

# Benchmarking Self-Supervised Learning on Diverse Pathology Datasets

Mingu Kang\* Heon Song\* Seonwook Park Donggeun Yoo Sérgio Pereira  
Lunit Inc.

{jeffkang, heon.song, spark, dgyoo, sergio}@lunit.io

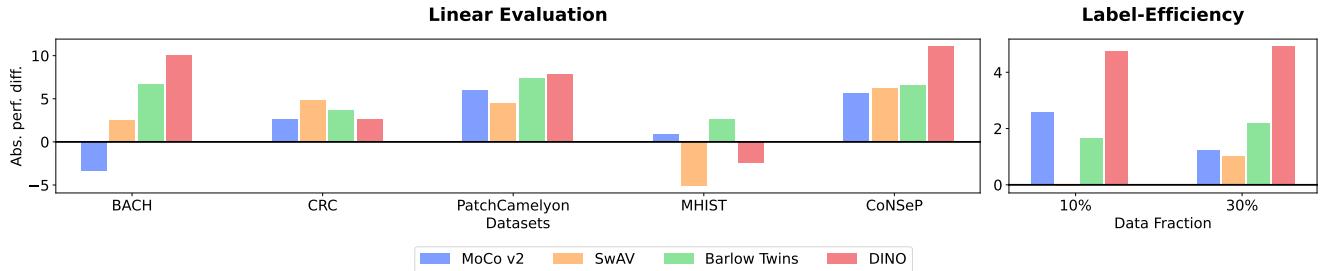


Figure 1. **Self-supervised pre-training on pathology data improves performance on pathology downstream tasks** compared to ImageNet-supervised baselines. The y-axes show absolute differences in downstream task performance (Top-1 Acc. or mPQ Score). Linear evaluation (**left**) is performed on 4 classification tasks (BACH, CRC, PatchCamelyon, and MHIST) and 1 nuclei instance segmentation task (CoNSeP). Label-efficiency (**right**) is assessed by fine-tuning using small fractions of labeled data from the CoNSeP dataset.

## Abstract

*Computational pathology can lead to saving human lives, but models are annotation hungry and pathology images are notoriously expensive to annotate. Self-supervised learning (SSL) has shown to be an effective method for utilizing unlabeled data, and its application to pathology could greatly benefit its downstream tasks. Yet, there are no principled studies that compare SSL methods and discuss how to adapt them for pathology. To address this need, we execute the largest-scale study of SSL pre-training on pathology image data, to date. Our study is conducted using 4 representative SSL methods on diverse downstream tasks. We establish that large-scale domain-aligned pre-training in pathology consistently out-performs ImageNet pre-training in standard SSL settings such as linear and fine-tuning evaluations, as well as in low-label regimes. Moreover, we propose a set of domain-specific techniques that we experimentally show leads to a performance boost. Lastly, for the first time, we apply SSL to the challenging task of nuclei instance segmentation and show large and consistent performance improvements. We release the pre-trained model weights<sup>1</sup>.*

## 1. Introduction

The computational analysis of microscopic images of human tissue – also known as computational pathology – has emerged as an important topic of research, as its clinical implementations can result in the saving of human lives

by improving cancer diagnosis [57] and treatment [50]. Deep Learning and Computer Vision methods in pathology allow for objectivity [18], large-scale analysis [23], and triaging [7] but often require large amounts of annotated data [60]. However, the annotation of pathology images requires specialists with many years of clinical residency [42], resulting in scarce labeled public datasets and the need for methods to train effectively on them.

When annotated data is scarce for a given Computer Vision task, one common and practical solution is to fine-tune a model that was pre-trained in a supervised manner using the ImageNet dataset [22, 39]. This paradigm of transfer learning [39] was recently challenged by self-supervised learning (SSL), which trains on large amounts of unlabeled data only, yet out-performs supervised pre-training on ImageNet [10, 13, 30]. In the field of pathology, large unlabeled datasets are abundant [6, 42, 43, 65] in contrast to the lack of annotated datasets [60]. If we were to apply SSL effectively to this huge amount of unlabeled data, downstream pathology tasks could benefit greatly even if they contain limited amount of annotated training data. Naturally, we ask the question: *How well does self-supervised learning help in improving the performance of pathology tasks?*

ImageNet pre-trained weights are commonly used in medical imaging and are known to be helpful in attaining high task performance [35, 37, 51, 67]. Due to the difference between natural images and medical images, large-scale domain-aligned pre-training has the potential to push performance beyond ImageNet pre-training [47]. Accordingly, recent works show that SSL pre-training on pathol-

\*The first two authors contributed equally.

<sup>1</sup>[https://lunit-io.github.io/research/publications/pathology\\_ssl](https://lunit-io.github.io/research/publications/pathology_ssl)

ogy data can improve performance on downstream pathology tasks [5, 19, 26, 63]. Our study aims to expand on these previous works by evaluating multiple SSL methods on diverse downstream pathology tasks. In addition, we propose techniques to adapt SSL methods that were designed for natural image data, to better learn from pathology data.

To understand how to adapt existing SSL methods to work on pathology image data, we must identify several key differences between natural and pathology imagery. Unlike natural images, pathology images can be rotated arbitrarily (impossible to determine a canonical orientation) and exhibit fewer variations in color. Also, pathology images can be interpreted differently depending on the field-of-view (FoV) due to the multiple hierarchies and contextual differences involved in each task. We propose to overcome these differences when adapting SSL methods for pathology data, via changes to the training data augmentation scheme in particular, during pre-training.

In this paper, we carry out an in-depth analysis of 4 recent and representative SSL methods; MoCo v2 [15], SwAV [9], Barlow Twins [70], and DINO [10], when applied to large-scale pathology data. For this purpose, we source 19 million image patches from Whole Slide Images (WSI) in The Cancer Genome Atlas (TCGA) dataset [65] and apply our domain-specific techniques in training the SSL methods on this data. The evaluations are conducted for 2 different downstream tasks over 5 datasets: (1) pathological image classification using BACH [2], CRC [36], MHIST [64], and PatchCamelyon [62] datasets, and (2) nuclei instance segmentation and classification using the CoNSEP dataset [29].

Our large-scale study yields several useful contributions: (a) we conduct the largest-scale study of SSL pre-training on pathology image data to date, and show its benefit over using ImageNet pre-trained weights on diverse downstream tasks (see Fig. 1), (b) we propose a set of carefully designed data curation and data augmentation techniques that can further improve downstream performance, (c) we demonstrate that SSL is label-efficient, and is therefore a practical solution in pathology where gathering annotation is particularly expensive, and (d) for the first time, we apply SSL to the dense prediction task of nuclei instance segmentation and show its value under diverse evaluation settings. We release our pre-trained model weights at [https://lunit-io.github.io/research/publications/pathology\\_ssl](https://lunit-io.github.io/research/publications/pathology_ssl) to further contribute to the research community.

## 2. Related Work

### 2.1. Self-supervised Representation Learning

SSL methods learn representations through pre-text tasks designed to exploit supervisory signals obtained from the unlabeled data itself. We describe the 4 major paradigms of SSL as commonly discussed in literature.

**Contrastive Learning.** Contrastive methods [31, 48, 49] such as SimCLR [13] and MoCo v2 [15] learn to discriminate each training data instance from all the others. The objective is to learn similar representations of positive pairs (perturbations by data augmentation) and discriminative representations in relation to negative pairs (other instances). A limitation is the need for diverse negative pairs, which is mitigated through large batch sizes [13] or memory banks [15]. In this work, we explore MoCo v2 [15].

**Non-contrastive Learning.** Methods such as BYOL [30], SimSiam [16], and Barlow Twins [70], share similarities with *contrastive learning* methods in that they learn representations of images under different augmented views. The fundamental difference is that these approaches do not rely on negative pairs, which allows them to work with small batch sizes. In this work, we explore Barlow Twins [70].

**Clustering.** This paradigm uses the concept of clustering and is shown in DeepCluster [8] and SwAV [9]. Clustering-based SSL discriminates between clusters of image representations instead of explicit pairs of images. In this work, we explore SwAV [9].

**SSL with VisionTransformer.** The effectiveness of Vision Transformers (ViT) [24] has been demonstrated on various computer vision tasks. Thus, the paradigm shift from CNN to ViT has recently emerged in the field of self-supervised learning. Consequently, recent studies [10, 17, 41] try to investigate techniques that facilitate SSL with ViT-based architectures. In this work, we explore DINO [10].

### 2.2. SSL in Medical Imaging

Recently, [47] investigates transfer learning in medical imaging and observes that using domain-aligned datasets for pre-training improves the transferability of models. Moreover, domain-specific SSL methods can further improve the performance of models fine-tuned on downstream medical image-related tasks [3, 5, 19, 26, 56, 63, 68]. In pathology, [63] employs BYOL and evaluates pre-trained weights learned from pathology data on image classification tasks. [26] adopts SimSiam, showing that SSL improves pathology image retrieval. Also recently, [19] uses SimCLR and observes that SSL consistently improves on downstream pathology tasks compared to ImageNet pre-training.

Unlike previous works that focus on just a single SSL approach [12, 19, 40], or either CNNs or ViTs only [68], we explore one representative method from each of the aforementioned SSL paradigms including ViT-based SSL. In this way, we establish a common and fair benchmark for comparing these methods in pathology. Furthermore, we evaluate the domain-specific pre-trained weights on various downstream tasks, including the challenging task of nuclei instance segmentation. Finally, we devise techniques for data augmentation that are specifically tailored to tackle

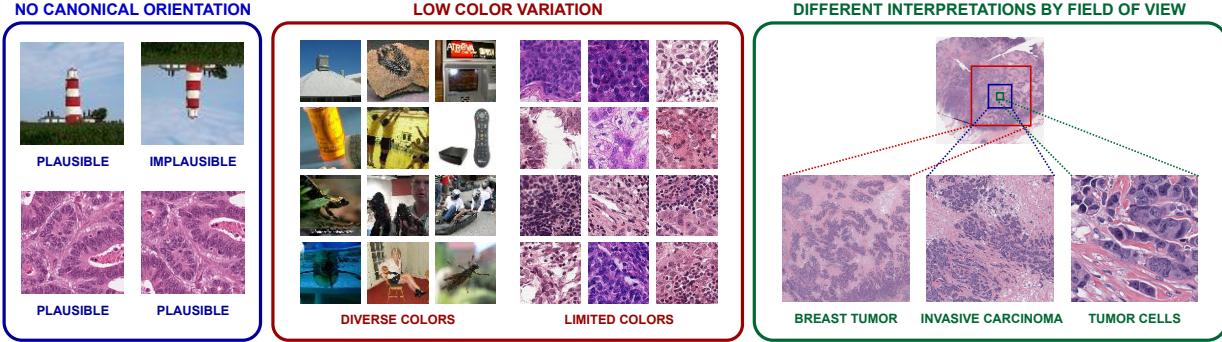


Figure 2. **Pathology images vs natural images.** Pathology images are different from natural images in 3 major ways. They have no canonical orientation (no “right way up”), have low color variation, and can be interpreted differently depending on field-of-view. Consequently, self-supervised learning methods need to be implemented differently when working in the domain of pathology.

pathology-specific challenges, thus leading to better representations and performance in the downstream tasks.

