# ScoreNet: Learning Non-Uniform Attention and Augmentation for Transformer-Based Histopathological Image Classification

Thomas Stegmüller[1]  Behzad Bozorgtabar[1,2,3]  Antoine Spahr[1]  Jean-Philippe Thiran[1,2,3]

[1]EPFL, Switzerland  [2]CHUV, Switzerland  [3]CIBM, Switzerland
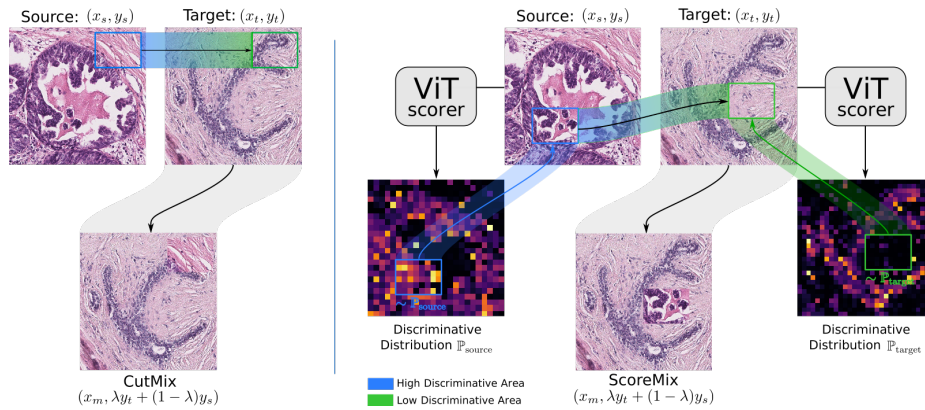
{firstname.lastname}@epfl.ch

**Figure 1:** CutMix (**left**) randomly mixes samples, yielding label misallocation, while our ScoreMix (**right**) creates a coherent artificial training pair $(x_m, y_m)$ by pasting a region of high semantic content from the source image, $x_s$ at a non-discriminative region of the target image $x_t$, and obtains a convex combination of the labels.

## Abstract

*Progress in digital pathology is hindered by high-resolution images and the prohibitive cost of exhaustive localized annotations. The commonly used paradigm to categorize pathology images is patch-based processing, which often incorporates multiple instance learning (MIL) to aggregate local patch-level representations yielding image-level prediction. Nonetheless, diagnostically relevant regions may only take a small fraction of the whole tissue, and current MIL-based approaches often process images uniformly, discarding the inter-patches interactions. To alleviate these issues, we propose ScoreNet, a new efficient transformer that exploits a differentiable recommendation stage to extract discriminative image regions and dedicate computational resources accordingly. The proposed transformer leverages the local and global attention of a few dynamically recommended high-resolution regions at an efficient computational cost. We further introduce a novel mixing data-augmentation, namely ScoreMix, by leveraging the image's semantic distribution to guide the data mixing and produce coherent sample-label pairs. ScoreMix is embarrassingly simple and mitigates the pitfalls of previous augmentations, which assume a uniform semantic distribution and risk mislabeling the samples. Thorough experiments and ablation studies on three breast cancer histology datasets of Haematoxylin & Eosin (H&E) have validated the superiority of our approach over prior arts, including transformer-based models on tumour regions-of-interest (TRoIs) classification. ScoreNet equipped with proposed ScoreMix augmentation demonstrates better generalization capabilities and achieves new state-of-the-art (SOTA) results with only 50% of the data compared to other mixing augmentation variants. Finally, ScoreNet yields high efficacy and outperforms SOTA efficient transformers, namely TransPath [43] and SwinTransformer [24], with throughput around $3\times$ and $4\times$ higher than the aforementioned architectures, respectively.*

## 1. Introduction

Due to the increasing availability of digital slide scanners enabling pathologists to capture high-resolution whole slide images (WSI), computational pathology is becoming a ripe ground for deep learning and recently witnessed a lot of advances. Nonetheless, the diagnosis from H&E stained WSIs remains challenging. The difficulty of the task is a consequence of two inherent properties of histopathology image datasets: *i)* the huge size for images and *ii)* the cost of exhaustive localized annotations, making the usage of most deep learning models computationally infeasible. Patch-based processing approaches [37, 27, 17] have become a *de facto* practice for high dimensional pathology images that aggregate individual patch representation/classification predictions by, e.g., a convolutional neural network (CNN) for image-level prediction. Nonetheless, patch-based methods increase the requirement of patch-level labeling and further regions of interest (RoI) detection as diagnostic-related tissue sections might only take a small fraction of the whole tissue, leading to considerable uninformative patches. Prior CNN methods [18, 22] have adopted multiple instance learning (MIL) [26] to address the above issues, which incorporates an attention-based aggregation operator to identify tissue sub-regions of high diagnostic value automatically. Nonetheless, these MIL methods embed all the patches independently and discard the inter-patches correlation or only incorporate it at a later stage.

Recently, self-supervised learning (SSL) methods [22, 21, 38, 9] aimed to construct semantically meaningful visual representations via pretext tasks for histopathological images. Despite their notable success using CNN backbones in improving classification performances, CNN's receptive field often restricts the learning of global context features. In another line of research, to compensate for the lack of diverse and large datasets, mixing augmentation techniques [42, 45, 46] have been developed to further enhance the performance of these models. While there have been substantial performance gains on natural image datasets, we argue that such data augmentations may not be helpful for histopathological images, as they risk creating locally ambiguous images or mislabelled samples. Furthermore, contrary to CNNs, vision transformer (ViT) models [13, 41] can capture long-range visual dependencies due to their flexible receptive fields via self-attention mechanisms. More recently, self-supervised ViTs method [43, 23] combined the advantages of ViT and SSL to efficiently learn visual representations from less curated pre-training data. Despite their usefulness, there is relatively little research on the impact of data augmentation design, efficiency and robustness of ViT for histopathological image classification. For example, can we train an efficient transformer by selecting only informative regions of high diagnostic value (RoIs) from high-resolution images? What data augmen-

tation strategies can improve the transformer's representation learning for TRoIs classification? This paper addresses these questions by uncovering insights about key aspects of data augmentation and exploits the self-attention maps to identify the most relevant regions for the end task and train an efficient transformer.

**Contributions.** Our contributions are as follows:

1. We propose ScoreNet, a new efficient transformer-based architecture for histopathological image classification. It combines a fine-grained local attention mechanism with a coarse-grained global attention module to extract cell- and tissue-level features. Benefiting from a differentiable recommendation module, the proposed architecture only processes the most discriminative regions of the high-resolution image, making it significantly more efficient than competitive transformer architectures without compromising accuracy;

2. A novel mixing data-augmentation, namely ScoreMix for histopathological images is presented. ScoreMix works in synergy with our architecture, as they build upon the same observation: the different regions of the images are not equally relevant for a given task. Using the learned self-attention w.r.t. the `[CLS]` token, we determine the distribution of the semantic regions in images during training to ensure sampling of informed cutting and pasting locations (see Fig. 1);

3. We empirically show consistent improvements of ScoreNet over SOTA methods for TRoIs classification on the BRACS dataset while we demonstrate ScoreNet's generalization capability on the CAMELYON16 and BACH datasets. The interpretability of ScoreNet behaviour is also investigated. Finally, we demonstrate ScoreNet throughput improvements over existing efficient transformers, making it an ideal candidate for applications on WSIs.

## 2. Related work

**TRoIs Classification.** Conventionally, deep convolutional neural networks [37, 36, 27, 17, 44] process pathology images in a patch-wise manner using a MIL formulation [26] and aggregate patch-level features extracted by CNNs. Nonetheless, current MIL methods discard the inter-patches interaction or only integrate it at the very end of the pipeline. Similarly, the computational resources dedicated to a specific region are independent of its pertinence for the task. Current methods rely on attention-based MIL techniques [18, 22, 19, 6, 33] to account for the non-uniform relevance of patches. On the contrary, the integration of contextual cues remains almost untouched, as all the
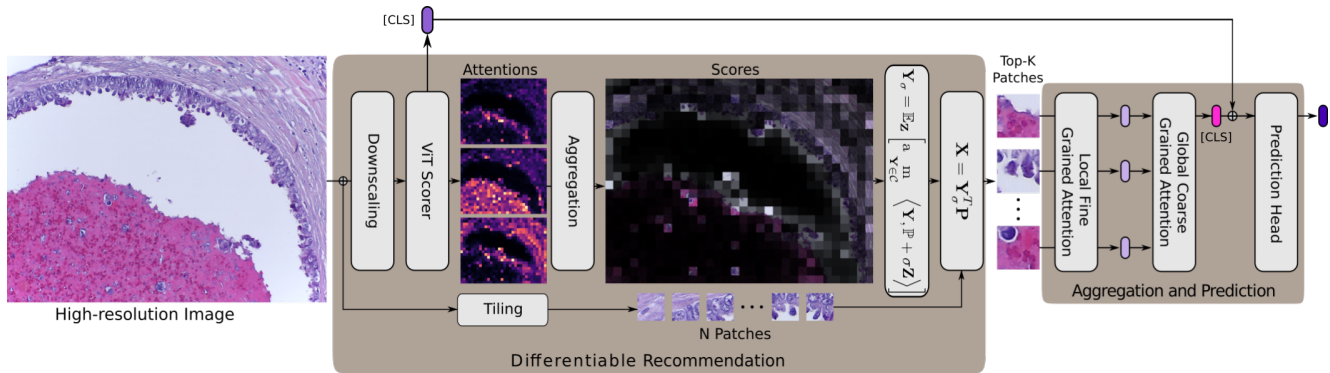
**Figure 2: An overview of the proposed ScoreNet**. The recommendation stage provides tissue-level features, and **differentiably selects** the most discriminative high-resolution patches. The aggregation stage independently extracts cell-level features and embeds the patches via a *local fine-grained attention* mechanism and endows them with contextual information with the *global coarse-grained attention* mechanism.

