



# Vision Transformers for Small Histological Datasets Learned through Knowledge Distillation

Neel Kanwal<sup>1\*</sup>, Trygve Eftestøl<sup>1</sup>, Farbod Khoraminia<sup>2</sup>, Tahlita CM Zuiverloon<sup>2</sup>, and Kjersti Engan<sup>1</sup>

<sup>1</sup> Department of Electrical Engineering and Computer Science, University of Stavanger, Norway

<sup>2</sup> Department of Urology, University Medical Center Rotterdam, Erasmus MC Cancer Institute, Rotterdam, The Netherlands

\*Corresponding author: neel.kanwal@uis.no

**Abstract.** Computational Pathology (CPATH) systems have the potential to automate diagnostic tasks. However, the artifacts on the digitized histological glass slides, known as Whole Slide Images (WSIs), may hamper the overall performance of CPATH systems. Deep Learning (DL) models such as Vision Transformers (ViTs) may detect and exclude artifacts before running the diagnostic algorithm. A simple way to develop robust and generalized ViTs is to train them on massive datasets. Unfortunately, acquiring large medical datasets is expensive and inconvenient, prompting the need for a generalized artifact detection method for WSIs. In this paper, we present a student-teacher recipe to improve the classification performance of ViT for the air bubbles detection task. ViT, trained under the student-teacher framework, boosts its performance by distilling existing knowledge from the high-capacity teacher model. Our best-performing ViT yields 0.961 and 0.911 F1-score and MCC, respectively, observing a 7% gain in MCC against stand-alone training. The proposed method presents a new perspective of leveraging knowledge distillation over transfer learning to encourage the use of customized transformers for efficient preprocessing pipelines in the CPATH systems.

**Keywords:** Artifact Detection · Computational Pathology · Deep Learning · Knowledge Distillation · Vision Transformer · Whole Slide Images

## 1 Introduction

Histological examination of tissue samples is conducted by studying thin slices from a tumor specimen mounted on a glass slide. During the laboratory procedures, the preparation of glass slides may introduce artifacts and variations causing loss of visual [15,27]. Artifacts, such as air bubbles, occur when air is trapped under the cover slip due to improper mounting procedure [16]. Eventually, the presence of air bubbles leaves an altered and faded appearance [16,27]. During the manual assessment, pathologists usually ignore regions containing artifacts as they are irrelevant for diagnosis.

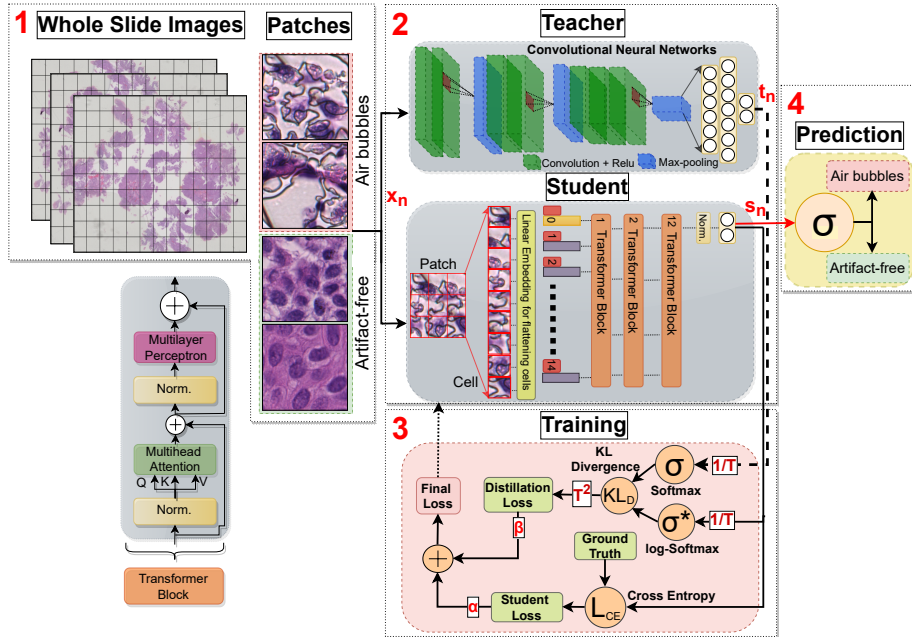
Computational Pathology (CPATH) systems are automated systems working with a digitized glass slide, called Whole Slide Image (WSI), as input. CPATH systems have the potential to automate diagnostic tasks and provide a second opinion or localize the Regions of Interest (ROIs) [14]. Different types of artifacts,

like air bubbles, might be present on the WSI [16] and can deteriorate diagnostic CPATH results if included in the analysis. Therefore it has been proposed to detect and exclude artifacts as a first step before using more relevant tissue in a diagnostic or prognostic system [15,16]. The detection and exclusion of artifacts can be regarded as (a part of) a *preprocessing pipeline*, which also might include color normalization and patching [16]. A complete preprocessing pipeline should detect folded tissue, damaged tissue, blood, and blurred (out of focus) areas, as well as air bubbles [16]. This might be done by an ensemble of models, one for each artifact, or by a multiclass model. In this paper, we consider detecting *air bubbles* artifact, which is not given much attention in the literature.