### 3. Self-supervised Pre-training for Pathology

The performance of SSL methods can vary greatly depending on the composition of training data and the selected set of data augmentation methods. SSL methods in the literature are commonly designed and evaluated in settings involving natural images and may benefit from further adaptation when applied to different domains, such as pathology. In this section, we discuss the differences between natural images and pathology images. We also propose a set of techniques that can be easily adopted to improve the performance of models pre-trained on pathology image data.

#### 3.1. Differences to Natural Images

Images contained in popular Computer Vision datasets (e.g. ImageNet [22]) are often denoted as “natural images”. Pathology images have several unique characteristics that make them distinct from natural images. We discuss these differences in this section and summarize them in Fig. 2.

**No canonical orientation.** Objects or scenes contained in natural images are oriented based on plausibility, i.e. how a human expects the objects to be oriented. Methods in Computer Vision can take advantage of such assumptions or patterns (e.g. Manhattan World assumption [20]) and thus SSL methods do not randomly augment the orientation of images at training time. However, pathology images can be oriented in any way and still remain plausible. Furthermore, objects (e.g. cells) are many and dispersed at arbitrary locations, making it impossible to define a “canonical orientation”, i.e. the correct standard orientation.

**Low color variation.** While natural images contain a large range of colors due to the diversity of represented objects, pathology images tend to display similar color distributions (e.g. purple and pink staining). Though the staining can vary across institutions and the same biological structures

have different appearances depending on the cancer type, pathology images are more consistent than natural images.

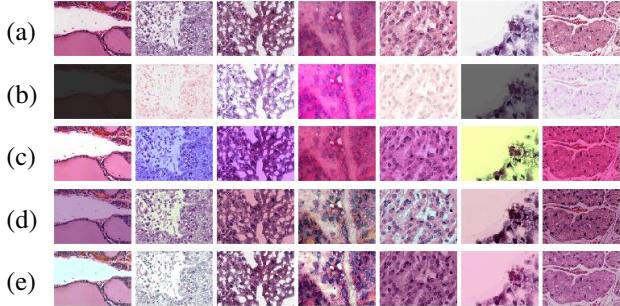
**Different FoVs.** To correctly analyze pathology images, different Field of Views (FoVs) must be considered. A larger FoV allows pathologists and algorithms to better understand the larger context of the tissue regions and cell classes to make high-level predictions such as the grading of prostate cancer [6]. In other tasks that require the classification of individual cells or communities of cells, a smaller FoV is required to increase the resolution on the objects of interest [29]. Thus, a pre-trained model for pathology should ideally be able to handle tasks from diverse FoVs.

#### 3.2. Techniques to Adapt SSL for Pathology

In this section, we introduce our techniques for adapting SSL methods for pathology imagery.

**Random vertical flips.** Unlike natural images, pathology images are no less plausible or realistic when they are vertically flipped. We therefore propose to randomly apply vertical flips during SSL pre-training.

**Stain augmentations.** The typical color distortion employed by SSL methods applies a strong amount of jitter to brightness, contrast, and saturation, resulting in pathology images that look highly unrealistic. [58] proposes to apply this jitter in pathology-specific color spaces such as the HED-space [54]. [55] points out that naive jittering can produce unrealistic images and proposes RandStainNA. RandStainNA fits unimodal Gaussian distributions to the channel-wise statistics of 3 color spaces (HSV, Lab, HED) using images from the training dataset. At training time, a color space is randomly chosen, then the target channel-wise mean and standard deviations for that color space are sampled from the fitted Gaussian distributions. Reinhard’s method [52] is used to re-normalize the input image to match the target statistics. RandStainNA is shown to improve supervised learning performance for pathology, and therefore we adopt it for our SSL pre-training.



**Figure 3. Color augmentations on pathology images.** (a) input images, (b) color jitter as used in typical SSL methods [13], (c) HED-light [58], (d) RandStainNA [55], and (e) RandStainNA<sub>GMM</sub>. RandStainNA and RandStainNA<sub>GMM</sub> can produce more realistic and plausible pathology images.

Data source	No. of WSIs	No. of patches			Task
		20x	40x	Total	
TCGA	20,994	9,497,768	9,502,301	19,000,069	
TULIP	15,672	7,084,130	6,494,358	13,578,488	
Total	36,666	16,581,898	15,996,659	32,578,557	

**Table 1. Unlabeled data for pre-training.** Amount of data used for pre-training in terms of the number of WSIs and the actual number of extracted patches (per FoV).

Furthermore, we attempt to improve the realism of RandStainNA by fitting a Gaussian Mixture Model (GMM) with 10 components, to the channel-wise statistics of each color space. The GMM can fit the covariances between variables and respect the multi-modality of the channel-wise mean and standard deviation values. We denote this as RandStainNA<sub>GMM</sub> and show the visual differences against alternative methods in Fig. 3.

Lastly, we observe that previous works in SSL [13, 16, 30, 70] highlight the importance of color distortion. We therefore propose to apply color distortion with a weaker jittering strength as done in [19]. Our main experiments adopt RandStainNA<sub>GMM</sub> along with random grayscale and a weaker color jittering augmentation.

**Using multiple FoVs.** As aforementioned, some pathology tasks require high FoV while others benefit from low FoV. We identify that pathology tasks (e.g., image classification and instance segmentation) are commonly defined at approximately 20× [2, 36] or 40× [29, 62, 64] objective magnification. Therefore, we build our large-scale unlabeled dataset using image patches from both magnifications.

## 4. Experiment Setup

### 4.1. Pre-training Dataset

Tab. 1 presents the scale of unlabeled data used for pre-training. We first collect 20,994 WSIs from The Cancer

Dataset	# Classes	# Patches	Patch size	FoV	Task
BACH	4	400	2048×1536	20×	Cls
CRC	9	107,180	224×224	20×	Cls
PCam	2	327,680	96×96	40×	Cls
MHIST	2	3,152	224×224	40×	Cls
CoNSeP	7	41	1000×1000	40×	Seg

**Table 2. Datasets for downstream tasks.** Note that, Cls indicates “image classification” and Seg is “nuclei instance segmentation”.

Genome Atlas (TCGA) and 15,672 WSIs from TULIP. Both datasets consist of Hematoxylin & Eosin (H&E) stained WSIs of various cancers. TCGA is publicly available and widely used for training deep learning models [21, 23, 50]. TULIP is an internally collected dataset. To increase diversity and keep our experimental setting practical, we extract at most 1,000 patches of resolution 512 × 512 pixels from each slide, resulting in a total of 32.6M patches (19M from TCGA and 13.6M from TULIP). The pre-training data covers two different FoVs; 20× (0.5μm/px) and 40× (0.25μm/px) objective magnification. All experiments, unless specified otherwise, present the results of pre-training on the TCGA dataset only.

### 4.2. Downstream Datasets

We validate the pre-trained models under classification and segmentation tasks using various downstream datasets described in Tab. 2. For image classification, the following four datasets are used: BACH (four-class breast cancer type classification) [2], CRC (nine-class human colorectal cancer type classification) [36], MHIST (binary class colorectal polyp type classification) [64], and PCam (binary class breast cancer type classification) [62]. The patches of the datasets are labeled according to the predominant cancer type or the presence of cancers. For nuclei instance segmentation, we use CoNSeP [29] which contains segmentation masks for each cell nucleus along with nuclei types. Further details of the downstream datasets are shown in the supplementary materials.

### 4.3. Pre-training Details

We learn representations using 4 different SSL methods. Unless otherwise mentioned, we use the ResNet-50 (1×) [33] architecture for MoCo v2 [15], Barlow Twins [70], and SwAV [9]. For DINO [10], we use ViT-Small [24] with different patch sizes, 16×16 and 8×8 (denoted DINO<sub>p=16</sub> and DINO<sub>p=8</sub>, respectively), as it has a comparable number of parameters to ResNet-50 (1×). We follow the proposed recipe of each SSL method and launch the pre-training, distributed across 64 NVIDIA V100 GPUs. The linear scaling rule [27] is applied to adjust the learning rate:  $lr = lr_{method} * batchsize / 256$ . We adopt the concept of the “ImageNet epoch” from [59] for ease of analysis and train models for 200 ImageNet epochs, across all experiments. We

Arch.	Method	BACH		CRC		PCam		MHIST	
		Linear	Fine-tune	Linear	Fine-tune	Linear	Fine-tune	Linear	Fine-tune
ResNet-50	Random	51.67	61.67	68.91	89.99	76.52	75.71	63.15	75.54
	Supervised	80.83	86.67	90.93	92.09	80.79	80.63	76.25	78.92
	MoCo v2	77.50	<b>90.83</b>	93.52	<b>96.21</b>	86.78	<b>87.62</b>	77.07	<b>85.88</b>
	SwAV	83.33	82.50	<b>95.78</b>	93.31	85.28	87.60	71.14	77.99
	BT	<b>87.50</b>	85.00	<b>94.60</b>	93.23	<b>88.15</b>	86.92	<b>78.81</b>	81.27
ViT-S	$Random_{p=16}$	45.00	57.50	69.90	86.10	74.43	75.42	63.46	62.13
	$Supervised_{p=16}$	75.83	85.83	91.56	<b>95.81</b>	80.96	88.30	<b>78.51</b>	<b>81.68</b>
	$DINO_{p=16}$	<b>85.83</b>	<b>87.50</b>	94.19	<b>95.81</b>	88.78	<b>90.40</b>	76.15	79.43
	$DINO_{p=8}$	83.33	<b>93.33</b>	<b>95.29</b>	<b>97.13</b>	<b>90.12</b>	<b>90.76</b>	77.89	78.40

Table 3. **Downstream evaluation of the image classification tasks.** We report Top-1 accuracy for both linear and fine-tuning experiment protocols for models trained using the TCGA data source. Note that,  $p$  represents the patch size used in ViT.

define 10 ImageNet epochs for the warmup, and the cosine scheduler is followed. The details of configurations can be found in supplementary materials.

#### 4.4. Downstream Training Details

For the downstream tasks, we follow the standard practice as introduced in various SSL papers [13, 30, 70]. For image classification tasks, the datasets are split into training, validation, and test sets. We perform the hyper-parameter search based on the validation set, reporting the performance on the test set. For the segmentation task, the Hover-Net [29] architecture – a standard architecture in the nuclei instance segmentation task – is adopted with the pre-trained backbone. We follow the same data pre-processing and training schemes as in [29] to enable reproducibility and fair comparison of results. Further details of downstream tasks can be found in the supplementary materials.

#### 4.5. Evaluation Metrics

For image classification, we report top-1 accuracy, while using panoptic quality (PQ) [38] for nuclei instance segmentation. PQ is a standard metric for assessing the performance of nuclear instance segmentation [29] that accounts for both detection and segmentation quality with respect to each instance. The PQ metric is defined as,

$$PQ = \frac{\sum_{(p,g) \in TP} \text{IoU}(p, g)}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|}, \quad (1)$$

where  $p$  denotes a predicted mask for each nuclei class and  $g$  denotes a corresponding ground truth. The numerator  $\sum_{(p,g) \in TP} \text{IoU}(p, g)$  represents the summation of correctly matched Intersection Over Union (IoU) over all pairs between predictions and ground truth. We count pairs between predictions and ground truths with IoU of more than 0.5 as True Positives (TP). False Positives (FP) and False Negatives (FN) represent wrongly predicted predictions and ground truth, respectively. Note that we use multi-class PQ

(mPQ) to measure the performance of instance segmentation and classification simultaneously.