aforementioned methods rely on a pipeline where the patch embedding and patch contextualization tasks are disconnected w.r.t. the gradient flow. For example, [19] processes representative patches extracted by an external tool [20]. Thus, their patch extraction is fixed and not data-driven as ours. Alternatively, [39] resort to using a multiple field-of-views/resolutions strategy to endow local patches with contextual information. In another line of research, graph neural network (GNN)-based methods [47, 31] have been proposed to capture global contextual information. These approaches build a graph model that operates on the cell-level structure or combines the cell-level and tissue-level context. However, graph generation can be cumbersome and computationally intensive, prohibiting its use in real-time applications. Recently, SSL methods [22, 21, 38] have demonstrated their capabilities to improve classification for histopathological images. Most of these methods harness pretext tasks, e.g., contrastive pre-training, to learn semantically meaningful features. Nonetheless, the CNN back-bone used in these approaches inevitably abandons learning of global context features. The transformer-based architectures [43, 23] can be an alternative solution for processing images as a de-structured patch sequence and capturing their global dependencies. More recently, hybrid-based vision transformer models [7, 34, 43] have been used in digital pathology, either based on MIL framework [34] or SSL pre-training [43] on unlabeled histopathological images. Nevertheless, these methods process the whole image uniformly and do not allow dynamic extraction of the region of interest.

**Mixing Data-Augmentation Methods.** Recently, mixing data augmentations strategies [42, 45, 45] have been proposed to enhance the generalization capabilities of deep network classifiers. These improvements are further exacerbated when the augmentations model the interactions be-

tween the classes [45]. These methods create a new augmented sample by cutting an image region from one image and pasting it on another image, while a convex combination of their labels gives the ground-truth label of the new sample. Despite the strong performances of the existing methods, none of them is genuinely satisfying as they either create samples that exhibit atypical local features as in MixUp [46] or produce potentially mislabeled samples as in CutMix [45]. CutMix approach has been improved by [6] via re-weighting the mixing factor w.r.t. the sum of the attention map values in the randomly sampled image region, which is still at risk of producing mislabelled samples. In addition, recent CutMix based augmentation methods [42, 40] bear additional disadvantages. For example, Attentive CutMix [42] requires an auxiliary pre-trained model to select the most salient patches from the source image and disregards the location of the informative regions in the target image. SaliencyMix [40] assumes that discriminative parts in an image are highly correlated with the saliency map, which is typically not the case for histopathological images.

## 3. Methods

**Model Overview.** An overview of the proposed training pipeline for H&E stained histology TRoIs' representation learning is illustrated in Fig. 2. Histopathological image classification requires capturing cellular and tissue-level microenvironments and learning their respective interactions. Motivated by the above, we propose an efficient transformer, ScoreNet that captures the cell-level structure and tissue-level context at the most appropriate resolutions. Provided sufficient contextual information, we postulate and empirically verify that a tissue's identification can be achieved by only attending to its sub-region in a high-resolution image. As a consequence, ScoreNet en-

compasses two stages. The former (*differentiable recommendation*) provides contextual information and selects the most informative high-resolution regions. The latter (*aggregation and prediction*) processes the recommended regions and the global information to identify the tissue and model their interactions simultaneously.

More precisely, the recommendation stage is implemented by a ViT and takes as input a downscaled image to produce a semantic distribution over the high-resolution patches. Then, the most discriminative high-resolution patches for the end task are **differentiably extracted**. These selected patches (tokens) are then fed to a second ViT implementing the *local fine-grained attention* module, which identifies the tissues represented in each patch. Subsequently, the **embedded patches attend to one another via a transformer encoder** (*global coarse grained attention*). This step concurrently refines the tissues' representations and model their interactions. As a final step, the concatenation of the [CLS] tokens from the recommendation's stage and that of the *global coarse-grained attention*'s encoder produces the image's representation. Not only does ScoreNet's workflow allows for a significantly increased throughput compared to SOTA methods (see Table 4), it further enables the independent pre-training and validation of its constituent parts.

### 3.1. ScoreNet

**Semantic Regions Recommendation.** Current MIL-based approaches [18, 22] based on patch-level features aggregation often process histopathological images uniformly and discard the inter-patches interactions. To alleviate these issues, we exploit **a differentiable recommendation stage to extract discriminative image regions relevant to the classification**. More specifically, we leverage the self-attention map of a ViT as a distribution of the semantic content. Towards that end, the high-resolution image is first downscaled by a factor $s$ and subsequently fed to the recommendation's stage ViT. The resulting self-attention map captures the contribution of each patch to the overall representation. Let's assume a ViT, that processes a low-resolution image $x_l \in \mathbb{R}^{C \times h \times w}$ encompassing $N$ patches of dimension $P_l \times P_l$. The attended patches (tokens) of the $(L-1)$ layer are conveniently represented as a matrix $\mathbf{Z} \in \mathbb{R}^{(N+1) \times d}$, where $d$ is the embedding dimension of the model, and the extra index is due to the [CLS] token. Up to the last MLP and for a single attention head, the representation of the complete image is given by:

$$y_{[\text{CLS}]} = \underbrace{\text{softmax}(a_1^T)}_{1 \times (N+1)} \underbrace{\mathbf{Z}\mathbf{W}_{\text{val}}}_{(N+1) \times d} \qquad (1)$$

where $\mathbf{W}_{\text{val}} \in \mathbb{R}^{d \times d}$ is the value matrix, and $a_1^T$ is the first row of the self-attention matrix $\mathbf{A}$:

$$\mathbf{A} = \mathbf{Z}\mathbf{W}_{\text{qry}} \left(\mathbf{Z}\mathbf{W}_{\text{key}}\right)^T \qquad (2)$$

where $\mathbf{W}_{\text{qry}}$ and $\mathbf{W}_{\text{key}}$ are the query and key matrices, respectively. The first row of the self-attention matrix captures the contribution of each token to the overall representation (Eq. 1). This is in line with the discriminative capacity of the [CLS] token that patches having the highest contribution are the ones situated in the highest semantic regions of the images. The distribution of the semantic content over the patches is therefore defined as:

$$\mathbb{P}_{\text{patch}} = \text{Softmax}(\tilde{a}_1^T) \in \mathbb{R}^N \qquad (3)$$

where $\tilde{a}_1$ stands for $a_1$ without the first entry, namely the one corresponding to the [CLS] token. Since ViTs typically encompasses multiple heads, we propose to add an extra learnable parameter, which weights the relative contributions of each head to the end task; after aggregation of the multiples self-attention maps, the formulation is identical to that of Eq. 3.