Deep Learning (DL) methods have shown promising results in various medical image analysis tasks [4,28], and can be used for detecting artifacts in a preprocessing pipeline. Supervised learning for generalized DL models requires a significant amount of data and labels. In CPATH literature, little effort has been made to annotate artifacts; thus, publicly available datasets for histological artifacts are unavailable. Transfer Learning (TL) has been widely used for medical images to deal with the lack of labeled training data [6,21]. TL methods use the existing knowledge, such as ImageNet [2] weights, and fine-tune the model for a different task. Although TL on ImageNet weights is useful to cope with a lack of data, ImageNet weights are mostly available for complicated Deep Convolutional Neural Networks (DCNN) architectures and carry a strong texture bias [5]. However, such DCNNs are typically computationally complex, whereas a preprocessing pipeline, being a first step prior to diagnostic or prognostic models, should have generalized and efficient DL models with high throughput. This is especially true with an ensemble of DCNN models for the different artifacts.

After the success in natural language processing tasks, *transformers* have been given attention for vision tasks [3,17]. Vision Transformers (ViTs), using a convolution-free approach, have surpassed DCNNs in accuracy and efficiency on image classification benchmarks [1,3]. Unlike the convolution layer in DCNNs, which applies the same filter weights to all inputs, the multi-head attention [30] in ViTs attends to image-wide structural information [20]. Interestingly, ViTs are also shown to be more robust and generalized than DCNNs [1,20]; Unfortunately, the robustness and generalizability come from training on extremely large datasets [1,3,29], which contrasts with the biomedical scenario. These limitations bring us to the question: *how can we train generalized ViTs on a small histopathological dataset?*

One possible answer lies in Knowledge Distillation (KD) [10], which transfers knowledge from a usually large teacher model to another, typically smaller, student model. Motivated by the KD idea, we present a student-teacher recipe, as shown in Fig 1. We propose to use KD in combination with TL for detecting air bubbles on WSIs using a small training set. In short, we let the teacher model be a complex ImageNet pretrained DCNN, and using KD, we train a small student model, which is a ViT. In the inference stage, we only need the small ViT, which is computationally efficient enough for a preprocessing pipeline implementation.



**Fig. 1. An overview of our proposed air bubbles detection method by knowledge distillation:** Predefined size patches for air bubbles and artifact-free classes are extracted from the WSI. A ViT student model is trained with the help of a DCNN teacher model by leveraging the transference of knowledge during the training process. The student-teacher recipe weights the teacher and student’s outputs by the temperature ( $T$ ). The overall training objective is to minimize the final loss, which is a linear combination of student loss and distillation loss. Finally, the student model is used to perform predictions for binary air bubbles detection task.

Our contributions in this paper can be summarized as follows:

- We train several state-of-the-art DCNNs and ViTs to compare their performance on a binary air bubbles detection task. Later, we choose suitable architectures to test our student-teacher framework.
- We conduct an in-depth comparison by initializing models with and without ImageNet weights and training ViT under a standalone vs. a student-teacher framework. We also assess the improvements in ViT’s generalization capability over ImageNet transfer learning.
- We run extensive experiments to test the student ViT’s performance under different teacher models and distillation configurations on unseen data.

## 2 Related Work

**Artifact and air bubbles detection:** The detection of histopathological artifacts has largely been overlooked during the development of CPATH systems, and the literature on air bubbles is scarce. Shakhawat et al. [11], in their quality evaluation method, detected air bubbles in two steps. First, the non-overlapping

affected patches were detected using a Support Vector Machine (SVM) classifier. Later, the remaining patches with faded appearance were separated using handcrafted Gray-level Co-occurrence Matrix (GLCM) features. This work was later extended in [24], where a pretrained VGG16 [25] network was used to compare the handcrafted features against the CNN-based method. Their experiments concluded that handcrafted features provide stable classification, but their evaluation was based on a relatively smaller dataset. Recently, Raipuria et al. [22] performed stress testing for common histological artifacts, including air bubbles, using a vision transformer [29] and a ResNet [9] model. Though, MobileNet [12] and VGG16 [25] have been popular DCNN choices for artifact detection [15]. DCNNs are found to be less robust than ViTs and exhibit strong texture bias [20,22].