### 5. Experimental Results

In this section, we carry out various experiments on the downstream tasks of image classification and nuclei instance segmentation. Through these experiments, we compare the utility of various SSL pre-training methods in the light of downstream pathology task performance. First, we present our quantitative results using well-established evaluation protocols. We then demonstrate the benefit of SSL pre-training with a limited number of labeled data and under different fine-tuning schedules. Finally, we perform an ablation study to quantitatively verify the efficacy of our techniques for adapting SSL methods to pathology data.

In evaluating downstream task performance, we stick to well-established evaluation protocols in SSL. The first is *linear evaluation* (denoted *Linear*), where we freeze the backbone and train the remaining parts of the model (e.g., linear classifier or decoders). The second is *full fine-tuning* (denoted *Fine-tune*), where all layers including the backbone are fine-tuned. The former protocol assesses the quality of learned representations, whereas the latter evaluates the transferability of learned weights. In our experiments, we compare against the ImageNet-supervised (denoted *Supervised*) pre-training baseline of the corresponding backbone type as well as a random initialization (denoted *Random*) baseline.

#### 5.1. Image Classification

We present our linear evaluation and fine-tuning results for 4 image classification benchmarks in Tab. 3.

**Linear evaluation.** In linear evaluation results, we find that self-supervised TCGA pre-training typically outperforms supervised ImageNet pre-training. Of the ResNet-50 based SSL methods, Barlow Twins performs consistently well, out-performing other methods on the BACH,

Arch.	Method	CoNSeP	
		Linear	Fine-tune
ResNet-50	<i>Random</i>	22.29	46.72
	<i>Supervised</i>	34.25	49.60
	MoCo v2	39.85	<b>51.75</b>
	SwAV	40.45	51.16
	BT	<b>40.79</b>	<u>51.61</u>
ViT-S	<i>Random</i> <sub>p=16</sub>	20.55	27.19
	<i>Supervised</i> <sub>p=16</sub>	21.43	36.70
	DINO <sub>p=16</sub>	32.54	38.43
	DINO <sub>p=8</sub>	<b>42.71</b>	<b>46.70</b>

Table 4. **Downstream evaluation for the nuclei instance segmentation task.** We report the mPQ score for both linear and fine-tuning experiment protocols for models trained using the TCGA data source.

PCam, and MHIST datasets. Of the ViT-based SSL methods, DINO<sub>p=16</sub> achieves competitive results, and DINO<sub>p=8</sub> performs even better on the CRC and PCam datasets. The improved performance of DINO<sub>p=8</sub> is in line with observations from [10] which shows that performance improves at the cost of computation with smaller patch size. One exception is on the MHIST dataset, where the supervised baseline shows good performance. Based on the linear evaluation results, we can certainly claim that domain-aligned pre-training improves representation quality.

**Fine-tuning.** Under fine-tuning, we find that the trends are similar but with notable changes. Firstly, as shown in other SSL works, the performance gap between ImageNet-supervised and TCGA-SSL reduces. Furthermore, MoCo v2 shows consistently high performance among CNN methods, showing that it may be the architecture of choice for transfer learning settings using CNNs. Regarding ViTs, we find that the trends are almost identical to linear evaluation except that the fine-tuned performances are often better than CNN counterparts trained using SSL on TCGA data. For classification tasks, then, SSL using ViT on large-scale pathology data is beneficial.

## 5.2. Nuclei Instance Segmentation

To the best of our knowledge, we show for the first time, the effect of self-supervised domain-aligned pre-training on a downstream dense prediction task. We run experiments on the CoNSeP dataset for the task of nuclei instance segmentation and report the mPQ score in Tab. 4.

**CNN experiments.** The performance of SSL using ResNet-50 shows similar trends to the case of image classification, where Barlow Twins performs well on the linear evaluation protocol and MoCo v2 performs well in fine-tuning. More consistently than in the case of classification, SSL pre-trained models out-perform supervised ImageNet pre-training by a large margin, especially considering the difficulty of increasing the mPQ score.

Method	Pre. Data	Top-1 Acc. (%)				mPQ CoNSeP
		BACH	CRC	PCam	MHIST	
<i>Random</i>	-	51.67	68.91	76.52	63.15	22.29
<i>Supervised</i>	IN	80.83	90.93	80.79	76.25	33.49
MoCo v2	IN	71.67	92.86	82.37	79.73	39.13
MoCo v2	TCGA	77.50	93.52	<b>86.78</b>	77.07	39.85
MoCo v2	TC+TU	<b>85.00</b>	<b>93.94</b>	86.53	<b>82.29</b>	<b>41.40</b>

Table 5. **Varying pre-training datasets under the linear evaluation protocol.** We consider ImageNet (IN), TCGA, and TCGA and TULIP combined (TC+TU) as pre-training datasets. Training with TC+TU results in consistent performance improvements.

**ViT experiments.** To the best of our knowledge, we are the first to integrate ViT backbones into the HoVer-Net architecture for nuclei instance segmentation. We find that DINO trained on TCGA data out-performs ImageNet-trained weights for DINO<sub>p=16</sub> by a large margin, showing again the importance of domain-aligned SSL. While DINO<sub>p=16</sub> does not work well in neither linear nor fine-tuning evaluations, DINO<sub>p=8</sub> out-performs even CNN-based methods in linear evaluation and performs reasonably well with fine-tuning. Future work may be able to further unlock the power of transformers as a backbone for nuclei instance segmentation.

## 5.3. Pre-training on Different Datasets

The experiment aims to investigate the impact on the downstream task in accordance with the pre-training data. We select MoCo v2 as it has previously shown robust performance in relation to various domain-specific data [61]. We show our linear evaluation results in Tab. 5 where we compare against MoCo v2 pre-training on ImageNet<sup>2</sup> as well as on the combined TCGA and TULIP data. We note that TCGA-pre-training out-performs supervised/SSL pre-training with ImageNet on BACH, CRC, PCAM, and CoNSeP. When adding TULIP data into the mix, we can use a total of 36K slides and 32.6M patches for pre-training, and we see that this results in the overall best performance. Through these experiments, we conclude that using a domain-aligned dataset such as TCGA is useful, and increasing the scale and diversity of data can further boost performance.

## 5.4. Fine-tuning with Limited Labeled Data

In the pathology domain, acquiring high-quality annotations require expert-derived labor and exhaustive refinement to maximize consensus across annotators, hindering the establishment of large-scale annotated datasets. Prior findings in computer vision and medical imaging show that SSL pre-trained models are label-efficient under fine-tuning evaluations [3, 14]. To evaluate the label-efficiency of pre-trained models in the pathology domain, we perform similar eval-

<sup>2</sup>We use a publicly available ImageNet pre-trained model for MoCo v2.

Method	CRC (Top-1 Acc.)			CoNSeP (mPQ)		
	1%	10%	100%	10%	30%	100%
ResNet-50						
<i>Supervised</i>	90.28	93.87	92.09	40.01	41.92	49.60
MoCo v2	<b>91.73</b>	<b>95.10</b>	<b>96.21</b>	<b>42.59</b>	<b>43.15</b>	<b>51.75</b>
SwAV	89.26	92.84	93.31	39.97	42.94	51.16
BT	<u>91.23</u>	92.84	<u>93.23</u>	<u>41.66</u>	<b>44.10</b>	<u>51.61</u>
ViT-S						
<i>Supervised</i> <sub>p=16</sub>	93.15	94.76	95.81	18.49	20.92	36.70
DINO <sub>p=16</sub>	94.03	94.92	95.81	23.22	25.85	38.43
DINO <sub>p=8</sub>	<b>95.03</b>	<b>96.27</b>	<b>97.13</b>	<b>35.53</b>	<b>37.82</b>	<b>46.70</b>

Table 6. **Label-efficiency.** Full fine-tuning results when using a limited number of training samples for the CRC and CoNSeP downstream benchmarks.

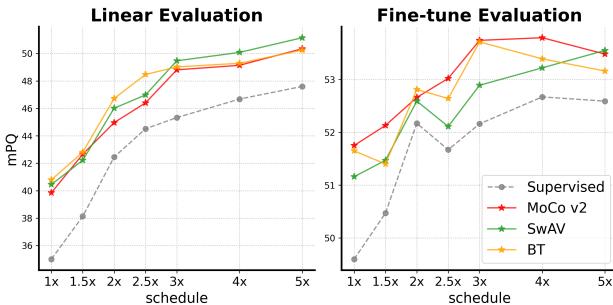


Figure 4. **Different learning schedules for the nuclei instance segmentation task using CoNSeP.** We scale up the learning schedule from 1x (20 epochs) to 5x (100 epochs).

uations and fine-tune our models while varying the labeled data fraction. Following prior works [13, 70], subsets of size 1% and 10% are sampled for the image classification dataset. We pick CRC dataset since it has sufficient amounts of data as well as a number of classes to conduct the experiment. On the other hand, the CoNSeP dataset has insufficient data to conduct the same experiment, and therefore, we use subsets of size 10% and 30% for nuclei instance segmentation. Further details can be found in our supplementary materials.

Tab. 6 presents fine-tuning evaluation results with varying amounts of downstream labeled data. Compared to the *Supervised* baselines, SSL pre-trained weights out-perform the ImageNet pre-trained weights. In particular, MoCo v2 and DINO<sub>p=8</sub> show the best performances for ResNet-50 and ViT-S backbones respectively, maintaining the performance gap to *Supervised* baselines even with increasing amounts of labeled data.

## 5.5. Longer Learning Schedules

When evaluating SSL methods on downstream dense prediction tasks, it is desirable to show results with longer fine-tuning schedules [31, 53, 66] but this is rarely shown in papers due to the diminishing performance gap between methods when fine-tuning for longer [32]. To better demon-

	S	V	G	ColorJitter		mPQ
				weak	strong	
Baseline	✓		✓		✓	50.71
	✓	✓	✓		✓	51.03
	✓	✓			✓	51.07
HED-light		✓				50.48
RandStainNA		✓				50.86
RandStainNA <sub>GMM</sub>	✓					50.71
		✓	✓			51.07
RandStainNA	✓	✓		✓		51.13
	✓	✓			✓	50.27
		✓	✓			50.99
RandStainNA <sub>GMM</sub>	✓	✓		✓		51.61
	✓	✓			✓	50.51

where S: Solarization, V: Vertical Flip, G: Grayscale

Table 7. **Augmentation ablation study on CoNSeP.** mPQ scores are computed after fine-tuning with a 1x schedule. *Baseline* refers to the original Barlow Twins setting.

strate the real-world implications of SSL pre-training on pathology data for nuclei instance segmentation, we scale the training schedule from 1x to 5x and show mPQ scores computed on the CoNSeP dataset in Fig. 4. In both linear and fine-tuning evaluations, we clearly observe the benefit of TCGA pre-trained weights compared to ImageNet-supervised weights and note that the performance gap is maintained even at longer training schedules.

## 5.6. Ablation Study

As described in Sec 3.2, we propose to use data augmentation techniques tailored for pathology data. We empirically assess whether these domain-adapted techniques are beneficial to SSL in pathology through an ablation study. We select Barlow Twins [70] to conduct our experiments and pre-train our models for 200 ImageNet epochs on the TCGA dataset. The fine-tuning evaluation protocol is adopted to better indicate real-world implications.