Concurrently with acquiring the above defined semantic distribution, the high-resolution image, $x_h \in \mathbb{R}^{C \times H \times W}$, is tiled in a regular grid of large patches ($P_h \times P_h$), stored in a tensor $\mathbf{P} \in \mathbb{R}^{N \times C \times P_h \times P_h}$. At inference time, a convenient way to select the $K$ most semantically relevant high-resolution regions is to encode the *top-K* indices as one-hot vectors: $\mathbf{Y} \in \mathbb{R}^{N \times K}$, and to extract the corresponding $K$ patches, $\mathbf{X} \in \mathbb{R}^{K \times C \times P_h \times P_h}$ via:

$$\mathbf{X} = \mathbf{Y}^T \mathbf{P} \qquad (4)$$

At training time, since the above formulation is not differentiable, we propose to adopt the differentiable approach of [10]. Following the perturbed optimizers scheme, the *top-K* operation is bootstrapped by applying a Gaussian noise, $\sigma\mathbf{Z}$, to the semantic distribution. The noisy indicators, $\mathbf{Y}_\sigma$, are subsequently computed as:

$$\mathbf{Y}_\sigma = \mathbb{E}_{\mathbf{Z}} \left[ \arg\max_{\mathbf{Y} \in \mathcal{C}} \left\langle \mathbf{Y}, \tilde{\mathbb{P}} + \sigma\mathbf{Z} \right\rangle \right] \qquad (5)$$

where $\sigma$ is the standard deviation of the noise, $\tilde{\mathbb{P}} \in \mathbb{R}^{N \times K}$ is obtained by broadcasting $\mathbb{P}_{\text{patch}}$ to match the dimension of $\mathbf{Y}$, and $\mathcal{C}$ is a restriction of the domain ensuring the equivalence between solving Eq. 5 and the *top-K* operation [10]. The extraction of the high-resolution regions follows the procedure described in Eq. 4. Similarly, the gradient of the indicators w.r.t. the semantic distribution, $\mathbb{P}_{\text{patch}}$ can be computed as:

$$\nabla\mathbf{Y}_\sigma = \mathbb{E}_{\mathbf{Z}} \left[ \arg\max_{\mathbf{Y} \in \mathcal{C}} \left\langle \mathbf{Y}, \tilde{\mathbb{P}} + \sigma\mathbf{Z} \right\rangle \mathbf{Z}^T / \sigma \right] \qquad (6)$$
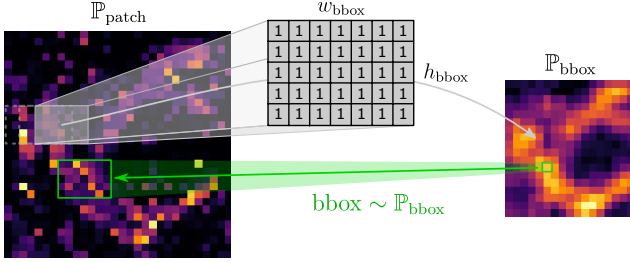
**Figure 3: The bounding box selection scheme for ScoreMix.** The score distribution for each bounding box ($\mathbb{P}_{\text{bbox}}$) is obtained by convolving the patch distribution map $\mathbb{P}_{\text{patch}}$ with a kernel of **1**s of the bounding box dimensions ($h_{\text{bbox}}, w_{\text{bbox}}$). The bbox is then sampled from $\mathbb{P}_{\text{bbox}}$, which we often refer to as $\mathbb{P}_{\text{source}}$ or $\mathbb{P}_{\text{target}}$.

**Computational Complexity.** Vision transformers heavily rely on the attention mechanism to learn a high-level representation from low-level regions. The underlying assumption is that the different sub-regions of the image are not equally important for the overall representation. Despite this key observation, the computation cost dedicated to a sub-region is independent of its contribution to the high-level representation, which is inefficient. Our ScoreNet attention mechanism overcomes this drawback by learning to attribute more resources to regions of high interest. For a high-resolution input image $x_h \in \mathbb{R}^{C \times H \times W}$, the asymptotical time and memory cost is $\mathcal{O}\left(\left(\frac{H}{s \cdot P_l} \cdot \frac{W}{s \cdot P_l}\right)^2\right)$, when the recommendation stage uses inputs downscaled by a factor $s$ and processes them with a patch size of $P_l$. The derivation of this cost, including that of the recommendation stage, which is independent of the input size, can be found in the **Appendix**.

### 3.2. ScoreMix

We propose a new mixing data augmentation for histopathological images by learning the **distribution of the semantic image regions** using the learned self-attention for `[CLS]` token of the ViT without requiring architectural changes or additional loss. More formally, let $x_s, x_t \in \mathbb{R}^{C \times H \times W}$ be the source and target images respectively and let $y_s$ and $y_t$ be their corresponding labels. We aim to mix the source and target samples to generate a new training example $(x_m, y_m)$. To do so, we first compute the semantic distributions using the current parameters of the model and the input samples; namely, we compute $\mathbb{P}_{\text{source}}(x_s, \theta)$ and $\mathbb{P}_{\text{target}}(x_t, \theta)$. Given these distributions and a randomly defined bounding box size, we sample the cutting and pasting locations from the source and target distributions, respec-

tively:

$$
\begin{aligned}
M_s &\sim \frac{1}{Z_s} \cdot \mathbb{P}_{\text{source}}(x_s, \theta, \lambda) \\
M_t &\sim \frac{1}{Z_t} \cdot (1 - \mathbb{P}_{\text{target}}(x_t, \theta, \lambda))
\end{aligned}
\tag{7}
$$

where $Z_s$ and $Z_t$ are normalization constants, and $1 - \lambda \sim \mathcal{U}([0, 1])$ defines the strength of the mixing, i.e. the size of the bounding box. The locations of the cutting and pasting regions are encoded as binary masks, i.e., $M_s, M_t \in \{0, 1\}^{H \times W}$, where a value of 1 encodes for a patch in the cutting/pasting region. Under the above formalism, the mixing operation can be defined as:

$$
\begin{aligned}
x_m &= (\mathbf{1} - M_t) \odot x_t \\
M_t \otimes x_m &\leftarrow M_s \otimes x_s \\
y_m &= \lambda y_t + (1 - \lambda) y_s
\end{aligned}
\tag{8}
$$

where **1** is a mask of ones, $\odot$ denotes the element-wise multiplication, and $\otimes$ indicates an indexing w.r.t. a mask. **Computing the Semantic Distributions.** Computing the semantic distributions of the target and source images is an essential part of the pipeline as it allows for a data-driven selection of the cutting/pasting sites, thereby avoiding the pitfalls of random selection. When the size of the bounding box matches that of a single patch, the distribution can be directly deduced from the self-attention map, as described in Sec. 3. As a consequence, and when the bounding box's size matches that of a single patch, the semantic distribution can be directly obtained from $\mathbb{P}_{\text{patch}}$ (see Eq. 3). In practice, we would typically use bounding boxes encompassing more than a single patch. In that case, the distribution of the semantic content at the bounding box resolution can be obtained by a local aggregation of the distribution above:

$$
\mathbb{P}_{\text{bbox}}(i) \propto \sum\nolimits_{j \in \mathcal{N}(i)} \mathbb{P}_{\text{patch}}(j)
\tag{9}
$$

where $\mathcal{N}(i)$ returns the indices of the patches situated in the bounding box whose top left corner is the patch $i$. In practice, this can be efficiently implemented by first unflattening the patch distribution $\mathbb{P}_{\text{patch}}$, and convolving it with a kernel of ones and of the same dimension as the desired bounding box (see Fig. 3).

## 4. Experiments

**Datasets.** The primary dataset used in our experiments is the BReAst Carcinoma Sub-typing (**BRACS**) [31]. BRACS consists of 4391 RoIs acquired from 325 H&E stained breast carcinoma WSI (at $0.25$ $\mu$m/pixel) with varying dimensions and appearances. Each RoI is annotated with one of the seven classes: Normal, Benign, Usual Ductal Hyperplasia (UDH), Atypical Ductal Hyperplasia (ADH), Flat Epithelial Atypia (FEA), Ductal Carcinoma In Situ (DCIS),

and Invasive. Our experiments follow the same data splitting scheme as [31] for training, validation, and test set at the WSI level to avoid test leakage. In addition, we use publicly available BreAst Cancer Histology (**BACH**) dataset [1] to show ScoreNet generalization capabilities. It contains 400 training and 100 test images from four different breast cancer types: Normal, Benign, In Situ, and Invasive. All images have a fixed size of $1536 \times 2048$ pixels and a pixel scale of $0.42 \times 0.42 \ \mu$m. To assess the interpretability of ScoreNet, we further evaluate our model on the **CAMELYON16** dataset [3] for binary tumour classification. We extract a class-balanced dataset of $1920 \times 1920$ pixels from high-resolution WSIs.

**Experimental Setup.** We base ScoreNet' ViTs, namely the one used by the recommendation stage and by the local fine-grained attention mechanism on a modified ViT-tiny architecture (see **Appendix**) and follow the self-supervised pre-training scheme of [5] for both of the aforementioned ViTs. Noteworthy that an end-to-end pre-training of ScoreNet is also feasible. After pre-training, the ScoreNet is optimized using the SGD optimizer (momentum=0.9) with a learning rate chosen with the linear scaling rule [14] ($lr = 10^{-2} \cdot$ batchsize$/256 = 3.125 \cdot 10^{-4}$) annealed with a cosine schedule until $10^{-6}$. ScoreNet is finetuned for 15 epochs with a batch-size of 8. We empirically determine the top $K = 20$ regions, and a downscaling factor $s = 8$ by a hyperparameter sweep (cf. ablation experiment in the **Appendix**). All experiments are implemented in PyTorch 1.9 [28] using a single GeForce RTX3070 GPU.