**Knowledge Distillation (KD):** Originally proposed by Hinton et al. [10] for model compression, KD sought to extract knowledge from an ensemble of CNN experts to a smaller two-layer CNN generalist network to make it perform equally well. In short, KD aims to train a small student model under the guidance of a complicated teacher model, where the student model optimizes its learning by absorbing the hidden knowledge from the teacher. This transference of knowledge can be accomplished by minimizing output logits of student and teacher networks through some distillation methods, such as logit-based, feature-based, and relationship-based distillation methods [19].

KD helps make computationally friendly deployment algorithms, making it interesting for many biomedical imaging algorithms. Lingmei et al. [18] proposed a CNN model for glioma classification. They used the KD approach to compress the model and make it suitable for deployment on medical equipment. Salehi et al. [23] used a VGG16 [25] cloner network to calculate multi-level loss from a source network for detecting anomalies. Their method relied on distilling intermediate knowledge from the ImageNet pretrained source network. In a similar approach, He et al. [8] used the KD technique to boost the performance of CNN for ocular disease classification. They used fundus images and clinical information to train a ResNet [9] teacher first and used only the fundus images to train a similar student network later. Guan et al. [7] detected Alzheimer’s disease by leveraging multi-modal data to train a teacher network. Their distillation scheme improved the prediction performance of the ResNet [9] student using a single imaging modality.

However, all these works focused on using only CNN as a student network and did not explore the effects of different configurations and teacher networks on the final classification outcome. In addition, the use of KD for histological artifacts has not been investigated yet.

### 3 Data Materials and Method

Fig. 1 provides an overview of our air bubbles detection method using KD [10] in a student-teacher recipe. We exploit KD for data-efficient training by leveraging the transference of knowledge from the teacher model to the student model. Our proposed method uses a complex DCNN as the pre-trained teacher and a small ViT as the student when a small histological dataset is available. We are

doing a logit-based distillation [19] since our teacher and student models are very different. The steps of our method are further described below.

### 3.1 Dataset

The air bubbles dataset was prepared from 55 bladder biopsy WSIs, provided by Erasmus Medical Center (EMC), Rotterdam, The Netherlands. The glass slides were stained with Hematoxylin and Eosin (H&E) dyes and scanned with Hamamatsu Nanozoomer at 40× magnification. WSIs are stored in *ndpi* format with a pixel size of 0.227 μm × 0.227 μm. These WSIs were manually annotated for air bubbles and artifact-free tissue by a non-pathologist who has received training for the task. To prevent data leakage, the dataset was later split into 35/10/10 training, validation, and test WSIs, respectively.

### 3.2 Foreground Segmentation and Patching

Let  $I_{\text{WSI}(i)}^{40x}$  correspond to a WSI at magnification level 40x (sometimes referred to as 400x).  $I_{\text{WSI}}^{40x}$  are very large gigapixel images, and it is not feasible to process the entire WSI at once. As such, all CPATH systems resort to patching or tiling of the image, or the ROI in the image, before further processing. Let  $\mathcal{T} : I_{\text{WSI}(i) \in R}^{40x} \rightarrow \{\mathbf{x}_j^i; j = 1 \dots J\}$  represent the process of patching a ROI denoted by  $R$  of the image  $I_{\text{WSI}(i)}^{40x}$  into a set of  $J$  patches, where  $\mathbf{x}_j^i \in \mathbb{R}^{W \times H \times C}$  and  $W, H, C$  present the width, height, and channels of the image, respectively. In the patching process, foreground-background segmentation was performed first by transforming (Red, Green, Blue) RGB images to (Hue, Saturation, Value) HSV color space. Later, Otsu thresholding was applied to the value channel to obtain the foreground with tissue. The extracted foreground was later divided over a non-overlapping square grid, and patches with at least 70% overlap to the annotation region ( $R$ ) were extracted.

Let  $\mathcal{D} = (\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$  denote our prepared dataset of  $N$  patches from a set of WSIs and  $y_n \in \{0, 1\}$  is the binary ground truth for the  $n$ -th instance, where 1 indicates a patch within a region marked as air bubbles. Fig 1 (step 1) shows the patches  $\mathbf{x}_n$  of  $224 \times 224 \times 3$  pixels with air bubbles and artifact-free classes obtained from a WSI at 40x magnification.

### 3.3 Selecting Student-Teacher Architectures

Let’s symbolize the student model  $\xi$  with parameters  $\theta$  providing the prediction output logits  $\mathbf{s}_n = \xi_{\theta}(\mathbf{x}_n)$ , and correspondingly, the teacher model  $\varphi$  parameterized by  $\phi$  providing the output logits  $\mathbf{t}_n = \varphi_{\phi}(\mathbf{x}_n)$ .