**Augmentation Ablation.** We select nuclei instance segmentation as a target task, as it is one of the most practical and challenging tasks in pathology. Starting from the default data augmentation of Barlow Twins [70] denoted as *Baseline* in Tab. 7, we add a random vertical flip augmentation in order to take advantage of the nature wherein pathology images have no canonical orientation. Based on prior work that claims that solarization can produce unrealistic and harmful images for pathology model training [25], we exclude the random solarization augmentation. With these two changes, we observe a gain of 0.36 in mPQ score.

We expect that stain augmentations serve to generate more domain-relevant augmented views, particularly, in terms of color variations. However, stain augmentation clashes with color distortion, as both alter the color statis-

FoV	Top-1 Acc. (%)				mPQ
	BACH	CRC	PCam	MHIST	
20×	79.17	<b>95.04</b>	85.24	79.32	50.66
40×	79.17	91.88	80.82	74.21	48.83
20×, 40×	<b>85.00</b>	93.23	<b>86.92</b>	<b>81.27</b>	<b>51.61</b>

Table 8. **Magnification ablation study.** Fine-tuning performance when using different FoVs during pre-training.

tics of images. Thus, we begin by disabling color distortion and then compare key stain augmentation methods first. We find that RandStainNA [55] and RandStainNA<sub>GMM</sub> outperform HED-light [58], confirming insights from supervised image classification [55].

Next, we bring back the components of color distortion (consisting of grayscale and color jitter) and evaluate them in detail. We find that random grayscale augmentation is surprisingly effective, given that the produced images may not be considered plausible in pathology. As the standard strength of color jittering produces highly unrealistic outputs, we evaluate a weaker jitter strength as well. Indeed, we find that while the performance drops substantially when using strong color jitter, using weak color jitter together with random grayscale results in the best performances. In particular, RandStainNA<sub>GMM</sub> shows high performance, motivating us to adopt it in our main experiments.

Through these observations, we substantiate our claim that existing augmentation schemes designed for SSL using natural images are sub-optimal for pathology data, necessitating pathology-specific augmentation schemes when training on pathology data such as TCGA and TULIP.

**Magnification Ablation.** In Sec. 3.2, we argue that pre-training using multiple magnifications or FoVs is beneficial as downstream pathology tasks occur at various magnifications. We do find experimentally that using multiple FoVs in the pre-training dataset is beneficial for overall downstream task performance (see Tab. 8).

Interestingly, we observe that using both 20× and 40× is best, while using only 20× is typically second-best. This is the case even for datasets such as PCam, MHIST, and CoNSeP which consist of images collected at approximately 40×. We hypothesize that the use of multiple magnifications is not valuable due to the matching of magnifications between upstream and downstream training, but rather due to the diversity of image appearances. Specifically, 20× images, due to the wider field-of-view, are visually and texture-wise more diverse than 40× images. Combining the two magnifications results in an even more diverse set of images. The more diverse data also results in better convergence during pre-training (see supplementary materials).

## 6. Discussion

In this section, we answer a few key questions that computational pathology researchers may naturally ask when considering self-supervised pre-training for their research.

**Should we pre-train on pathology data?** Yes – We have consistently demonstrated that pre-training on pathology data out-performs supervised pre-training on ImageNet by performing comprehensive experiments on many SSL methods and datasets. Interestingly, SSL pre-trained weights can maintain the performance gap on CoNSeP even for longer training schedules. Our experiments demystify and confirm the effectiveness of domain-aligned SSL pre-training on the pathology domain.

**Which SSL method is best?** We find that there is *no clear winner*. All SSL methods applied with domain-aligned pre-training generally perform well. Thus, instead of focusing on selecting a specific SSL method, we recommend that practitioners focus on curating large-scale domain-aligned datasets for SSL pre-training. Yet, some initial observations may be useful to future research. For example, (a) Barlow Twins tends to perform well in linear evaluations and MoCo v2 in fine-tuning evaluations, and (b) ViTs benefit more from domain-aligned SSL compared to CNNs.

**What is a key ingredient for successful self-supervised pre-training?** Domain knowledge – our proposed set of techniques are fully based on observations in pathology, and are experimentally shown to be effective. By incorporating domain-specific knowledge into the pre-training step, e.g., using stain augmentation and extracting patches from multiple FoVs, we go beyond the performance one can get from naively applying SSL to a new dataset.

## 7. Conclusion and Future Work

In this paper, we conduct the largest and most comprehensive study of SSL in the pathology domain, to date, using up to 33 million image patches during pre-training and evaluating 4 representative SSL methods (both CNNs and ViTs) on 2 downstream tasks and 5 datasets. Our study confirms that large-scale domain-aligned pre-training is helpful for pathology, showing its value in scenarios with limited labeled data, longer fine-tuning schedules, and when using larger and more diverse datasets for pre-training (such as TCGA + TULIP). Furthermore, we propose a set of techniques that are carefully designed by leveraging pathology-specific knowledge, and integrate them into the self-supervised pre-training stage, resulting in performance improvements. We believe that further exploration of domain-specific augmentation strategies will yield improved techniques for pathology-specific SSL in the future.

## References

- [1] Alaaeldin Ali, Hugo Touvron, Mathilde Caron, Piotr Bojanowski, Matthijs Douze, Armand Joulin, Ivan Laptev, Natalia Neverova, Gabriel Synnaeve, Jakob Verbeek, et al. Xcit: Cross-covariance image transformers. *NeurIPS*, 34:20014–20027, 2021. 14
- [2] Guilherme Aresta, Teresa Araújo, Scotty Kwok, Sai Saketh Chennamsetty, Mohammed Safwan, Varghese Alex, Bahram Marami, Marcel Prastawa, Monica Chan, Michael Donovan, et al. Bach: Grand challenge on breast cancer histology images. *Medical Image Analysis*, 56:122–139, 2019. 2, 4, 11
- [3] Shekoofeh Azizi, Basil Mustafa, Fiona Ryan, Zachary Beaver, Jan Freyberg, Jonathan Deaton, Aaron Loh, Alan Karthikesalingam, Simon Kornblith, Ting Chen, et al. Big self-supervised models advance medical image classification. In *ICCV*, pages 3478–3488, 2021. 2, 6
- [4] Babak Ehteshami Bejnordi, Mitko Veta, Paul Johannes Van Diest, Bram Van Ginneken, Nico Karssemeijer, Geert Litjens, Jeroen AWM Van Der Laak, Meyke Hermsen, Quirine F Manson, Maschenka Balkenhol, et al. Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer. *Jama*, 318(22):2199–2210, 2017. 12, 17
- [5] Joseph Boyd, Mykola Liashuba, Eric Deutsch, Nikos Paragios, Stergios Christodoulidis, and Maria Vakalopoulou. Self-supervised representation learning using visual field expansion on digital pathology. In *ICCV Workshops*, pages 639–647, 2021. 2
- [6] Wouter Bulten, Kimmo Kartasalo, Po-Hsuan Cameron Chen, Peter Ström, Hans Pinckaers, Kunal Nagpal, Yuannan Cai, David F Steiner, Hester van Boven, Robert Vink, et al. Artificial intelligence for diagnosis and gleason grading of prostate cancer: the panda challenge. *Nature medicine*, 28(1):154–163, 2022. 1, 3
- [7] Gabriele Campanella, Matthew G Hanna, Luke Geneslaw, Allen Miraflor, Vitor Werneck Krauss Silva, Klaus J Busam, Edi Brogi, Victor E Reuter, David S Klimstra, and Thomas J Fuchs. Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nature medicine*, 25(8):1301–1309, 2019. 1
- [8] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *ECCV*, pages 132–149, 2018. 2
- [9] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *NeurIPS*, 33:9912–9924, 2020. 2, 4, 13
- [10] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, pages 9650–9660, 2021. 1, 2, 4, 6, 13
- [11] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, pages 801–818, 2018. 12
- [12] Richard J Chen, Chengkuan Chen, Yicong Li, Tiffany Y Chen, Andrew D Trister, Rahul G Krishnan, and Faisal Mah-
- mood. Scaling vision transformers to gigapixel images via hierarchical self-supervised learning. In *CVPR*, pages 16144–16155, 2022. 2, 17
- [13] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607. PMLR, 2020. 1, 2, 4, 5, 7, 14
- [14] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E Hinton. Big self-supervised models are strong semi-supervised learners. *NeurIPS*, 33:22243–22255, 2020. 6
- [15] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 2, 4, 13
- [16] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *CVPR*, pages 15750–15758, 2021. 2, 4
- [17] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *ICCV*, pages 9640–9649, 2021. 2
- [18] Sangjoon Choi, Soo Ick Cho, Minuk Ma, Seonwook Park, Sergio Pereira, Brian Jaehong Aum, Seunghwan Shin, Kyunghyun Paeng, Donggeun Yoo, Wonkyung Jung, et al. Artificial intelligence-powered programmed death ligand 1 analyser reduces interobserver variation in tumour proportion score for non–small cell lung cancer with better prediction of immunotherapy response. *European Journal of Cancer*, 170:17–26, 2022. 1
- [19] Ozan Ciga, Tony Xu, and Anne Louise Martel. Self supervised contrastive learning for digital histopathology. *Machine Learning with Applications*, 7:100198, 2022. 2, 4
- [20] James Coughlan and Alan L Yuille. The manhattan world assumption: Regularities in scene statistics which enable bayesian inference. In *NeurIPS*, volume 13, 2000. 3
- [21] Angel Cruz-Roa, Hannah Gilmore, Ajay Basavanhally, Michael Feldman, Shridhar Ganesan, Natalie NC Shih, John Tomaszewski, Fabio A González, and Anant Madabhushi. Accurate and reproducible invasive breast cancer detection in whole-slide images: A deep learning approach for quantifying tumor extent. *Scientific reports*, 7(1):1–14, 2017. 4
- [22] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. IEEE, 2009. 1, 3
- [23] James A Diao, Jason K Wang, Wan Fung Chui, Victoria Mountain, Sai Chowdary Gullapally, Ramprakash Srinivasan, Richard N Mitchell, Benjamin Glass, Sara Hoffman, Sudha K Rao, et al. Human-interpretable image features derived from densely mapped cancer pathology slides predict diverse molecular phenotypes. *Nature communications*, 12(1):1–15, 2021. 1, 4
- [24] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020. 2, 4
- [25] Khrystyna Faryna, Jeroen van der Laak, and Geert Litjens. Tailoring automated data augmentation to h&e-stained