## 4.1. TRoIs Classification Results and Discussion

In Table 1, we compare the TRoIs classification performance of ScoreNet on the BRACS dataset against the state-of-the-arts, including MIL-based [27, 34, 25], GNN-based, e.g., [31], and self-supervised transformer-based [43] approaches. The first MIL-based baseline [27] aggregates independent patch representations from the penultimate layer of a ResNet-50 [16] pre-trained on ImageNet [11]. The patch model is further finetuned on $128 \times 128$ patches at different magnification, e.g., $10\times$, $20\times$ or $40\times$. The latter operate either on multi- or single-scale images to benefit from varying levels of context and resolution. Similarly, we report the performances of the recent MIL-based methods, TransMIL [34], and CLAM [25] using the original implementations and setup. Both methods are tested with different magnifications (see Table 1). Additionally, the single-head (-SB) and multi-head (-MB) variants of CLAM are used with the small (-S) and big (-B) versions of the models (see CLAM's implementation). We further use various GNN-based baselines, particularly HACT-Net [31], the current SOTA approach for TRoIs classification on the BRACS. Finally, we report the performance of the recent self-supervised transformer approach, TransPath [43],

which is a hybrid transformer/convolution-based architecture. ScoreNet reaches a new state-of-the-art weighted F1-score of 64.4% on the BRACS TRoIs classification task outperforming the second-best method, HACT-Net, by a margin of 2.9% (Table 1). The results are reported for two variants of ScoreNet, namely ScoreNet/4/1 and ScoreNet/4/3, which use the four last [CLS] tokens of the scorer and the last or the three last [CLS] tokens from the coarse attention mechanism (aggregation stage). ScoreNet/4/3 variant puts more emphasis on the features available at ($40\times$), whereas ScoreNet/4/1 is more biased towards the global representation available at ($5\times$) (with a downscaling factor $s = 8$). One can observe that both model variants significantly outperform the existing baseline in terms of weighted F1-scores and for almost every class. More interestingly, the architectural differences directly translate to differences in the classification results. ScoreNet/4/3 is more suitable for classes where the **discriminative features are at the cell level than** ScoreNet/4/1, **which is more suited when the tissue organization is the discriminative criterion**. Nonetheless, both of these architectures indeed benefit from the information available at each scale. This observation is well supported by the classification results obtained when a linear layer is trained independently on the scorer's [CLS] tokens (Lin. scorer's [CLS] in Table 1) or using only the [CLS] tokens from the aggregation stage (Lin. encoder's [CLS] in Table 1). Despite the difference in results between the two model variants, it is clear that they both perform worse when separated, which indicates that the representations of both stages are complementary. In brief, ScoreNet allows for an easily tuning to meet prior inductive biases on the ideal scale for a given task.

**ScoreMix & Data-Regime Sensitivity.** We also show that ScoreNet equipped with the proposed ScoreMix augmentation achieves superior TRoIs classification performances compared to CutMix [45] and SaliencyMix [40] augmentations for different data regimes, e.g., low-regime with only 10% of the data. **Our proposed ScoreMix outperforms SOTA methods with only 50% of the data** and is on-par or better than most baselines with only 20% of the data (Table 2). We argue that these improvements are primarily due to the generation of more coherent sample-label pairs under the guidance of the learned semantic distribution. This alleviates randomly cutting and pasting non-discriminative patches, as is the case with CutMix. Our results further support that image saliency used in the SaliencyMix is not correlated with discriminative regions.

**Generalization Capabilities.** To gauge the generalization capabilities of ScoreNet compared to other current SOTA methods, e.g., HACT-Net [29], we leverage two external evaluation datasets, namely CAMELYON16 and

**Table 1: Comparison with the prior art for TRoIs classification** using weighted and class-wise F1-scores averaged over three independent runs on the BRACS dataset. The best results are in **bold**. ScoreNet/x/y refers to an instance of ScoreNet using the recommendation module's last x [CLS] tokens and the last y tokens from the global coarse-grained attention.

| | Method | Normal | Benign | UDH | ADH | FEA | DCIS | Invasive | Weighted F1 |
|---|---|---|---|---|---|---|---|---|---|
| MILs | Agg-Penultimate (10×) [36] | 48.7 ± 1.7 | 44.3 ± 1.9 | 45.0 ± 5.0 | 24.0 ± 2.8 | 47.0 ± 4.3 | 53.3 ± 2.6 | 86.7 ± 2.6 | 50.8 ± 2.6 |
| | Agg-Penultimate (20×) [36] | 42.0 ± 2.2 | 42.3 ± 3.1 | 39.3 ± 2.0 | 22.7 ± 2.5 | 47.7 ± 1.2 | 50.3 ± 3.1 | 77.0 ± 1.4 | 46.8 ± 2.2 |
| | Agg-Penultimate (40×) [36] | 32.3 ± 4.6 | 39.0 ± 0.8 | 23.7 ± 1.7 | 18.0 ± 0.8 | 37.7 ± 2.9 | 47.3 ± 2.0 | 70.7 ± 0.5 | 39.4 ± 1.9 |
| | Agg-Penultimate (10× + 20×) [36] | 48.3 ± 2.0 | 45.7 ± 0.5 | 41.7 ± 5.0 | 32.3 ± 0.9 | 46.3 ± 1.4 | 59.3 ± 2.0 | 85.7 ± 1.9 | 52.3 ± 1.9 |
| | Agg-Penultimate (10× + 20× + 40×) [36] | 50.3 ± 0.9 | 44.3 ± 1.2 | 41.3 ± 2.5 | 31.7 ± 3.3 | 51.7 ± 3.1 | 57.3 ± 0.9 | 86.0 ± 1.4 | 52.8 ± 1.9 |
| | CLAM-SB/S (10×) [25] | 39.6 ± 4.6 | 45.5 ± 4.9 | 34.7 ± 2.0 | 30.4 ± 6.7 | 68.8 ± 1.9 | 64.3 ± 0.8 | 84.2 ± 2.6 | 53.9 ± 1.9 |
| | CLAM-SB/S (20×) [25] | 50.2 ± 3.2 | 45.5 ± 1.8 | 32.2 ± 1.6 | 25.5 ± 4.2 | 69.6 ± 1.0 | 60.8 ± 2.7 | 84.2 ± 1.6 | 54.0 ± 0.7 |
| | CLAM-SB/S (40×) [25] | 47.0 ± 5.2 | 38.8 ± 1.8 | 30.0 ± 7.7 | 29.4 ± 2.9 | 65.9 ± 1.2 | 52.2 ± 1.3 | 76.7 ± 1.6 | 49.9 ± 0.8 |
| | CLAM-SB/B (10×) [25] | 46.4 ± 6.0 | 42.4 ± 2.8 | 33.1 ± 1.0 | 29.3 ± 2.1 | 67.4 ± 1.4 | 63.0 ± 4.5 | 84.4 ± 2.1 | 53.7 ± 1.9 |
| | CLAM-SB/B (20×) [25] | 56.2 ± 1.2 | 42.3 ± 4.4 | 27.4 ± 2.4 | 30.1 ± 4.0 | 68.5 ± 2.1 | 60.9 ± 2.1 | 84.6 ± 1.2 | 54.3 ± 1.5 |
| | CLAM-SB/B (40×) [25] | 42.8 ± 1.1 | 43.3 ± 2.8 | 33.8 ± 0.7 | 29.6 ± 3.6 | 64.1 ± 2.6 | 52.0 ± 3.8 | 78.8 ± 2.2 | 50.5 ± 0.9 |
| | CLAM-MB/S (10×) [25] | 42.5 ± 3.3 | 43.4 ± 3.6 | 31.4 ± 3.2 | 32.1 ± 4.8 | 67.5 ± 2.2 | 59.7 ± 2.4 | 83.8 ± 2.0 | 52.9 ± 1.7 |
| | CLAM-MB/S (20×) [25] | 56.6 ± 0.8 | 47.4 ± 0.9 | 33.5 ± 5.2 | 17.0 ± 1.5 | 70.3 ± 1.1 | 56.9 ± 1.6 | 84.9 ± 1.2 | 53.8 ± 0.6 |
| | CLAM-MB/S (40×) [25] | 50.2 ± 7.7 | 39.3 ± 2.9 | 38.6 ± 2.4 | 26.5 ± 8.9 | 69.4 ± 2.6 | 54.1 ± 3.3 | 82.9 ± 2.5 | 52.9 ± 0.8 |
| | CLAM-MB/B (10×) [25] | 39.7 ± 1.6 | 41.0 ± 2.6 | 34.5 ± 1.0 | 29.8 ± 4.7 | 66.8 ± 1.5 | 63.4 ± 1.0 | 83.5 ± 0.4 | 52.7 ± 0.9 |
| | CLAM-MB/B (20×) [25] | 59.4 ± 2.0 | 47.7 ± 1.2 | 31.7 ± 0.7 | 20.1 ± 3.4 | 68.3 ± 0.4 | 59.9 ± 1.7 | 86.8 ± 0.6 | 54.8 ± 1.0 |
| | CLAM-MB/B (40×) [25] | 47.3 ± 3.2 | 39.5 ± 1.5 | 38.8 ± 4.5 | 30.2 ± 6.3 | 68.2 ± 1.9 | 59.2 ± 2.9 | 82.1 ± 2.7 | 53.5 ± 1.3 |
| GNNs | CGC-Net [47] | 30.8 ± 5.3 | 31.6 ± 4.7 | 17.3 ± 3.4 | 24.5 ± 5.2 | 59.0 ± 3.6 | 49.4 ± 3.4 | 75.3 ± 3.2 | 43.6 ± 0.5 |
| | Patch-GNN (10×) [2] | 52.5 ± 3.3 | 47.6 ± 2.2 | 23.7 ± 4.6 | 30.7 ± 1.8 | 60.7 ± 5.3 | 58.8 ± 1.1 | 81.6 ± 2.2 | 52.1 ± 0.6 |
| | Patch-GNN (20×) [2] | 43.9 ± 4.2 | 43.4 ± 3.2 | 19.5 ± 2.3 | 25.7 ± 2.9 | 55.6 ± 2.1 | 52.9 ± 1.8 | 79.2 ± 1.1 | 47.1 ± 0.7 |
| | Patch-GNN (40×) [2] | 41.7 ± 3.1 | 32.9 ± 1.0 | 25.1 ± 3.7 | 25.6 ± 2.0 | 49.5 ± 3.5 | 48.6 ± 4.2 | 71.6 ± 5.1 | 43.2 ± 0.6 |
| | TG-GNN [29] | 58.8 ± 6.8 | 40.9 ± 3.0 | 46.8 ± 1.9 | 40.0 ± 3.6 | 63.7 ± 10.5 | 53.8 ± 3.9 | 81.1 ± 3.3 | 55.9 ± 1.0 |
| | CG-GNN [29] | 63.6 ± 4.9 | 47.7 ± 2.9 | 39.4 ± 4.7 | 28.5 ± 4.3 | 72.1 ± 1.3 | 54.6 ± 2.2 | 82.2 ± 4.0 | 56.6 ± 1.3 |
| | CONCAT-GNN | 61.0 ± 4.5 | 43.1 ± 2.3 | 42.0 ± 4.7 | 26.1 ± 3.7 | 71.3 ± 2.1 | 60.8 ± 3.7 | 85.4 ± 2.7 | 57.0 ± 2.3 |
| | HACT-Net [29] | 61.6 ± 2.1 | 47.5 ± 2.9 | 43.6 ± 1.9 | 40.4 ± 2.5 | 74.2 ± 1.4 | **66.4 ± 2.6** | 88.4 ± 0.2 | 61.5 ± 0.9 |
| Transformers | TransPath [43] | 58.5 ± 2.5 | 43.1 ± 1.8 | 34.9 ± 5.2 | 38.3 ± 6.0 | 66.9 ± 0.8 | 61.4 ± 1.2 | 85.0 ± 1.4 | 56.7 ± 2.0 |
| | TransMIL (10×) [34] | 38.7 ± 5.4 | 44.0 ± 2.9 | 30.5 ± 4.1 | 31.0 ± 11.8 | 68.1 ± 2.6 | 61.8 ± 1.9 | 87.3 ± 2.6 | 53.2 ± 1.1 |
| | TransMIL (20×) [34] | 51.0 ± 0.1 | 44.5 ± 2.9 | 31.6 ± 2.1 | 31.4 ± 10.3 | 71.3 ± 4.8 | 63.0 ± 2.8 | 89.9 ± 1.6 | 56.2 ± 1.6 |
| | TransMIL (40×) [34] | 47.6 ± 9.8 | 42.9 ± 3.6 | 41.5 ± 5.3 | 38.4 ± 5.9 | 72.7 ± 2.6 | 62.7 ± 2.9 | 87.1 ± 3.9 | 57.5 ± 0.7 |
| | Lin. encoder's [CLS] | 52.7 ± 9.4 | 35.6 ± 3.4 | 34.5 ± 6.7 | 25.1 ± 3.6 | 53.5 ± 9.8 | 38.7 ± 2.8 | 63.3 ± 7.6 | 43.8 ± 3.4 |
| | Lin. scorer's [CLS] | 57.5 ± 4.2 | 48.8 ± 5.5 | 42.7 ± 3.5 | 42.7 ± 7.4 | 74.3 ± 5.2 | 60.5 ± 2.4 | 90.6 ± 0.2 | 60.9 ± 3.1 |
| | ScoreNet/4/1 | **64.6 ± 2.2** | 52.6 ± 2.8 | **48.4 ± 2.2** | **47.4 ± 2.4** | 77.9 ± 0.7 | 59.3 ± 1.1 | 90.6 ± 1.5 | 64.1 ± 0.7 |
| | ScoreNet/4/3 | 64.3 ± 1.5 | **54.0 ± 2.2** | 45.3 ± 3.4 | 46.7 ± 1.0 | **78.1 ± 2.8** | 62.9 ± 2.0 | **91.0 ± 1.4** | **64.4 ± 0.9** |