Our student model is a ViT, similar to the pioneering work [3], which leverages multi-head self-attention mechanism [30] to capture content-dependant relations across the input patch. At the image pre-processing layer, the patches of  $224 \times 224$  pixels are split into the non-overlapping cells of  $16 \times 16$  pixels. Later, the linear embedding layer flattens these cells, and positional encodings are added before feeding the embeddings to the pile of transformer blocks, as illustrated in Fig. 1 (step 2). Since convolutional networks have shown their efficacy in image recognition tasks, transferring knowledge from a DCNN network

can help the ViT absorb inductive biases. Therefore, we rely on popular state-of-the-art DCNNs for selecting teacher architecture. Nevertheless, we systemically discover appropriate student and teacher candidates during the experiments later to demonstrate the approach’s effectiveness over TL.

### 3.4 Training Student under Knowledge Distillation

After selecting student and teacher architectures, we begin the process of training the student  $\xi$ . The goal is to train  $\xi$  with the assistance of a  $\varphi$  to improve the  $\xi$ ’s generalization performance using additional knowledge beyond the labels. Our approach is similar to Hinton et al. [10] where model outputs  $\mathbf{s}$ , and  $\mathbf{t}$  are normalized by a temperature  $T$  parameter before using the softmax function  $\sigma$ . The increasing value of  $T$  softens the impact of the fluctuations in the output probability distribution; therefore, more knowledge can be devolved with each input  $\mathbf{x}_n$ . Instead of using softmax on  $\mathbf{s}_n$ , we take advantage of the log-softmax function  $\sigma^*$ , which stabilizes the distillation process by penalizing for incorrect class.  $\sigma^*$  also adds efficiency by optimizing gradient calculations.

The output logits for input patch  $\mathbf{x}_n$  can be written as;

$$\mathbf{s}_n = \xi_\theta(\mathbf{x}_n) \quad \text{and} \quad \mathbf{t}_n = \varphi_\phi(\mathbf{x}_n) \quad (1)$$

Let the log-softmax and softmax on logits,  $\sigma^*(\mathbf{s}/T)$  and  $\sigma(\mathbf{t}/T)$ , for each element can be defined as (see Eq. (2));

$$\sigma^*(s_i/T) = \log\left(\frac{\exp(s_i/T)}{\sum_{j=1}^c \exp(s_j/T)}\right) \quad \text{and} \quad \sigma(t_i/T) = \frac{\exp(t_i/T)}{\sum_{j=1}^c \exp(t_j/T)} \quad (2)$$

where  $c$  is the total number of classes and  $T$  is the temperature. The class probabilities at the output of the  $\xi$  and  $\varphi$  model can thus be written as;

$$p_\xi = \sigma^*(\mathbf{s}/T) = \sigma^*(\xi_\theta(\mathbf{x})) \quad \text{and} \quad p_\varphi = \sigma(\mathbf{t}/T) = \sigma(\varphi_\phi(\mathbf{x})) \quad (3)$$

The student loss  $L_{student}$  (Eq. (4)) provides hard targets and is obtained by applying cross entropy  $L_{CE}$  on ground truth  $y$ , and  $\mathbf{s}$  when  $T$  is set to 1;

$$L_{student} = L_{CE}(y, \mathbf{s}) = - \sum_{i=1}^c y_i \cdot \log(\sigma^*(s_i)) \quad (4)$$

Distillation loss  $L_{distillation}$  provides the soft targets and is computed from the  $p_\xi$  and  $p_\varphi$  by applying Kullback-Leibler divergence  $KL_D$ . Since the outputs from  $\xi$  and  $\varphi$  were normalized by  $T$ , we multiply the loss with  $T^2$  to maintain their relative contribution;

$$L_{distillation} = T^2 \times KL_D(p_\xi || p_\varphi) = T^2 \cdot \sum_{i=1}^c p_{\xi_i} \log \frac{p_{\xi_i}}{p_{\varphi_i}} \quad (5)$$

The final loss function, as shown in Eq. (6), is a weighted average of student and distillation losses where  $\alpha \in [0, 1]$ ;

$$L_{Final} = \alpha \times L_{student} + \beta \times L_{distillation} \quad \cdot \cdot \cdot \beta = 1 - \alpha \quad (6)$$

High entropy in soft targets offers significantly more information per training patch than hard targets [10], allowing the student ViT to train with fewer data and a higher learning rate. Therefore, using a smaller alpha can be beneficial if the  $\xi$  is trained from scratch. Our standalone training setup for baseline comparison can be obtained by putting  $\alpha$  and  $T$  equal to one and replacing log softmax with softmax function.