- histopathology. In *Medical Imaging with Deep Learning*, 2021. 7
- [26] Jacob Gildenblat and Eldad Klaiman. Self-supervised similarity learning for digital pathology. In *MICCAI Workshops*, 2019. 2
- [27] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large mini-batch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017. 4
- [28] Priya Goyal, Quentin Duval, Jeremy Reizenstein, Matthew Leavitt, Min Xu, Benjamin Lefauzeux, Mannat Singh, Vini-cius Reis, Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Ishan Misra. Vissl. <https://github.com/facebookresearch/vissl>, 2021. 13
- [29] Simon Graham, Quoc Dang Vu, Shan E Ahmed Raza, Ayesha Azam, Yee Wah Tsang, Jin Tae Kwak, and Nasir Rajpoot. Hover-net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images. *Medical Image Analysis*, 58:101563, 2019. 2, 3, 4, 5, 12, 14, 16
- [30] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *NeurIPS*, 33:21271–21284, 2020. 1, 2, 4, 5, 14
- [31] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, pages 9729–9738, 2020. 2, 7
- [32] Kaiming He, Ross Girshick, and Piotr Dollár. Rethinking imagenet pre-training. In *ICCV*, pages 4918–4927, 2019. 7
- [33] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 4, 14
- [34] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *ECCV*, pages 630–645. Springer, 2016. 14
- [35] Vladimir Iglovikov and Alexey Shvets. Ternausnet: U-net with vgg11 encoder pre-trained on imagenet for image segmentation. *arXiv preprint arXiv:1801.05746*, 2018. 1
- [36] Jakob Nikolas Kather, Niels Halama, and Alexander Marx. 100,000 histological images of human colorectal cancer and healthy tissue, Apr. 2018. 2, 4, 11
- [37] 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. In *Proceedings of the Conference on Health, Inference, and Learning*, pages 116–124, 2021. 1
- [38] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *CVPR*, pages 9404–9413, 2019. 5
- [39] Simon Kornblith, Jonathon Shlens, and Quoc V Le. Do better imagenet models transfer better? In *CVPR*, pages 2661–2671, 2019. 1
- [40] Bin Li, Yin Li, and Kevin W Eliceiri. Dual-stream multiple instance learning network for whole slide image classification with self-supervised contrastive learning. In *CVPR*, pages 14318–14328, 2021. 2
- [41] Chunyuan Li, Jianwei Yang, Pengchuan Zhang, Mei Gao, Bin Xiao, Xiyang Dai, Lu Yuan, and Jianfeng Gao. Efficient self-supervised vision transformers for representation learning. *arXiv preprint arXiv:2106.09785*, 2021. 2
- [42] Geert Litjens, Peter Bandi, Babak Ehteshami Bejnordi, Oscar Geessink, Maschenka Balkenhol, Peter Bult, Altuna Halilovic, Meyke HermSEN, Rob van de Loo, Rob Vogels, et al. 1399 h&e-stained sentinel lymph node sections of breast cancer patients: the camelyon dataset. *GigaScience*, 7(6):giy065, 2018. 1
- [43] John Lonsdale, Jeffrey Thomas, Mike Salvatore, Rebecca Phillips, Edmund Lo, Saboor Shad, Richard Hasz, Gary Walters, Fernando Garcia, Nancy Young, et al. The genotype-tissue expression (gtex) project. *Nature genetics*, 45(6):580–585, 2013. 1
- [44] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. 2018. 13
- [45] Ming Y Lu, Drew FK Williamson, Tiffany Y Chen, Richard J Chen, Matteo Barbieri, and Faisal Mahmood. Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature biomedical engineering*, 5(6):555–570, 2021. 17
- [46] Marc Macenko, Marc Niethammer, James S Marron, David Borland, John T Woosley, Xiaojun Guan, Charles Schmitt, and Nancy E Thomas. A method for normalizing histology slides for quantitative analysis. In *ISBI*, pages 1107–1110. IEEE, 2009. 11
- [47] Christos Matsoukas, Johan Fredin Haslum, Moein Sorkhei, Magnus Söderberg, and Kevin Smith. What makes transfer learning work for medical images: Feature reuse & other factors. In *CVPR*, pages 9225–9234, 2022. 1, 2
- [48] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *CVPR*, pages 6707–6717, 2020. 2
- [49] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. In *NeurIPS*, 2018. 2
- [50] Sehhoon Park, Chan-Young Ock, Hyojin Kim, Sergio Pereira, Seonwook Park, Minuk Ma, Sangjoon Choi, Seokhwi Kim, Seunghwan Shin, Brian Jaehong Aum, et al. Artificial intelligence-powered spatial analysis of tumor-infiltrating lymphocytes as complementary biomarker for immune checkpoint inhibition in non-small-cell lung cancer. *Journal of Clinical Oncology*, 40(17):1916, 2022. 1, 4
- [51] Maithra Raghu, Chiyuan Zhang, Jon Kleinberg, and Samy Bengio. Transfusion: Understanding transfer learning for medical imaging. *NeurIPS*, 32, 2019. 1
- [52] Erik Reinhard, Michael Adikhmin, Bruce Gooch, and Peter Shirley. Color transfer between images. *IEEE Computer graphics and applications*, 21(5):34–41, 2001. 3, 13
- [53] Byungseok Roh, Wuhyun Shin, Ildoo Kim, and Sungwoong Kim. Spatially consistent representation learning. In *CVPR*, pages 1144–1153, 2021. 7
- [54] Arnout C Ruifrok, Dennis A Johnston, et al. Quantification of histochemical staining by color deconvolution. *Analyti-*

- cal and quantitative cytology and histology*, 23(4):291–299, 2001. 3
- [55] Yiqing Shen, Yulin Luo, Dinggang Shen, and Jing Ke. Randstainna: Learning stain-agnostic features from histology slides by bridging stain augmentation and normalization. In *MICCAI*, pages 212–221. Springer, 2022. 3, 4, 8, 13
- [56] Hari Sowrirajan, Jingbo Yang, Andrew Y Ng, and Pranav Rajpurkar. Moco pretraining improves representation and transferability of chest x-ray models. In *Medical Imaging with Deep Learning*, pages 728–744. PMLR, 2021. 2
- [57] Fabio A Spanhol, Luiz S Oliveira, Caroline Petitjean, and Laurent Heutte. A dataset for breast cancer histopathological image classification. *Ieee transactions on biomedical engineering*, 63(7):1455–1462, 2015. 1
- [58] David Tellez, Geert Litjens, Péter Bárdi, Wouter Bulten, John-Melle Bokhorst, Francesco Ciompi, and Jeroen Van Der Laak. Quantifying the effects of data augmentation and stain color normalization in convolutional neural networks for computational pathology. *Medical image analysis*, 58:101544, 2019. 3, 4, 8
- [59] Yonglong Tian, Olivier J Henaff, and Aäron van den Oord. Divide and contrast: Self-supervised learning from uncultured data. In *ICCV*, pages 10063–10074, 2021. 4
- [60] Jeroen Van der Laak, Geert Litjens, and Francesco Ciompi. Deep learning in histopathology: the path to the clinic. *Nature medicine*, 27(5):775–784, 2021. 1
- [61] Wouter Van Gansbeke, Simon Vandenhende, Stamatis Georgoulis, and Luc V Gool. Revisiting contrastive methods for unsupervised learning of visual representations. *NeurIPS*, 34:16238–16250, 2021. 6
- [62] Bastiaan S Veeling, Jasper Linmans, Jim Winkens, Taco Cohen, and Max Welling. Rotation equivariant cnns for digital pathology. In *MICCAI*, pages 210–218. Springer, 2018. 2, 4, 12
- [63] Xiye Wang, Sen Yang, Jun Zhang, Minghui Wang, Jing Zhang, Junzhou Huang, Wei Yang, and Xiao Han. Transpath: Transformer-based self-supervised learning for histopathological image classification. In *MICCAI*, pages 186–195. Springer, 2021. 2
- [64] Jerry Wei, Arief Suriawinata, Bing Ren, Xiaoying Liu, Mikhail Lisovsky, Louis Vaickus, Charles Brown, Michael Baker, Naofumi Tomita, Lorenzo Torresani, et al. A petri dish for histopathology image analysis. In *International Conference on Artificial Intelligence in Medicine*, pages 11–24. Springer, 2021. 2, 4, 12
- [65] John N Weinstein, Eric A Collisson, Gordon B Mills, Kenna R Shaw, Brad A Ozenberger, Kyle Ellrott, Ilya Shmulevich, Chris Sander, and Joshua M Stuart. The cancer genome atlas pan-cancer analysis project. *Nature genetics*, 45(10):1113–1120, 2013. 1, 2
- [66] Tete Xiao, Colorado J Reed, Xiaolong Wang, Kurt Keutzer, and Trevor Darrell. Region similarity representation learning. In *ICCV*, pages 10539–10548, 2021. 7
- [67] Yiting Xie and David Richmond. Pre-training on grayscale imagenet improves medical image classification. In *ECCV Workshops*, pages 0–0, 2018. 1
- [68] Jiawei Yang, Hanbo Chen, Yuan Liang, Junzhou Huang, Lei He, and Jianhua Yao. Concl: Concept contrastive learning for dense prediction pre-training in pathology images. In *ECCV*, pages 523–539, 2022. 2
- [69] Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017. 13
- [70] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *ICML*, pages 12310–12320. PMLR, 2021. 2, 4, 5, 7, 13, 14

## Appendix

**Overview.** In this supplementary material, we describe the details of the downstream datasets adopted in the main paper and show some example images. This document also contains further implementation details regarding the pre-training and downstream training steps, including fine-tuning with limited labeled data. Last but not least, we provide further analyses, such as the effectiveness of pre-training for longer epochs and pre-training stability when using data from different magnifications.

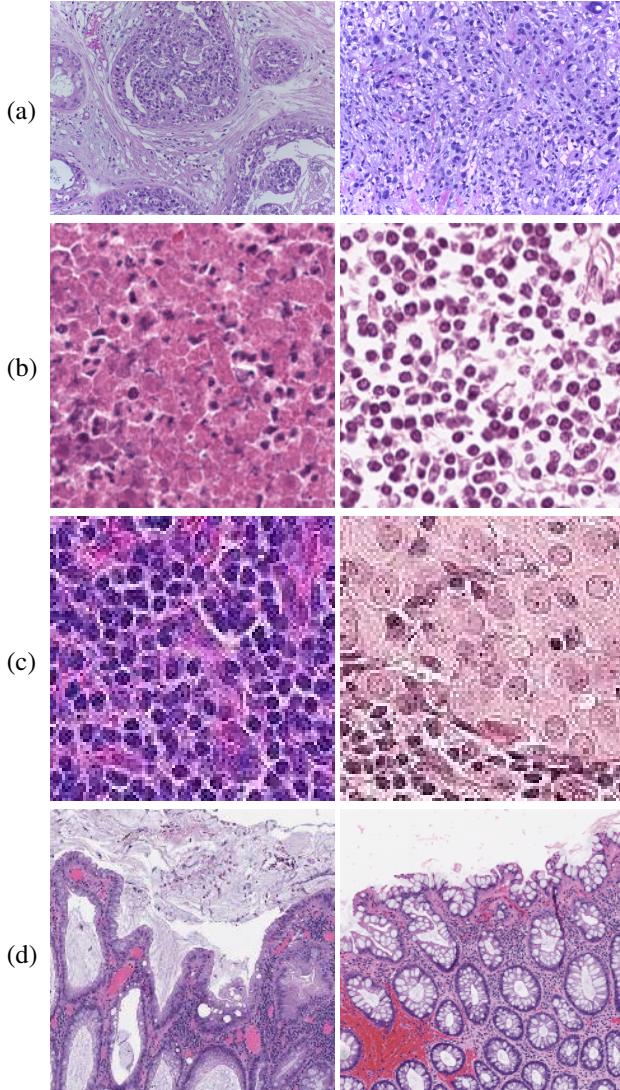
Note that the corresponding or relevant sections from the main paper are referenced Note that the corresponding or relevant sections from the main paper are referenced **in blue text** in the section titles.

