_out = self.dino_head(x['dino'])
30     meta_x = torch.cat([simclr_out, swav_out, dino_out], dim=1)
31     # reshape the meta-emb into (batch, tokens, dim)
32     meta_x = meta_x.view(meta_x.size(0), -1, 1)
33     out = self.attn(meta_x)
34     out = self.normlayer(out.view(out.size(0), -1))
35     out = self.proj_head(out).relu()
36     out = self.dropout(out)
37     out = self.classifier(out)
38     return out

```

Listing 1: PyTorch code snippet of DVME