### 3.5 Prediction

Once the final loss is minimized based on the experimental setup (defined in Sec. 4), we find predictions from the student  $\xi$  by setting  $T$  equal to one. For an unseen test patch  $\mathbf{x}_*$ , output can be defined as (7);

$$\hat{y}_s = \arg \max(\sigma(\mathbf{s}_*)) = \arg \max(\sigma(\xi_\theta(\mathbf{x}_*))) \in \{0, 1\} \tag{7}$$

## 4 Experimental Setup

**Implementation Details:** The patch extraction was accomplished using the HistoLab library. Extracted patches were normalized to ImageNet [2] mean and standard deviation. We augmented data at every training epoch using random geometric transformations, such as rotations, horizontal and vertical flips. ViTs were borrowed from Timm Library, and the experimental setup was built on the Pytorch. We used four variants of ViTs with different parametric depths from [3,29], where the classifier was replaced by a fully connected (FC) layer. We used four state-of-the-art DCNNs with varying parametric complexity. All DCNN backbones were initialized with ImageNet [2] weights, and classifiers were replaced with three-layer FC classifiers. All classifiers were initialized with random weights. After hyper-parameter exploration, the final parameters were set to a batch size of 64, SGD optimizer, ReduceLROnPlateau scheduler with a learning rate of 0.001, dropout of 0.2, cross-entropy loss, and early stopping with the patience of 20 epochs on validation loss to prevent over-fitting. For KD parameters, values of  $T \in \{2, 5, 10, 20, 40\}$  and  $\alpha \in \{0.3, 0.5, 0.7\}$  were explored. The best model weights are used to report the results. The NVIDIA GeForce A100 SXM 40GB GPU was utilized for training all models.

**Evaluation Metrics:** We evaluate the presented method using accuracy, F1-score, and Mathew Correlation Coefficient (MCC). Let TP, FN, FP, and TN describe true positive, false negative, false positive, and false negative predictions. The accuracy, termed as  $(TP + TN)/(TP + FN + FP + TN)$ , is the ratio of correct predictions by the model. F1 is the harmonic mean, defined as  $2 \cdot (\text{precision} \cdot \text{recall})/(\text{precision} + \text{recall})$  where  $\text{Recall} = TP/(TP + FN)$  and  $\text{Precision} = TP/(TP + FP)$ . MCC is an informative measure in binary classification over imbalanced datasets and is defined as Eq. (8).

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP) \cdot (TP + FN) \cdot (TN + FP) \cdot (TN + FN)}} \in [-1, 1] \tag{8}$$

**Table 1. Results from Exp. 1:** Four variants of Deep Convolutional Neural Networks (DCNNs) and Vision Transformers (ViTs), with increasing parametric complexity, are trained for the air bubbles detection task. The best outcomes in every section are bolded. ViT-tiny and MobileNet architectures provide the best results on the test set.

Architecture	Param. (#)	Validation Set			Test Set		
		Acc.(%)	F1	MCC(↑)	Acc.(%)	F1	MCC(↑)
<b>Deep Convolutional Neural Networks (DCNNs)</b>							
MobileNetv3 [12]	3.52M	98.28	0.983	0.965	<b>93.88</b>	<b>0.945</b>	<b>0.876</b>
EfficientNet [26]	20.89M	96.52	0.966	0.931	92.54	0.935	0.851
DenseNet161 [13]	27.66M	98.12	0.982	0.962	91.32	0.925	0.828
VGG16 [25]	136.42M	<b>98.34</b>	<b>0.984</b>	<b>0.966</b>	92.31	0.932	0.846
<b>Vision Transformers (ViTs)</b>							
ViT-tiny [29]	5.52M	<b>98.67</b>	<b>0.987</b>	<b>0.973</b>	<b>92.35</b>	<b>0.933</b>	<b>0.847</b>
ViT-small [29]	21.66M	97.01	0.971	0.941	91.16	0.922	0.822
ViT-large [3]	303.30M	98.12	0.982	0.962	92.08	0.928	0.839
ViT-huge [3]	630.76M	95.85	0.962	0.918	91.43	0.925	0.829
<b>Results from Literature (Validation Accuracy (%))</b>							
DeiT-S in [22]	91.5-92	ResNet-50 in [22]	88-89	VGG16 in [24]	87.33		

## 5 Results and Discussion

### 5.1 Exp. 1: Baseline Experiments for Architecture Decision

In this experiment, we only apply TL to a set of architectures. We evaluate state-of-the-art DCNNs, namely MobileNetv3 [12], EfficientNet [26], DenseNet161 [13] and VGG16 [25] architectures and a family of four ViTs [3,29], with increasing architecture size. Exp 1 provides a baseline as well as helps to choose architectures for the KD setup in later experiments. Table 1 reports the results of the validation and test set. DCNNs largely exceed the performance of ViTs, where top-performing ViT lags the generalization performance of top-performing DCNNs by 3% in MCC. Moreover, architectures with sizeable parameters like VGG16 and ViT-tiny and MobileNet, despite being architectures with fewer parameters, emerge as appropriate student and teacher candidates, respectively, based on the test results and outperform the results from the literature.