## A. Downstream Dataset Details (Section 4.2)

In this section, we describe the details of the datasets used in our analysis. We use BACH, CRC, PCam, and MHIST for the image classification task, and CoNSeP for the nuclei instance segmentation task. We sample a few training images from each dataset and present them in Fig. A.1 and Fig. A.2. **in blue text** in the section titles.

**BACH.** The goal of the Grand Challenge on BreAst Cancer Histology (BACH) [2] is to classify pathology images into four classes: normal, benign, in situ carcinoma, and invasive carcinoma. The dataset is composed of 400 training images and 100 test images. The test images are collected from a different set of patients from the training images. All images are collected from Hospital CUF Porto, Centro Hospitalar do Tâmega e Sousa, and Centro Hospitalar Cova da Beira.

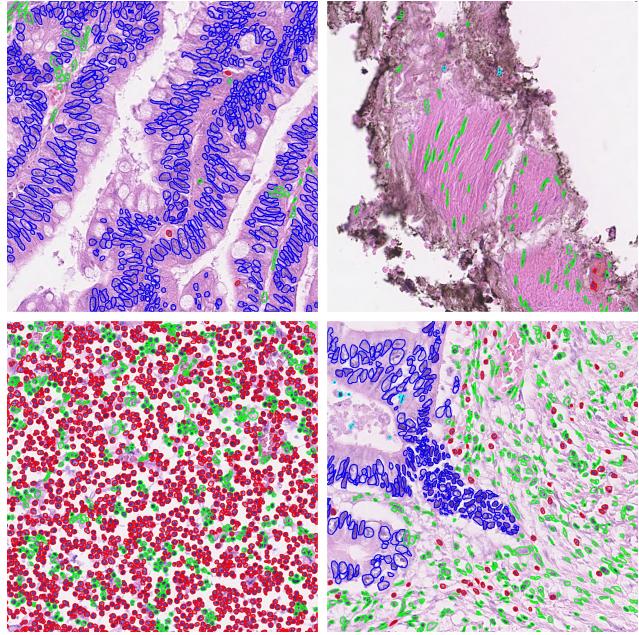
**CRC.** This dataset [36] consists of 100,000 training images and 7,180 test images from H&E stained WSIs of human colorectal cancer (CRC) and normal tissue. The training and test images are extracted from 86 WSIs and 25 WSIs, respectively. The slides are collected from the NCT Tissue Bank and the University Medical Center Mannheim. The task is the identification of nine tissue classes: adipose tissue, background, debris, lymphocytes, mucus, smooth muscle, normal colon mucosa, cancer-associated stroma, and CRC epithelium. All images are color normalized with the Macenko method [46].



**Figure A.1. Example training images from the classification datasets:** (a) BACH, (b) CRC, (c) PCam, and (d) MHIST.

**PCam.** The PatchCamelyon (PCam) [62] dataset is derived from the Camelyon16 [4] dataset that contains 400 H&E stained WSIs from two hospitals: Radboud University Medical Center (RUMC), and University Medical Center Utrecht (UMCU). The PCam dataset includes 262,144 training images, 32,768 validation images, and 32,768 test images. Each image is annotated with a binary label for determining the presence of metastases.

**MHIST.** The minimalist histopathology image analysis (MHIST) [64] dataset is comprised of 2,175 training images and 977 test images. The images are extracted from 328 H&E stained Formalin Fixed Paraffin-Embedded (FFPE) WSIs of colorectal polyps from Dartmouth-Hitchcock Medical Center. The task is the binary classification between hyperplastic polyps (HPs) and sessile serrated adenomas



**Figure A.2. Example training images from the CoNSeP dataset.** The dataset provides annotated nuclei masks along with cell type labels. Following the original HoVer-Net paper [29], we use the following nuclei types for training and evaluation: ■ epithelial, ■ inflammatory, ■ spindle-shaped, and ■ miscellaneous.

(SSAs), where HPs are benign and SSAs are precancerous lesions.

**CoNSeP.** The Colorectal Nuclear Segmentation and Phenotypes (CoNSeP) dataset [29] consists of 41 H&E images and is split into 27 images and 14 images for training and test sets, respectively. The data comes from University Hospitals Coventry and Warwickshire, UK (UHCW). The annotation contains segmentation masks of each nucleus along with its class (See Fig. A.2). Note that the healthy epithelial and dysplastic/malignant epithelial are considered general epithelial types. Fibroblast, muscle, and endothelial are matched into a spindle-shaped nuclei type. In total, 24,319 unique nuclei masks along with 4 major types out of 7 cell types are used during training.

## B. Implementation Details

In the interest of improving the reproducibility of our study, we provide further details regarding our pre-training data, setup, as well as details on how we conducted our downstream evaluations. Furthermore, we discuss the details of the limited labeled data experiments.

### B.1. Preparation of Pre-training Data (Section 4.1)

In selecting image patches to compose the TCGA dataset, we first use an internal model with a DeepLab v3+ architecture [11] to segment the foreground regions of WSI. From

the candidate patches that are located in areas predicted as foreground, we select up to 500 patches per magnification, per slide, with equal spacing between them. To ensure that we have informative image patches in our pre-training dataset, we filter out patches that are too white (mean saturation value below 5) or too smooth (mean squared Laplacian below 15). For TULIP, we do not apply such filtering logic due to the relatively smaller foreground area (too many patches are lost otherwise).

## B.2. Calculating Statistics of the Pre-training Data (Section 4.1)

For the purpose of input image standardization during SSL pre-training, we collect the per-channel mean and standard deviation of intensities in RGB space, using 10% of the full unlabeled image data. This subsampling is done per WSI, to maintain diversity and reduce computational cost.

In a similar manner, we compute the per-channel means and variances in 3 color spaces (HSV, Lab, HED) for use with the RandStainNA method, using 10% of the full image data. Specifically, we compute per color space, and per channel, the mean and standard deviation of per-image mean intensity, as well as the mean and standard deviation of the per-image standard deviation of intensity. Please refer to [55] for further details.

For RandStainNA<sub>GMM</sub>, we similarly compute per color space, and per channel, the per-image mean intensity and its standard deviation. However, instead of simply finding the mean and standard deviation of those values independently (fitting individual unimodal Gaussian distributions 18 times as in RandStainNA), we fit a 10-component Gaussian Mixture Model (GMM) for each color space, yielding 3 models. This is done to fit the covariance between the input variables (6 variables exist for each color space) and respect their multi-modal nature.

## B.3. Augmentation Details (Section 3.2)

Unless otherwise stated, in our experiments, we pre-train by applying the following changes to the default method-specific augmentation scheme:

- Random vertical flip ( $p=0.5$ ).
- Color dropping ( $p=0.2$ ): the color of images are converted randomly to grayscale.
- Weak color jittering ( $p=0.8$ ): the brightness, contrast, saturation, and hue of images are randomly adjusted with a strength of 0.2, 0.2, 0.2, 0.1, respectively.
- RandStainNA<sub>GMM</sub> ( $p=0.8$ ): per image, a color space is randomly selected (from HSV, Lab, or HED), then channel-wise mean and standard deviation values are sampled from a GMM which is fitted on statistics

from part of the pre-training data (10%). The input image is re-normalized based on these values, using Reinhard’s method [52].

## B.4. SSL Methods (Section 4.3)

We provide implementation details of each SSL method used in our analysis. We use the VISSL [28] library to pre-train the the 4 studied SSL methods, and follow the same configurations as originally proposed in [9, 10, 15, 70]. All representations are trained for 200 ImageNet epochs, distributed over 64 V100 16GB GPUs. A linear warmup schedule is applied for the first 10 epochs and a cosine learning rate decay is applied subsequently. Each method was originally proposed with its specific augmentation schemes, and we follow those original data augmentation pipelines while adding our proposed techniques on top. Regarding the RandStainNA augmentation, it requires the statistics of 3 color spaces (HSV, Lab, HED) to produce augmented images. To compute the statistics, we randomly sample 10% of the unlabeled image patches from the corresponding pre-training dataset.

**MoCo v2.** We use the SGD optimizer with an initial learning rate of 0.3. The learning rate is linearly scaled up based on  $lr = lr * batchsize/256$ , where  $batchsize$  is 4,096. The memory bank size is fixed to 65,536, and a momentum coefficient  $m$  of 0.999 is used. Weight decay of  $10^{-4}$  is utilized for regularization.

**SwAV.** We use the SGD optimizer with an initial learning rate of 0.3. The learning rate is linearly scaled up based on  $lr = lr * batchsize/256$ , where  $batchsize$  is 2,048. The number of prototypes is 3,000 to avoid intractable computational costs from the Sinkhorn algorithm.  $2 \times 224 + 6 \times 96$  multi-crop augmentation is employed as done in the original paper.

**Barlow Twins.** The LARS optimizer [69] is adopted for Barlow Twins pre-training. Note that, as in the original work [70], we apply different learning rates for weights and biases, 0.2 and 0.0048, respectively. The biases and batch normalization layers are excluded from LARS optimization to follow the original implementation. The learning rates of weights and biases are linearly scaled up based on  $lr = lr * batchsize/256$ , where  $batchsize$  is 2,048. The dimension of the embeddings is 8,192, and training is conducted with a coefficient of off-diagonal term  $\lambda = 5 \cdot 10^{-3}$  and a weight decay of  $1.5 \cdot 10^{-6}$ .

**DINO.** We train the model with the AdamW [44] optimizer. The learning rate of 0.0005 is used for stability during pre-training. The learning rate is linearly scaled up based on  $lr = lr * batchsize/256$ , where  $batchsize$  is 1,024. Similar to the learning rate decay, the weight decay also follows a cosine schedule from 0.04 to 0.4. For DINO<sub>p=16</sub>,  $2 \times 224 + 8 \times 96$

multi-crop augmentation is employed, while  $2 \times 224 + 6 \times 96$  multi-crop augmentation is used for  $\text{DINO}_{p=8}$ .

## B.5. Downstream Training Details (Section 4.4)

**Image Classification.** We split each downstream dataset into training, validation, and test sets. The learning rate and weight decay values are optimized using training and validation, only. In the BACH dataset, the labels for the test set are not provided. Hence, we split the training set by a 6:1:3 ratio (training, validation, test). For the CRC and MHIST datasets, the test set is provided with labels, and the training set is split by a 7:3 ratio (training, validation). For the PCam dataset, we follow the original data split. When splitting the data, we do it randomly but in a class-balanced manner. Based on the performance measured on the validation sets, we perform a grid search of learning rates from  $\{1, 0.1, 0.01, 0.001\}$  and weight decay values from  $\{0.1, 0.01, 0.001, 0\}$ .

As data augmentation for ResNet-50, the input image is randomly flipped both horizontally and vertically, at training time. For the BACH dataset, we apply random cropping and resizing to  $1024 \times 768$  at training time; at test time, we resize the images to  $1024 \times 768$ . For ViT-S, the same augmentation is used but all images are resized to  $224 \times 224$ . We train the models with the SGD optimizer with a momentum of 0.9 and a cosine learning rate decay. The ResNet-50 based models are trained for 200, 20, 20, and 90 epochs on the BACH, CRC, PCam, and MHIST datasets, respectively. The Transformer-based models are trained for 30 epochs on the CRC and PCam datasets and for 200 and 90 epochs on the BACH and MHIST datasets, respectively. During fine-tuning, the backbone layers (i.e., ResNet-50 and ViT-S) are trained with a learning rate 100 times lower than that of the last classification layer.