**Table 2: Comparison with SOTA Mixup-based augmentation methods [45, 40] and the standard random augmentation strategy** using various fractions of the BRACS dataset and identical distribution for the bounding boxes' sizes.

| Dataset | Random Aug. | CutMix [45] | SaliencyMix [40] | **ScoreMix** |
|---|---|---|---|---|
| BRACS 10% | 52.9 ± 2.4 | 53.7 ± 2.9 | 53.5 ± 2.7 | **55.9 ± 1.9** |
| BRACS 20% | 57.6 ± 1.8 | 58.0 ± 1.4 | 57.8 ± 1.0 | **58.7 ± 0.8** |
| BRACS 50% | 60.4 ± 1.8 | 61.2 ± 2.5 | 59.8 ± 2.4 | **62.3 ± 0.6** |
| BRACS 100% | 62.7 ± 1.6 | 63.1 ± 1.1 | 62.8 ± 1.2 | **64.0 ± 0.7** |

**BACH.** After training on the BRACS dataset, the weights of ScoreNet are frozen. To evaluate the quality of the learned features, we either train a linear classifier on top of the frozen features or apply a $k$-nearest-neighbor classifier ($k = 1$) without any finetuning. We perform stratified 5-fold cross-validation. For HACT-Net, we use the available pre-trained weights and follow the implementation of [31]. As HACT-Net sometimes fails to generate embeddings and to have a fair comparison, we only evaluate the samples for which HACT-Net could successfully produce embeddings (around 95% of the BACH and 80% of the CAMELYON16 datasets). Experimental results in Table 3 demonstrate the superiority of ScoreNet in learning generalizable features. It further demonstrates the robustness of ScoreNet to changes in magnification. Indeed, the model is pre-trained on BRACS (40×), while BACH's images were acquired at a magnification of 20×. Furthermore, the CAMELYON16 dataset contains WSIs collected from lymph nodes in the vicinity of the breast, while BRACS contains WSIs collected by mastectomy or biopsy (i.e., directly in the breast). The excellent knowledge transfer between the two datasets highlights the transferability of features learned by ScoreNet in various use cases.

**Interpretability?** To probe the internal behavior of ScoreNet, we finetune the model on CAMELYON16 images using image-level labels only. At test time, we scrutinize the learned semantic distributions of the tumour-positive images. The semantic distributions depicted in Fig. 4 seems to indicate that ScoreNet learns to identify the tumour area and interpret **cancer-related morphological information**, while never having been taught to do so. Quantitatively, we observe that, on average, 74.6% of

**Table 3: Generalization capabilities** of **ScoreNet** compared to HACT-Net trained on BRACS and evaluated on the BACH's annotated images and 1000 images from CAMELYON16, respectively. The weighted F1-scores over a stratified 5-fold cross-validation fold is reported.

| | BRACS → BACH | | BRACS → CAMELYON16 | |
|---|---|---|---|---|
| | Linear | $k$-NN | Linear | $k$-NN |
| TransPath [43] | $61.8 \pm 4.8$ | $72.0 \pm 2.9$ | $58.1 \pm 4.8$ | $69.9 \pm 2.5$ |
| TransMIL [34] | $46.5 \pm 10.2$ | $74.0 \pm 4.8$ | $59.8 \pm 3.0$ | $60.8 \pm 5.3$ |
| CLAM-SB/S [25] | $53.3 \pm 13.0$ | $69.8 \pm 4.5$ | $56.7 \pm 1.9$ | $68.0 \pm 3.5$ |
| CLAM-SB/B [25] | $57.5 \pm 3.6$ | $75.3 \pm 3.1$ | $55.5 \pm 4.1$ | $68.0 \pm 1.5$ |
| HACT-Net [29] | $40.2 \pm 2.8$ | $32.8 \pm 5.8$ | $60.0 \pm 4.6$ | $61.0 \pm 4.2$ |
| **ScoreNet** | $\mathbf{73.4 \pm 3.5}$ | $\mathbf{76.9 \pm 6.1}$ | $\mathbf{81.1 \pm 3.5}$ | $\mathbf{77.0 \pm 4.6}$ |

**Table 4: Inference throughput comparison of ScoreNet, HACT-Net, and SOTA transformer-based architectures**. All models were tested with the same image size and a single GeForce RTX 3070 GPU.