### 5.2 Exp. 2: How Important is Teacher’s Knowledge?

This experiment evaluates the impact of existing teacher knowledge in the KD process to assess the real-life analogy where good teachers make good students. Therefore, we initialize MobileNet teachers with no knowledge (scratch), knowledge from a general domain (ImageNet), knowledge from another WSI artifact (damaged tissue [15]), and finally, domain-relevant knowledge (air bubbles) from the previous experiment. In addition, we also select VGG16 with air bubble knowledge as a teacher to assess the effect of highly parametric DCNN in the



**Table 2. Results from Exp. 2:** Knowledge Distillation (KD) outcome for selected teacher and student candidates from Exp.1. The values of  $\alpha, T$  are fixed at 0.5 and 10, respectively. The best results in every part are marked in bold, and the second best is underlined. ViT-tiny, with two scratch and ImageNet initialization, is used for baseline comparisons. Two teachers (MobileNet and VGG16) with air bubbles knowledge are used. While MobileNet is also initialized with knowledge of other domains to evaluate the importance of teachers’ knowledge.

Architecture (Initial.)	Validation Set			Test Set		
	Acc.(%)	F1	MCC( $\uparrow$ )	Acc.(%)	F1	MCC( $\uparrow$ )
<b>Baseline (Initial.) - Standalone training</b>						
ViT-tiny (Scratch)	96.13	0.963	0.922	91.51	0.925	0.829
ViT-tiny (ImageNet [2])	<b>98.67</b>	<b>0.987</b>	<b>0.973</b>	<b>92.35</b>	<b>0.933</b>	<b>0.847</b>
<b>Teacher (Initial.) - Student [ViT-tiny (Scratch)]</b>						
MobileNet (Scratch)	96.13	0.962	0.924	87.92	0.889	0.756
MobileNet (ImageNet [2])	95.58	0.957	0.914	92.31	0.927	0.848
MobileNet (Damaged [15])	76.8	0.785	0.533	49.23	0.608	-0.075
MobileNet (Air bubbles)	<b>98.01</b>	<b>0.981</b>	<b>0.960</b>	<b>95.25</b>	<b>0.957</b>	<b>0.904</b>
VGG16 (Air bubbles)	<u>97.18</u>	<u>0.973</u>	<u>0.944</u>	<u>93.42</u>	<u>0.940</u>	<u>0.867</u>
<b>Teacher (Initial.) - Student [ViT-tiny (ImageNet)]</b>						
MobileNet (Scratch)	<b>98.73</b>	0.983	0.971	93.38	0.941	0.866
MobileNet (ImageNet [2])	98.62	<b>0.987</b>	<u>0.972</u>	93.40	0.942	0.867
MobileNet (Damaged [15])	50.08	0.211	0.09	35.51	0.116	-0.294
MobileNet (Air bubbles)	98.61	<b>0.987</b>	<b>0.973</b>	<b>95.60</b>	<b>0.961</b>	<b>0.911</b>
VGG16 (Air bubbles)	<u>98.67</u>	<u>0.986</u>	<u>0.972</u>	<u>94.19</u>	<u>0.948</u>	<u>0.882</u>

KD process. For this experiment, the values of  $\alpha, T$  are fixed at 0.5 and 10, respectively. The student is a ViT-tiny architecture initialized with random and ImageNet weights separately.

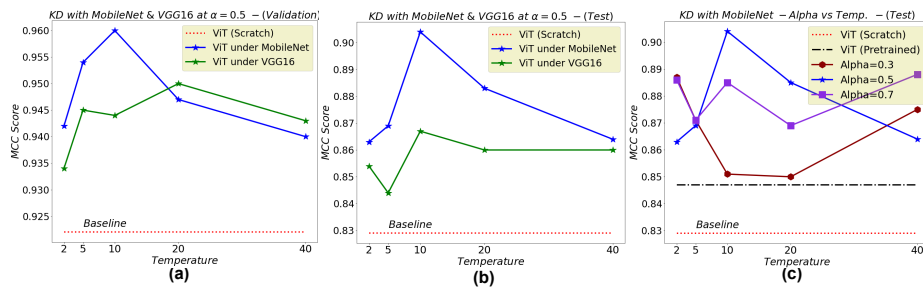
Table 2 exhibits that KD remarkably improves ViT’s classification ability. Even without ImageNet knowledge, ViT-tiny, under the KD framework, surpasses all metrics under both MobileNet and VGG16 teachers. However, the best results are obtained using the MobileNet teacher, ascertaining that hidden knowledge can be easily distilled from a simpler architecture. Interestingly, teachers with knowledge other than the relevant domain (air bubbles) produce poorly performing student. Although the student with ImageNet knowledge does not indicate gain on the validation results relative to the baseline, it achieves 3% and 7% improvement in F1 and MCC scores on the test set, respectively.