**Nuclei Instance Segmentation.** We follow the standard pipeline of HoVer-Net [29], as provided in its open-source implementation<sup>3</sup>, including data augmentation and patch extraction. HoVer-Net defines a two-stage training procedure. At the first stage, only the decoders are trained while freezing the backbone layers. With the trained decoders, all layers are then fine-tuned at the second stage. Technically, Preact-ResNet-50 [34] is employed in the original implementation, but we replace it with the standard ResNet-50 [33] and reproduce the results for a fair comparison. This is done to perform SSL pre-training in a standard manner while permitting this nuclei instance segmentation downstream task. Since we change the backbone, we perform Grid Search to find a proper learning rate. We use  $5 \cdot 10^{-4}$  learning rate for both stages of HoVer-Net. Moreover, based on the open-source implementation, the authors of HoVer-Net fine-tune the first convolutional layer of ResNet at the

<sup>3</sup>[https://github.com/vqdang/hover\\_net](https://github.com/vqdang/hover_net)

first stage, but we keep them frozen.

Similar to the architecture of FPN-based instance segmentation, HoVer-Net requires features from multiple scales in the encoder. However, the outputs of the ViT-based encoder are not compatible with the existing decoders of Hover-Net without further modifications. In order to provide multi-scale features to the decoder, we refer to the protocol from [1] where the feature scales are interpolated using several transposed convolution layers with kernel size  $k = 2$  and stride  $s = 2$ . More specifically, features from the 4<sup>th</sup>, 6<sup>th</sup>, 8<sup>th</sup>, and 12<sup>th</sup> layers are extracted from the ViT-S architecture, which consists of 12 layers in total. With this design, the decoders remain unchanged. For the sake of a fair comparison, we also perform Grid Search on the ViT-S architecture. The learning rate of  $5 \cdot 10^{-4}$  is used for both stages of HoVer-Net.

## B.6. Fine-tuning with Limited Labeled Data (Section 5.4)

**Image Classification.** Following prior works [13, 30, 70], we randomly sample 1% and 10% of the CRC training set by balancing classes. Based on the fine-tuning procedure, we train the models for 60 and 90 epochs for 1% and 10% labeled data, respectively.

**Nuclei Instance Segmentation using CoNSeP.** In our limited labeled data experiments using CoNSeP, we control the number of H&E images instead of the number of extracted patches to mimic the real-world setting where one H&E image corresponds to one unique patient. Since assuming 1% of training data is unreasonable in the current setting (i.e., 0.27 H&E image), instead, we define the ratio of 10% and 30% for nuclei instance segmentation. Note that the reported values in the experiments are the averaged number from 3 repetitive experiments with different seed values for image/patient selection. This is necessitated by the smaller dataset size (compared to CRC) and is done for a fair comparison between methods.

## C. Pre-training for more epochs (Section 5)

Typically, increasing the number of pre-training epochs has shown to be effective in improving the learned representations in various SSL methods. To investigate the effectiveness of the longer pre-training in the pathology domain, we pre-train the model for 800 ImageNet epochs using MoCo v2, SwAV, and Barlow Twins. Note that, due to computational costs, we report the results from  $\text{DINO}_{p=16}$  trained for 400 ImageNet epochs.

Tab. B.1 and Tab. C.1 present the performance of image classification and nuclei instance segmentation, respectively. Compared to the results from 200 ImageNet epochs, SwAV is the only method that benefits from the longer pre-training in the fine-tuning protocol, especially in BACH,

Arch.	Method	BACH		CRC		PCam		MHIST	
		Linear	Fine-tune	Linear	Fine-tune	Linear	Fine-tune	Linear	Fine-tune
ResNet-50	Random	51.67	61.67	68.91	89.99	76.52	75.71	63.15	75.54
	Supervised	80.83	86.67	90.93	92.09	80.79	80.63	76.25	78.92
	<b>Epoch 200</b>								
	MoCo v2	77.50	<u>90.83</u>	93.52	<b>96.21</b>	86.78	<b>87.62</b>	77.07	<b>85.88</b>
	SwAV	<u>83.33</u>	82.50	<u>95.78</u>	93.31	85.28	<u>87.60</u>	71.14	77.99
	BT	<b>87.50</b>	85.00	94.60	93.23	<b>88.15</b>	86.92	<b>78.81</b>	81.27
	<b>Epoch 800</b>								
	MoCo v2	79.17	<b>91.67</b>	95.01	<u>95.45</u>	87.84	86.90	72.77	84.95
	SwAV	82.50	85.83	<b>96.46</b>	<u>92.74</u>	<u>86.16</u>	87.05	75.54	<u>85.47</u>
	BT	86.67	<b>91.67</b>	94.48	94.99	86.26	86.75	<u>78.20</u>	<u>80.25</u>
ViT-S	$Random_{p=16}$	45.00	57.50	69.90	86.10	74.43	75.42	63.46	62.13
	$Supervised_{p=16}$	75.83	85.83	91.56	<u>95.81</u>	80.96	88.30	<b>78.51</b>	<b>81.68</b>
	<b>Epoch 200</b>								
	$DINO_{p=16}$	<u>85.83</u>	<u>87.50</u>	<u>94.19</u>	<u>95.81</u>	<b>88.78</b>	<u>90.40</u>	<u>76.15</u>	79.43
	<b>Epoch 400</b>								
	$DINO_{p=16}$	<b>86.67</b>	<b>88.33</b>	<b>95.13</b>	<b>96.48</b>	<u>88.60</u>	89.50	75.44	<u>81.06</u>

Table B.1. **Downstream evaluation of image classification tasks under a different number of pre-training epochs.** We report Top-1 accuracy for both linear and fine-tuning experiment protocols trained using the TCGA data source. Note that  $p$  represents the patch size used in ViT. We compare results column-wise and mark the best results in **bold** and the second-best results in underline for ResNet-50 based methods and ViT-S methods separately.

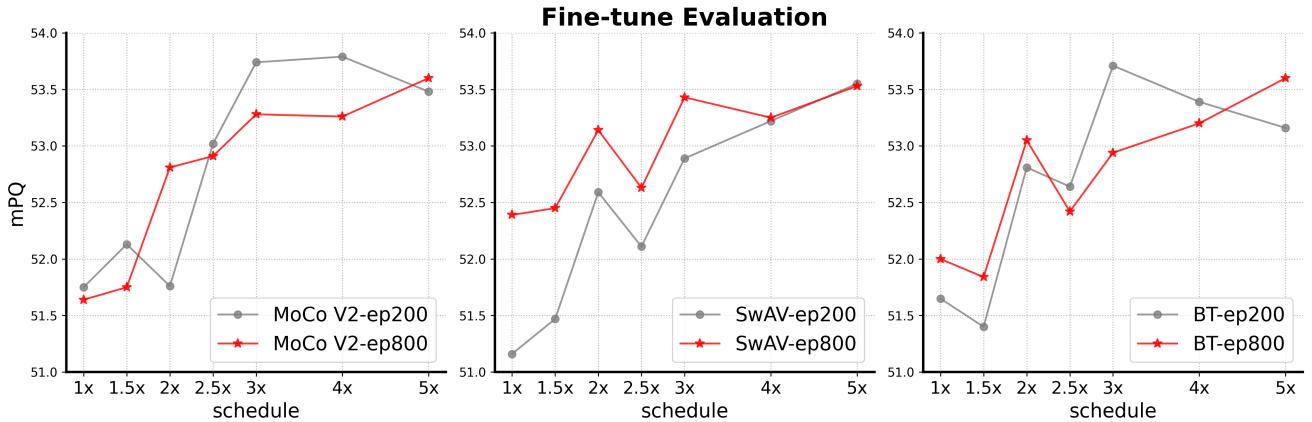


Figure C.1. **The effectiveness of longer pre-training according to learning schedules.** We present fine-tuning evaluation results for the nuclei instance segmentation task using the CoNSeP dataset. We see that there are few differences between the 200 epoch models and 800 epoch models, except that SwAV benefits from longer pre-training when the downstream task is fine-tuned with a limited learning schedule.

MHIST, and CoNSeP datasets. In contrast, the other methods show marginal improvements or are on par with the 200 ImageNet epoch counterparts.  $DINO_{p=16}$  shows a slightly improved performance on image classification, while nuclei instance segmentation remains on par. Even in the different learning schedules illustrated in Fig. C.1, we observe that no clear benefit of the longer pre-training stands out in MoCo v2 and Barlow Twins, yet SwAV consistently maintains the

benefit of the longer pre-training.

Across all experiments, we confirm that certain SSL methods (e.g., SwAV) may require more pre-training iterations, but generally increasing the number of pre-training epochs shows marginal improvements on both image classification and nuclei instance segmentation tasks. In other words, pre-training for 200 ImageNet epochs can be sufficient to achieve satisfactory downstream performance, espe-

Arch.	Method	CoNSeP	
		Linear	Fine-tune
ResNet-50	<i>Random</i>	22.29	46.72
	<i>Supervised</i>	34.25	49.60
	<b>Epoch 200</b>		
	MoCo v2	39.85	51.75
	SwAV	40.45	51.16
	BT	<u>40.79</u>	51.61
	<b>Epoch 800</b>		
	MoCo v2	<b>40.93</b>	51.64
	SwAV	40.59	<b>52.39</b>
ViT-S	BT	<u>40.65</u>	52.00
	<i>Random</i> <sub>p=16</sub>	20.55	27.19
	<i>Supervised</i> <sub>p=16</sub>	21.43	36.70
	<b>Epoch 200</b>		
	DINO <sub>p=16</sub>	<u>32.54</u>	<u>38.43</u>
<b>Epoch 400</b>	DINO <sub>p=16</sub>	<b>32.93</b>	<b>39.03</b>

Table C.1. **Downstream evaluation for the nuclei instance segmentation task under a different number of pre-training epochs.** We report the mPQ score for both linear and fine-tuning experiment protocols for models trained using the TCGA data source. We compare results column-wise and mark the best results in **bold** and the second-best results in underline for ResNet-50 based methods and ViT-S methods separately.

cially for MoCo v2, Barlow Twins, and DINO. We therefore suggest that using 200 ImageNet epochs would be adequate to study the potential of SSL pre-training in the pathology domain.

## D. Pre-training Stability with Different Magnifications (Section 5.6)

In the main paper, we show that it is beneficial to train on image data from a combination of 20× and 40× objective magnifications. Here, we show that pre-training stability is also affected by the choice of magnification. In Fig. D.1, we present the loss trajectory during the pre-training stage using Barlow Twins. As shown in the graph, using a single magnification produces unstable losses and the loss begins to converge after approximately 4,000 and 7,000 iterations for magnifications of 20× and 40×, respectively. The loss values at the end of the pre-training stage are also higher in the case of using a single magnification. In contrast, using multiple magnifications results in stable pre-training and fast convergence, in addition to improved downstream task performance.

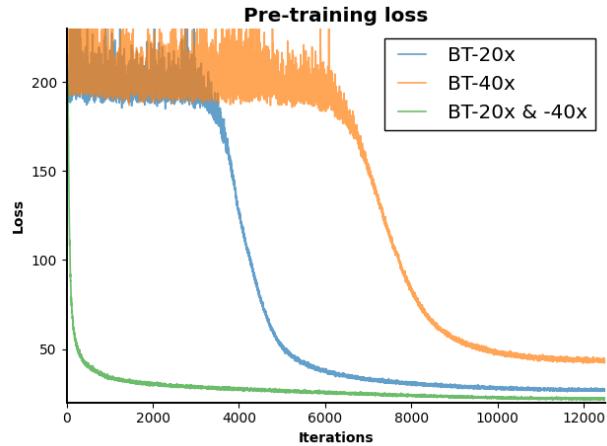


Figure D.1. **Loss progression while pre-training Barlow Twins on different magnifications.** Training on a combination of 20× and 40× results in quick convergence and stable pre-training.