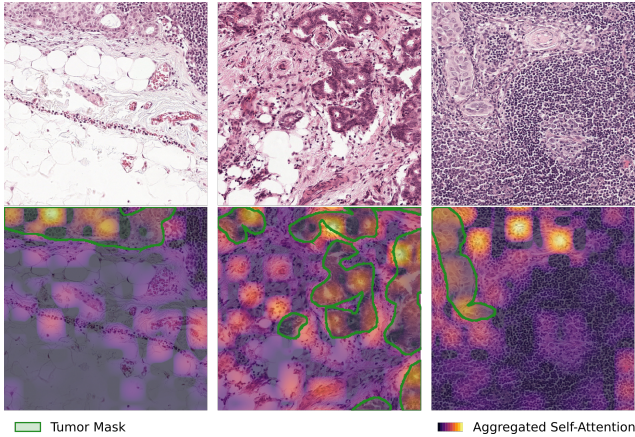| | Image size | Throughput (im./s) | Pre-processing |
|---|---|---|---|
| HACT-Net [29] | $1536 \times 2048$ | $4.95e\text{-}4 \pm 1.40e\text{-}3$ | ✓ |
| Vanilla ViT [13] | $1536 \times 2048$ | $3.8 \pm 0.1$ | - |
| SwinTransformer [24] | $1536 \times 2048$ | $76.8 \pm 0.4$ | ✗ |
| TransPath [43] | $1536 \times 2048$ | $97.6 \pm 3.1$ | ✗ |
| **ScoreNet** | $1536 \times 2048$ | $\mathbf{335.0 \pm 7.9}$ | ✗ |



**Figure 4: ScoreNet Interpretability**. Visualization of the semantic distribution, overlaid with the tumour ground-truth mask on a few samples of the CAMELYON16 dataset. The semantic distributions are obtained from the recommendation stage, i.e., at low-resolution. ScoreNet is pre-trained on BRACS and finetuned on CAMELYON16.

the 20 patches selected from positive images are tumour-positive. Furthermore, we report an average image-wise AuC of 73.6% when interpreting the probability of the recommendation stage to sample a patch as the probability of it being tumour-positive.

**Ablation on Efficacy of ScoreNet.** The critical aspect of ScoreNet is its improved efficiency compared to other transformer-based architectures. This improvement is due to the choice of a hierarchical architecture and the exploitation of redundancy in histology images. At inference time, we expect a gain in throughput compared to the vanilla ViT of the order of the squared downscaling factor, $s$, (see **Appendix**), typically $s^2 = 64$, which is well reflected in practice, as shown in Table 4. Due to the self-supervised pre-training, ScoreNet does not require any stain normalization or pre-processing, unlike its competitor HACT-Net. Similarly, ScoreNet yields higher throughput than other SOTA efficient transformers architectures, namely TransPath [43],

and SwinTransformer [24], with throughput around $3\times$ and $4\times$ higher than these methods. The latter observation is interesting considering the linear asymptotic time and memory cost of the SwinTransformer, which is probably a consequence of the fact that SwinTransformers process a lot of uninformative high-resolution patches in the first layer(s).

**Ablation on Shape Cues and Robustness.** We investigate ScoreNet's ability to learn shape-related features. To do so, we study shape cues extracted by the recommendation model via the concatenated `[CLS]` tokens (see Fig. 2). Consequently, we implement shape removal by applying a random permutation of the downscaled image's tokens at test time. With this setup, a weighted F1-score of $59.8 \pm 0.8\%$ is reached, representing a significant drop in performance compared to $64.4 \pm 0.9\%$ without permutation. It demonstrates that *i)* the recommendation stage's concatenated `[CLS]` tokens contribute positively to the overall representation and *ii)* the latter is **not permutation invariant** and thus shape-dependent. In a second experiment, we show the whole recommendation stage is also shape-dependent. To that end, we repeat the same experiment, but the patches are extracted from the permuted images, reaching a weighted F1-score of $59.5 \pm 0.6\%$. We further observe that for a given image, the overlap of the selected patches with and without permutation is, on average, only 15.7%, which indicates that the semantic distribution learned by ScoreNet is shape-dependent.

## 5. Conclusion and Future Work

We have introduced ScoreNet, an efficient transformer-based architecture that dynamically recommends discriminative regions from large histopathological images, yielding rich generalizable representations at an efficient computational cost. In addition, we propose ScoreMix, a new attention-guided mixing augmentation that produces coherent sample-label pairs. We achieve new SOTA results on the BRACS dataset for TRoIs classification and demonstrate ScoreNet's superior throughput improvements compared to previous SOTA efficient transformer-based architectures.

# References

[1] 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] Bulut Aygüneş, Selim Aksoy, Ramazan Gökberk Cinbiş, Kemal Kösemehmetoğlu, Sevgen Önder, and Ayşegül Üner. Graph convolutional networks for region of interest classification in breast histopathology. In *Medical Imaging 2020: Digital Pathology*, volume 11320, page 113200K. International Society for Optics and Photonics, 2020.

[3] 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.

[4] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In *Thirty-fourth Conference on Neural Information Processing Systems (NeurIPS)*, 2020.

[5] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021.

[6] Jie-Neng Chen, Shuyang Sun, Ju He, Philip Torr, Alan Yuille, and Song Bai. Transmix: Attend to mix for vision transformers. *arXiv preprint arXiv:2111.09833*, 2021.

[7] Richard J Chen, Ming Y Lu, Wei-Hung Weng, Tiffany Y Chen, Drew FK Williamson, Trevor Manz, Maha Shady, and Faisal Mahmood. Multimodal co-attention transformer for survival prediction in gigapixel whole slide images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4025, 2021.

[8] Xiangxiang Chu, Bo Zhang, Zhi Tian, Xiaolin Wei, and Huaxia Xia. Do we really need explicit position encodings for vision transformers? *CoRR*, abs/2102.10882, 2021.

[9] Ozan Ciga, Tony Xu, and Anne Louise Martel. Self supervised contrastive learning for digital histopathology. *Machine Learning with Applications*, 7:100198, 2022.

[10] Jean-Baptiste Cordonnier, Aravindh Mahendran, Alexey Dosovitskiy, Dirk Weissenborn, Jakob Uszkoreit, and Thomas Unterthiner. Differentiable patch selection for image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2351–2360, 2021.

[11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.

[12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

[13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021.

[14] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet in 1 hour, 2018.

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.

[16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[17] Le Hou, Dimitris Samaras, Tahsin M Kurc, Yi Gao, James E Davis, and Joel H Saltz. Patch-based convolutional neural network for whole slide tissue image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2424–2433, 2016.

[18] Maximilian Ilse, Jakub Tomczak, and Max Welling. Attention-based deep multiple instance learning. In *International conference on machine learning*, pages 2127–2136. PMLR, 2018.

[19] Shivam Kalra, Mohammed Adnan, Sobhan Hemati, Taher Dehkharghanian, Shahryar Rahnamayan, and Hamid R Tizhoosh. Pay attention with focus: A novel learning scheme for classification of whole slide images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 350–359. Springer, 2021.

[20] Shivam Kalra, Hamid R Tizhoosh, Charles Choi, Sultaan Shah, Phedias Diamandis, Clinton JV Campbell, and Liron Pantanowitz. Yottixel–an image search engine for large archives of histopathology whole slide images. *Medical Image Analysis*, 65:101757, 2020.

[21] Navid Alemi Koohbanani, Balagopal Unnikrishnan, Syed Ali Khurram, Pavitra Krishnaswamy, and Nasir Rajpoot. Self-path: Self-supervision for classification of pathology images with limited annotations. *IEEE Transactions on Medical Imaging*, 2021.

[22] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14318–14328, 2021.

[23] 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.

[24] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer:

Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021.

[25] 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.

[26] Oded Maron and Tomás Lozano-Pérez. A framework for multiple-instance learning. *Advances in neural information processing systems*, pages 570–576, 1998.

[27] Caner Mercan, Selim Aksoy, Ezgi Mercan, Linda G Shapiro, Donald L Weaver, and Joann G Elmore. From patch-level to roi-level deep feature representations for breast histopathology classification. In *Medical Imaging 2019: Digital Pathology*, volume 10956, page 109560H. International Society for Optics and Photonics, 2019.

[28] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32:8026–8037, 2019.

[29] Pushpak Pati, Maria Frucci, and Maria Gabrani. Hact-net: A hierarchical cell-to-tissue graph neural network for histopathological image classification. In *Uncertainty for Safe Utilization of Machine Learning in Medical Imaging, and Graphs in Biomedical Image Analysis: Second International Workshop, UNSURE 2020, and Third International Workshop, GRAIL 2020, Held in Conjunction with MICCAI 2020, Lima, Peru, October 8, 2020, Proceedings*, volume 12443, page 208. Springer Nature, 2020.