Overall, the test results demonstrate that the KD is promising to train generalized ViT-tiny with little data, even without pretrained weights. ViT significantly enhances its generalization against the baseline when trained in a standalone setting. Especially when the teacher is enriched with the knowledge related to the task. KD, on top of ImageNet TL, provides a marginal gain in the performance of ViT-tiny, overcoming the reliance on pretrained weights.

### 5.3 Exp. 3: Influence of KD Parameters

Since the initialization of teachers with air bubbles knowledge has been shown to improve the learning process, it would be interesting to assess the influence of DCNN teachers under the different KD parameters ( $T$  and  $\alpha$ ). In this experiment, we chose  $T \in \{2, 5, 10, 20, 40\}$  and  $\alpha \in \{0.3, 0.5, 0.7\}$  to estimate the influence of teacher’s output on ViT student, trained from scratch. The baseline experiment corresponds to  $\alpha$  and  $T = 1$  and uses sigmoid on ViT outputs. Fig. 2 (a) and (b) show MCC values as the effect of temperature on simple DCNN like MobileNet and complex DCNN like VGG16. Though the ViT-tiny student trained under the VGG16 teacher scores better on the validation set when  $T$  is high, the MobileNet teacher reveals better transference of hidden knowledge on all  $T$  values on the test set. Fig. 2 (c) depicts the effect of  $\alpha$  on ViT’s generalization results. All  $\alpha$  values give better results than the baseline, concluding that including distillation loss improves training compared to only student loss.

To sum up, the teacher’s outcome strongly influences the student’s generalizability in the KD process. Most of the  $T$  and  $\alpha$  values deliver a noticeable gain over the standalone training in our case. However, *intermediate*  $T$  values and assigning *equal weight* to student and distillation loss is the most advantageous.



**Fig. 2. Results from Exp. 3:** Knowledge Distillation (KD) improves the performance of the Vision Transformer (ViT-tiny) under the supervision of both MobileNet and VGG16 teachers. (a) and (b) shows an improved performance from the baseline (standalone training from scratch), under all temperature ( $T$ ) values, on validation and test set. (c) depicts the influence of giving higher/lower weightage to distillation loss from the teacher network (see Sec. 3). The MobileNet teacher, despite being simpler architecture, enriches ViT-tiny’s generalization capability on all chosen  $\alpha$  and  $T$  values.

## 6 Conclusion and Future Work

This paper presents the Knowledge Distillation (KD) to boost the generalization performance of small Vision Transformers (ViTs) on a small histopathological dataset. The main motivation is to create a well-performing and efficient preprocessing pipeline that requires a generalized and computationally-friendly model. We evaluated various pretrained DCNNs and ViTs for the air bubbles artifact detection task. ViTs, trained in a standalone setting, underperform DCNNs on unseen data. Our approach exploits the KD, in the absence of pretrained weights, to enhance the performance of ViT by training under the guidance of a DCNN

teacher. Our analysis found that KD provides significant gain under most distillation settings when the teacher holds the knowledge of the same task. In conclusion, the ViT, when trained under KD, outperforms its state-of-the-art DCNN teacher and its counterpart in standalone training.

In future work, the method can be developed and tested on larger cohorts of histological data with stain variations and to detect multiple artifacts. Moreover, artifact detection by ViT trained under the student-teacher recipe can be combined as a preprocessing step with a diagnostic or prognostic algorithm in the computational pathology system.

## 7 Acknowledgment

This research is supported by the European Horizon 2020 program under Marie Skłodowska-Curie grant agreement No. 860627 (CLARIFY). The authors have no relevant financial or non-financial interests to disclose.