## E. Larger Inputs for ViT (Section 5.2)

The implementation of the standard HoVer-Net [29] method involves the fine-tuning of a pre-trained ResNet, using images with a resolution of 270 × 270. However, by design, ViT expects input images of 224 × 224 resolution. Given the potential advantages that larger input resolutions can bring to the task of nuclei instance segmentation, we adopt a positional embedding interpolation technique to increase the input image size to 272 × 272, which is divisible by both 16 and 8. Through this technique, we aim to maintain consistent input resolutions across the ResNet and ViT backbones being evaluated. Tab. E.1 presents the result according to the input size. We observe that the larger input size improves performance for DINO<sub>p=16</sub>, while the performance of DINO<sub>p=8</sub> reduces.

## F. Further Data Augmentation Ablation Study (Section 5.6)

To provide a compelling demonstration of the effectiveness of the proposed techniques, we opted for the most practical, yet challenging fine-tuning setting of nuclei instance segmentation. Through the application of the linear evaluation protocol, we further validate the effectiveness of our techniques by showcasing improvements across all datasets. Notably, our set of techniques consistently and significantly improves the performance compared to the baseline approach that relies on augmentations designed for natural images. The improvement presented in Tab. F.1 serves as a clear signal of the effectiveness of our proposed techniques, which were carefully designed with the aid of domain-specific knowledge.

Arch.	Method	CoNSeP	
		Linear	Fine-tune
<b><u>224 input</u></b>			
ViT-S	<i>Supervised</i> <sub>p=16</sub>	21.43	36.70
	DINO <sub>p=16</sub>	<u>32.54</u>	<u>38.43</u>
	DINO <sub>p=8</sub>	<b>42.71</b>	<b>46.70</b>
<b><u>272 input</u></b>			
	<i>Supervised</i> <sub>p=16</sub>	28.60	34.50
	DINO <sub>p=16</sub>	<u>35.81</u>	<u>41.13</u>
	DINO <sub>p=8</sub>	<b>40.08</b>	<b>44.24</b>

Table E.1. **Downstream evaluation for the nuclei instance segmentation task under a different input resolution.** We report the mPQ score for both linear and fine-tuning experiment protocols for models trained using the TCGA data source. We compare results column-wise and mark the best results in **bold** and the second-best results in underline.

	BACH	CRC	PCam	MHIST	CoNSeP
BT trained on TCGA	84.2	94.2	84.5	78.0	40.9
+ our aug. techniques	<b>87.5</b>	<b>94.7</b>	<b>87.6</b>	<b>79.5</b>	<b>41.3</b>

Table F.1. **Benefit of our augmentation techniques.** Linear evaluation results show that our proposed augmentation techniques consistently and significantly improve performance.

## G. Intriguing Properties of Self-supervised ViT (Section 5.2)

As part of an effort to explore the potential of domain-aligned pre-training, we visualize the attention maps of self-supervised ViT and supervised ViT pre-trained on ImageNet. Our results, as illustrated in Fig. G.1, demonstrate that SSL ViT interestingly identifies and locates cells while also recognizing morphological phenotypes, which is aligned with recent observations [12]. Specifically, attention heads 1 ~ 4 attend to epithelial and inflammatory cells, whereas heads 5 ~ 6 focus on fibroblast cells. In contrast, supervised ViT pre-trained on ImageNet fails to generate interpretable signals due to the domain gap, highlighting the effectiveness of domain-aligned pre-training in generating informative signals for downstream tasks. We believe that this intriguing property can be leveraged to enable future potentials in the field of histopathology.

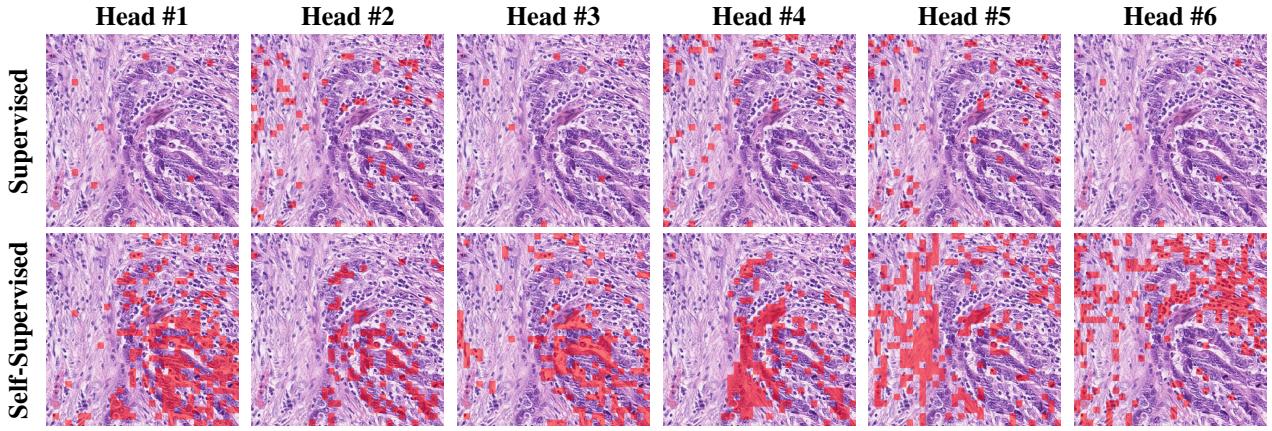
## H. Qualitative Results of Nuclei Instance Segmentation (Section 5.2)

In order to perform a qualitative assessment of the effect of domain-aligned pre-training on nuclei instance segmentation, we compare the predictions of models using supervised ImageNet pre-training and self-supervised TCGA pre-training, adapted under the linear evaluation protocol. The result presented in Fig. H.1 shows that domain-aligned

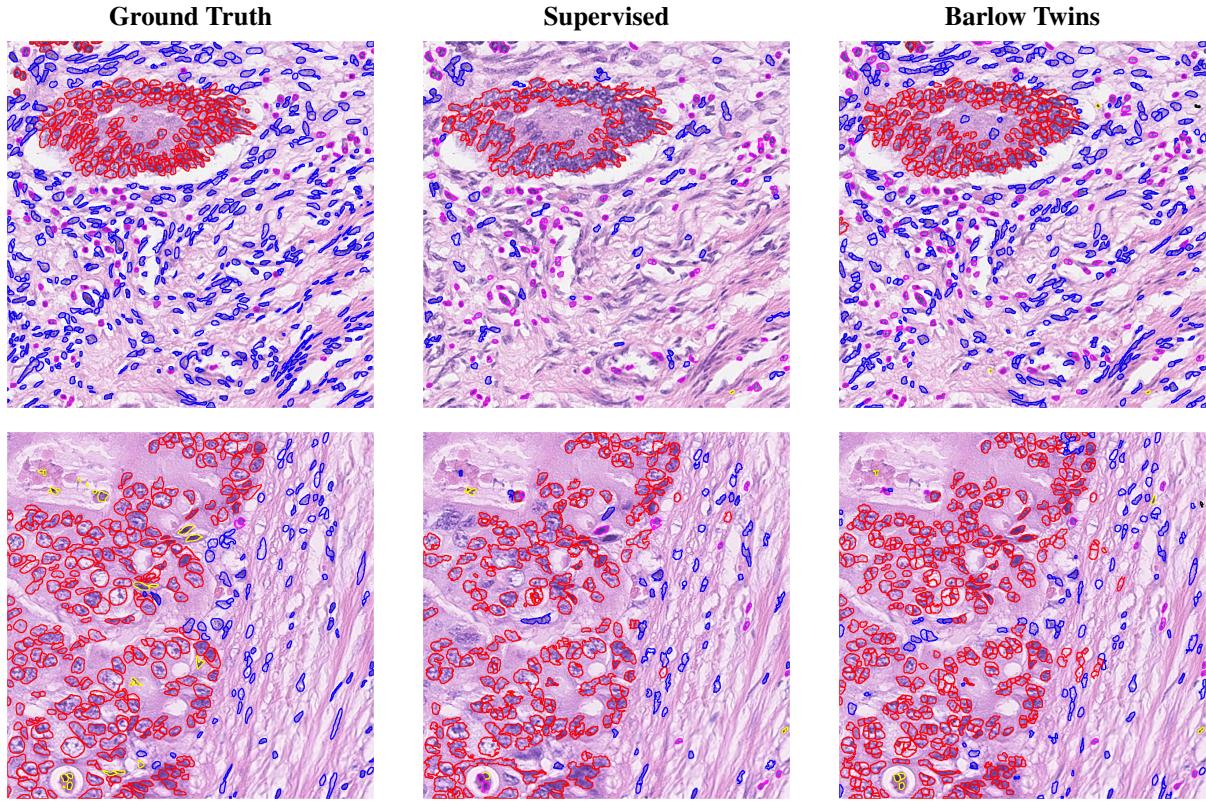
pre-training can offer the benefit on downstream tasks effectively, resulting in capturing foreground cells and accurately classifying them, in contrast to the model trained using ImageNet pre-trained weight.

## I. Slide-level Evaluation

The slide-level classification task is outside of the scope of our work. Nonetheless, we conduct a preliminary experiment to demonstrate the usefulness of the features learned through SSL for this task, too. We train and test models for the classification of breast cancer metastases in WSIs, using the same configuration as CLAM [45] but on the Camelyon16 [4] dataset. To extract features from the WSIs, we use two pre-trained weights: “Supervised (IN)” and “MoCo v2 (TC+TU)”. We find that models achieve an AUROC of 0.986 when using “MoCo v2 (TC+TU)” pre-trained weights, while models achieve an AUROC of 0.927 when using “Supervised (IN)” pre-trained weights. This result indicates that domain-aligned pre-training also can be beneficial to the slide-level task.



**Figure G.1. Visualizing multi-head self-attentions of ViT.** We visualize the attention map of several pre-trained ViT-S. Specifically, ViT-S has 6 attention heads. We visualize each head from the last layer of ViT. Our visualizations are presented in rows, with each row displaying attention maps alongside their corresponding overlayed image. The first two rows showcase the qualitative result of the supervised ViT pre-trained on ImageNet, while the next two rows display the qualitative result of the self-supervised ViT ( $\text{DINO}_{p=16}$ ) pre-trained on TCGA. Note that, the input image is resized to  $480 \times 480$  resolution, and overlaid in "red" are visual tokens whose attention weight  $> 0.5$  and span  $16 \times 16$  pixels.



**Figure H.1. Visualizing predictions of models.** We visualize the overlay predictions of different models on CoNSeP. A linear evaluation protocol is adopted to more directly assess the quality of the representations learned during pre-training. We selected the best-performing pre-trained model, Barlow Twins, based on the results of Table 4., obtained from the linear evaluation protocol. We find that predictions from Barlow Twins are similar to the ground-truth, whereas the “Supervised” alternative produces poor nuclei boundaries and merges cells incorrectly.