[30] Pushpak Pati, Guillaume Jaume, Antonio Foncubierta, Florinda Feroce, Anna Maria Anniciello, Giosuè Scognamiglio, Nadia Brancati, Maryse Fiche, Estelle Dubruc, Daniel Riccio, et al. Hierarchical cell-to-tissue graph representations for breast cancer subtyping in digital pathology. *arXiv e-prints*, pages arXiv–2102, 2021.

[31] Pushpak Pati, Guillaume Jaume, Antonio Foncubierta-Rodríguez, Florinda Feroce, Anna Maria Anniciello, Giosue Scognamiglio, Nadia Brancati, Maryse Fiche, Estelle Dubruc, Daniel Riccio, Maurizio Di Bonito, Giuseppe De Pietro, Gerardo Botti, Jean-Philippe Thiran, Maria Frucci, Orcun Goksel, and Maria Gabrani. Hierarchical graph representations in digital pathology. *Medical Image Analysis*, 75:102264, 2022.

[32] David W Romero and Jean-Baptiste Cordonnier. Group equivariant stand-alone self-attention for vision. In *International Conference on Learning Representations*, 2020.

[33] Dawid Rymarczyk, Adriana Borowa, Jacek Tabor, and Bartosz Zielinski. Kernel self-attention for weakly-supervised image classification using deep multiple instance learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1721–1730, 2021.

[34] Zhuchen Shao, Hao Bian, Yang Chen, Yifeng Wang, Jian Zhang, Xiangyang Ji, and Yongbing Zhang. Transmil: Transformer based correlated multiple instance learn-ing for whole slide image classication. *arXiv preprint arXiv:2106.00908*, 2021.

[35] Krishna Kumar Singh, Dhruv Mahajan, Kristen Grauman, Yong Jae Lee, Matt Feiszli, and Deepti Ghadiyaram. Don't judge an object by its context: Learning to overcome contextual bias. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11070–11078, 2020.

[36] Korsuk Sirinukunwattana, Nasullah Khalid Alham, Clare Verrill, and Jens Rittscher. Improving whole slide segmentation through visual context-a systematic study. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 192–200. Springer, 2018.

[37] Chetan L Srinidhi, Ozan Ciga, and Anne L Martel. Deep neural network models for computational histopathology: A survey. *Medical Image Analysis*, 67:101813, 2021.

[38] Chetan L Srinidhi, Seung Wook Kim, Fu-Der Chen, and Anne L Martel. Self-supervised driven consistency training for annotation efficient histopathology image analysis. *arXiv preprint arXiv:2102.03897*, 2021.

[39] Hiroki Tokunaga, Yuki Teramoto, Akihiko Yoshizawa, and Ryoma Bise. Adaptive weighting multi-field-of-view cnn for semantic segmentation in pathology. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12597–12606, 2019.

[40] AFM Shahab Uddin, Mst Sirazam Monira, Wheemyung Shin, TaeChoong Chung, and Sung-Ho Bae. Saliencymix: A saliency guided data augmentation strategy for better regularization. In *International Conference on Learning Representations*, 2020.

[41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.

[42] Devesh Walawalkar, Zhiqiang Shen, Zechun Liu, and Marios Savvides. Attentive cutmix: An enhanced data augmentation approach for deep learning based image classification. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3642–3646. IEEE, 2020.

[43] Xiyue 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 *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 186–195. Springer, 2021.

[44] Yan Xu, Zhipeng Jia, Yuqing Ai, Fang Zhang, Maode Lai, I Eric, and Chao Chang. Deep convolutional activation features for large scale brain tumor histopathology image classification and segmentation. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 947–951. IEEE, 2015.

[45] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6023–6032, 2019.

[46] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *International Conference on Learning Representations*, 2018.

[47] Yanning Zhou, Simon Graham, Navid Alemi Koohbanani, Muhammad Shaban, Pheng-Ann Heng, and Nasir Rajpoot. Cgc-net: Cell graph convolutional network for grading of colorectal cancer histology images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019.

## A. Appendix Overview

In this appendix, we provide additional ablation studies and experimental details. The remaining of this appendix is organized as follows. In Sec. B we detail the architectural and training details, e.g., parameters choices. Additional ablations are detailed in Sec. C . A detailed derivation of the computational cost is presented in Sec. D. We discuss, in Sec. E, some properties of ScoreMix and present some examples of our proposed ScoreMix augmentation. Finally, the suitability of ScoreNet to learn from uncurated data is evaluated in Sec. F.

## B. Experimental Setup & Datasets

### B.1. Networks Architectures

**ScoreNet.** The proposed ScoreNet architecture comprises two stages: the recommendation and aggregation stages. The former leverages a modified ViT-Tiny to produces the semantic distribution. Similarly, the latter relies on an identical ViT-Tiny to independently embed the selected high-resolution patches (*local fine-grained attention*) and on a transformer encoder to mix the embedded patches (*global coarse-grained attention*). The following parameters of the two identical ViT-Tiny were modified to be tailored for the task:

- `embed_dim=96`.
- `depth=8`.
- `num_heads=4`.
- `mlp_ratio=2`.

These modifications were brought to allow for a self-supervised pre-training with a sufficiently large batch size ($bs \geq 128$), which was reported to be of significant importance to reach good performance [5]. The parameters of the transformer encoder implementing the *global coarse-grained attention* mechanism are:

- `embed_dim=96`.
- `depth=4`.
- `num_heads=4`.
- `mlp_ratio=2`.

Overall ScoreNet's model totals approximately 1.79M parameters.

**SwinTransformer.** SwinTransformers [24] relies on hierarchical architecture attention mechanism, namely intra- and inter-window attentions. The patch-merging operation reduces the time, and memory cost of SwinTransformers [24] significantly, which decreases the total number of tokens by 4, while increasing the embedding by 2. The architecture is modified to accept non-square windows, allowing Swin-Transformers to process non-square images images. The resulting parameters are:

- `patch_size=16`.
- `input_embed_dim_size=24`.
- `output_embed_dim_size=192`.
- `depths=[2, 2, 6, 2]`.
- `num_heads=[3, 6, 12, 24]`.
- `window_size=(6, 8)`.
- `mlp_ratio=4`.

Overall the SwinTransformer model totals approximately 1.77M parameters.

**TransPath.** As described in [43], TransPath's architecture leverages a CNN encoder to jointly reduce the input image's size, extract relevant features, and tile the image in pre-embedded patches. Subsequently, a transformer encoder processes the CNN encoder's features to capture global interactions. The CNN encoder's architecture is as follows:

- `n_convolutions=4`.
- `n_filters=[8, 32, 128, 512]`.
- `kernel_sizes=[(3, 3), (3, 3), (3, 3), (3, 3)]`.
- `pooling_kernel_sizes=[(4, 4), (2, 2), (4, 4), (4, 4)]`.
- `activation=ReLU` [15].

A projection convolution is used to match the desired embedding dimension of the transformer encoder. Its parameters are:

- `n_filters=192`.
- `kernel_sizes=(1, 1)`.

The parameters of the transformer encoder are:

- `embed_dim=192`.
- `depth=4`.
- `num_heads=4`.
- `mlp_ratio=2`.

Each transformer block rely on TransPath's customized token-aggregating and excitation multi-head self-attention (MHSA-TAE) [43]. Overall, TransPath's model totals approximately 1.93M parameters.

**TransMIL.** We adopt the original implementation as provided by the authors [34]. It relies on a ResNet-50 [16] pre-trained on ImageNet [11] to embed the individual $256 \times 256$ patches. Overall, TransMIL's model totals approximately 3.19M parameters (not counting the parameters of the ResNet-50).

**CLAM.** The implementation of CLAM follows the code provided by the authors [25]. It relies on a ResNet-50 [16] pre-trained on ImageNet [11] to embed the individual $256 \times 256$ patches. Overall, the variations of CLAM-(SB/MB)/(S/B) total from 1.32M to 1.46M parameters (not counting the parameters of the ResNet-50).

## B.2. Self-Supervised Pre-training

**Modular Pre-training.** Our modular architecture allows for independent self-supervised pre-trainings of the recommendation stage's ViT and that of the local fine-grained attention mechanism. A two steps pre-training can be beneficial, as it provides the possibility to validate each part independently. Similarly, one of the modules, typically the one implementing the fine-grained local attention, can be trained on an auxiliary annotated dataset or be replaced by a standard pre-trained model.

The self-supervised pre-training follows the guidelines of [5]. Apart from the differences in architectures described in Sec. B.1, minor modifications were made in the projection head to account for the reduced heterogeneity in our datasets compared to that in ImageNet [12]. The modifications are:

- `hidden_dim=1024`.
- `bottleneck_dim=128`.
- `out_dim=1024`.