## References

1. Bhojanapalli, S., Chakrabarti, A., Glasner, D., Li, D., Unterthiner, T., Veit, A.: Understanding robustness of transformers for image classification. In: Proc. of the IEEE International Conf. on Computer Vision (ICCV). pp. 10231–10241 (2021)
2. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE ICCV. pp. 248–255. Ieee (2009)
3. Dosovitskiy, A., Beyer, L., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)
4. Fuster, S., Khoraminia, F., et al.: Invasive cancerous area detection in non-muscle invasive bladder cancer whole slide images. In: IEEE 14th Image, Video, and Multidimensional Signal Processing Workshop (IVMSP). pp. 1–5. IEEE (2022)
5. Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F.A., Brendel, W.: Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. arXiv preprint arXiv:1811.12231 (2018)
6. Golatkar, A., Anand, D., et al.: Classification of breast cancer histology using deep learning. In: International conf. image analysis and recognition. pp. 837–844. Springer (2018)
7. Guan, H., Wang, C., Tao, D.: Mri-based alzheimer’s disease prediction via distilling the knowledge in multi-modal data. *NeuroImage* **244**, 118586 (2021)
8. He, J., Li, C., Ye, J., Qiao, Y., Gu, L.: Self-speculation of clinical features based on knowledge distillation for accurate ocular disease classification. *Biomedical Signal Processing and Control* **67**, 102491 (2021)
9. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE CVPR. pp. 770–778 (2016)
10. Hinton, G., Vinyals, O., Dean, J., et al.: Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531 **2**(7) (2015)
11. Hossain, M.S., Nakamura, T., Kimura, F., Yagi, Y., Yamaguchi, M.: Practical image quality evaluation for whole slide imaging scanner. In: Biomedical Imaging and Sensing Conference. vol. 10711, pp. 203–206. SPIE (2018)

12. Howard, A., Sandler, M., Chu, G., Chen, L.C., Chen, B., Tan, M., Wang, W., Zhu, Y., Pang, R., Vasudevan, V., et al.: Searching for mobilenetv3. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1314–1324 (2019)
13. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4700–4708 (2017)
14. Kanwal, N., Amundsen, R., Hardardottir, H., Janssen, E.A., Engan, K.: Detection and localization of melanoma skin cancer in histopathological whole slide images. arXiv preprint arXiv:2302.03014 (2023)
15. Kanwal, N., Fuster, S., et al.: Quantifying the effect of color processing on blood and damaged tissue detection in whole slide images. In: IEEE 14th Image, Video, and Multidimensional Signal Processing Workshop (IVMSP). pp. 1–5. IEEE (2022)
16. Kanwal, N., Pérez-Bueno, F., Schmidt, A., Engan, K., Molina, R.: The devil is in the details: Whole slide image acquisition and processing for artifacts detection, color variation, and data augmentation. *IEEE Access* **10**, 58821–58844 (2022)
17. Kanwal, N., Rizzo, G.: Attention-based clinical note summarization. In: Proceedings of the 37th ACM Symposium on Applied Computing. pp. 813–820 (2022)
18. Lingmei, A., et al.: Noninvasive grading of glioma by knowledge distillation base lightweight convolutional neural network. In: IEEE 2021 AEMCSE. pp. 1109–1112
19. Meng, H., Lin, Z.e.a.: Knowledge distillation in medical data mining: A survey. In: 5th International Conf. on Crowd Science and Engineering. pp. 175–182 (2021)
20. Naseer, M., Ranasinghe, K., Khan, S., Hayat, M., Shahbaz Khan, F., Yang, M.H.: Intriguing properties of vision transformers. *NeurIPS* **34**, 23296–23308 (2021)
21. Noorbakhsh, J., Farahmand, S., et al.: Deep learning-based cross-classifications reveal conserved spatial behaviors within tumor histological images. *Nature communications* **11**(1), 1–14 (2020)
22. Raipuria, G., Singhal, N.: Stress testing vision transformers using common histopathological artifacts. In: Medical Imaging with Deep Learning (2022)
23. Salehi, M., Sadjadi, N., Baselizadeh, S., Rohban, M.H., Rabiee, H.R.: Multiresolution knowledge distillation for anomaly detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 14902–14912 (2021)
24. Shakhawat, H.M., Nakamura, T., Kimura, F., Yagi, Y., Yamaguchi, M.: Automatic quality evaluation of whole slide images for the practical use of whole slide imaging scanner. *ITE Trans. on Media Technology and Applications* **8**(4), 252–268 (2020)
25. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)
26. Tan, M., Le, Q.: Efficientnet: Rethinking model scaling for convolutional neural networks. In: International conf. on machine learning. pp. 6105–6114. PMLR (2019)
27. Taqi, S.A., Sami, S.A., Sami, L.B., Zaki, S.A.: A review of artifacts in histopathology. *Journal of oral and maxillofacial pathology: JOMFP* **22**(2), 279 (2018)
28. Tomasetti, L., Khanmohammadi, M., Engan, K., Høllesli, L.J., Kurz, K.D.: Multi-input segmentation of damaged brain in acute ischemic stroke patients using slow fusion with skip connection. arXiv preprint arXiv:2203.10039 (2022)
29. Touvron, H., et al.: Training data-efficient image transformers & distillation through attention. In: Int. Conf. on Machine Learning. pp. 10347–10357 (2021)
30. Vaswani, A., Shazeer, N., et al: Attention is all you need. *Advances in neural information processing systems* (2017)