These modifications are in line with the interpretation of [4] which considers the last linear layer as a projection on a set of learnable centroids and that their number should reflect the level of diversity present in the dataset. For this interpretation to hold, it is required that both the last layer's input and its weights are normalized, which is the case in our setup. The remaining parameters, aside from the position encoding which is discussed in Sec. C, are set to the default values (see [5] for details).

**End-to-end Pre-training.** In some cases, an end-to-end pre-training of ScoreNet is preferable. For that purpose, we experimented with two approaches: DINO and a variant of it for that purpose. The former uses the default values for all parameters but those of the projection head described above. On the contrary, the latter benefits from different augmentations and another pretext task and thereby avoid a potential pitfall of DINO: encouraging contextual bias [35], which occurs when the similarity between the representations of views depicting distinct tissue types is enforced.

In this regard, the set of admissible augmentations are constrained to those that change the pixels' values, but not their locations. Consequently, a given image's different views are bounded to bear identical semantic content. A key aspect of DINO's strong performance is due to enforcing the local-to-global correspondence between the student's local crops and that of the teacher's global crop. To mimic that knowledge distillation mechanism, we encourage the student, which only processes the most discriminative high-resolution patches, to match the teacher's representation, which on the contrary, is based on all high-resolution patches. One can observe that this pretext task enforces local-to-global correspondence while providing a strong supervisory signal to the student's scorer, which has to highlight the most relevant regions for the task to be successfully accomplished.

In that setting, ScoreNet's representation is obtained by the concatenation of the `[CLS]` tokens of the *global coarse-grained attention* module's last two transformer blocks. This representation benefits from global contextual information through the teacher, which processes the whole high-resolution image. The projection head's parameters are identified as described above.

## B.3. Datasets

In addition to the annotated TRoIs from two datasets, namely BRACS [30] and BACH [1], additional sets of unlabeled of images are used to pre-train the models and for various ablations. The sets of unlabeled images are detailed here.

**BRACS.** The BRACS dataset encompasses both the annotated TRoIs and the 547 whole-slide images from which they were extracted. We use the WSIs to create an unlabeled pre-training set. More precisely, two types of auxiliary datasets are extracted from BRACS's WSIs: tiles set at $40\times$ and low-resolution thumbnails set at $\frac{40}{s}\times$, where $s$ is the down-scaling ratio. The former set is used to pre-train the *local fine-grained attention* module, whereas the latter serves to pre-train the recommendation stag's scorer. We experimented with two variants of these paired sets. The first variant is designed for a version of ScoreNet, where the dimension of the finely attended regions is $P_h = 224$, the recommendation stage processes low-resolution patches of dimension $P_l = 16$ and consequently a down-scaling ratio $s = 14$. The second variant is designed for a version of ScoreNet, where the dimension of the finely attended regions is $P_h = 128$, the recommendation stage processes low-resolution patches of dimension $P_l = 16$ and consequently a down-scaling ratio $s = 8$. The resulting sets contain approximately 150k images (for a fair comparison of the two versions, see Sec. C).

The last images are extracted from BRACS to conduct TransPath's self-supervised pre-training. From the WSIs, an unlabeled set of approximately 100k images at $40\times$ are extracted. The images have dimensions $1536 \times 1536$, which is approximately the median dimensions of the annotated TRoIs.

**BACH.** Similarly, the BACH dataset comprises an annotated set of TRoIs and the accompanying 40 whole-slide images. From the WSIs, an unlabeled pre-training set of approximately 11k images at $20\times$ are extracted. The images have the exact dimensions as the annotated TRoIs, $1536 \times 2048$.

**CAMELYON16.** Finally, additional tiles set is extracted from CAMELYON16, which is, to our knowledge, the only

one with patch-level annotations. This set is used to evaluate the pre-training of the fine-grained attention module. The latter is composed of 10k images at $40\times$, of dimensions $128 \times 128$ or $224 \times 224$. It is class-balanced, and any patch which contains tumorous tissue is considered tumour positive. This set is also used to measure the effectiveness of the position encoding on the fine-grained attention module in Sec. C.

## C. Additional Ablations

**Down-Scaling Ratio & Dimensions of the Attended Regions.** A key component of the proposed pipeline is to determine the down-scaling ratio, $s$, and the dimension of the square patches in low-resolution, $P_l \times P_l$, and in high-resolution, $P_h \times P_h$. Considering the well-studied nature of the ViTs scorers, we use the standard patch dimension $P_l = 16$ for the patches in low-resolution. It has been shown that smaller patches ($P_l = 8$ or $P_l = 5$) improve the quality of the learned representations [5], nonetheless the incurred increase in computational and memory cost is unsuitable for our application. For the high-resolution patches, we experiment with two standard patch dimensions: $P_h = 128$ and $P_h = 224$. As the self-attention of the recommendation stage is used as a learnable distribution of the semantic content, there should exist a 1-to-1 mapping between the low-resolution patches and the high-resolution regions that can be extracted. As a consequence, the down-scaling ratio is fully determined by the dimensions of the patches: $s = P_h/P_l$. In our case, it translates to down-scaling ratio of either $s = 8$, or $s = 14$.

To find out which of these two setups is the most suitable for our application, we compare the models obtained by each of them via a weighted $k$ Nearest Neighbours classifier, which has the advantage of being fast and not requiring any finetuning. In Table 6, we compare the classification results on the low-resolution ($\frac{40}{s}\times$) BRACS dataset. We report the results of both the teacher and the student models as well as those obtained by a CNN with comparable capacity and identical pre-training. We do not observe significant differences between the two scales. On the other hand, these differences are much more emphasized when evaluating the same models on the low-resolution ($\frac{20}{s}\times$) BACH dataset (see Table 5). These promising results on the BACH dataset, despite the mismatched scales, are to be credited to the local to global views pre-training method [5].

The quality of the fine-grained attention module is assessed with the aforementioned method on the tile CAMELYON16 dataset introduced in Sec. B.3, and the hereby obtained results are reported in Table 7. In conclusion, we observe that the difference is either marginal (Table 6 & 7) or significantly in favor of the setup where $s = 8$ (Table 5) and therefore we choose this setup for the remaining experiments and architectures. As a side note, the CNN ar-

chitecture performs substantially worth, but it is most likely due to the fact that the DINO [5] method is biased towards ViT architectures. **Positional Encoding.** Without position encoding (PE), a ViT processes tokens as a set and hence completely discards the global shape information; therefore, position encoding is essential. The typical approach is to learn a single matrix of absolute and additive position encoding jointly during the training phase. This approach suffers from two drawbacks: *i)* the absolute encoding of each token's position implies that a pattern is different at every location it occurs, which reduces the sample efficiency [32], and *ii)* as a consequence of the storage of the position encoding in a single matrix, the model treats the input tokens as a 1D sequence and thus mislays the multi-dimensionality of the inputs. The latter is not an issue as long as the input images have the same aspect ratio, as is the case with the local/global crops strategy of DINO [5]. Nonetheless, and as depicted in Fig. 5, this approach fails when the model is fed an image of a different aspect ratio than those used to train the position encoding. As illustrated in Fig. 5, the 2D sine-cos position encoding does not introduce any artifacts when used with images of different resolutions. On the other hand, any absolute position encoding is not a translation equivariant operation, an undesired property for planar images. For these reasons, we experiment with Conditional Position encoding Vision Transformer (CPVT) [8]. This PE is input-dependent and convolution-based; consequently, it is suitable for any input resolution and patch-wise translation-equivariant. Fig. 6 reveals that the PE of border tokens is slightly different due to the needed zero-padding. This finding suggests that the absolute position encoding can be inferred from zero-paddings [8]. We argue that CPVT is well suited to be used conjointly with SCOREMIX as the local processing of the token is convenient for detecting local discontinuity caused by the pasting operation, which is needed to incorporate the added content to the global representation (see Sec. E). In Table 8 and Table 9, we evaluate the discriminability of the features obtained by a pre-training under the DINO framework and with various position encoding methods. Table 9, which reports results on the tile CAMELYON16 dataset (see Sec. B.3), does not provide substantial shreds of evidence in favor of one PE or the other; we postulate that this lack of significant differences is due to the lessened importance of position encoding for the tile dataset. Indeed, at $40\times$ and with tiles of dimension $128 \times 128$, the available features are mostly texture-based, and the relative organization of the patches is less relevant. This claim is well supported by the substantial differences in performance obtained by distinct PE when evaluated on the low-resolution BACH and BRACS datasets (see Table 9). These differences are further exacerbated by the fact that images on which performance is evaluated are either of varied size (BRACS) or at least of a

**Table 5: A weighted $k$ Nearest Neighbors classifier assesses the learned features' discriminability (weighted F1-score) on the low-resolution BACH dataset**. The performances of CNN and ViT-based architectures are compared, and similarly for two down-scaling ratios ($s = 8$ or $s = 14$). We use a 4-fold scheme with $75\%/25\%$ train/test splits.

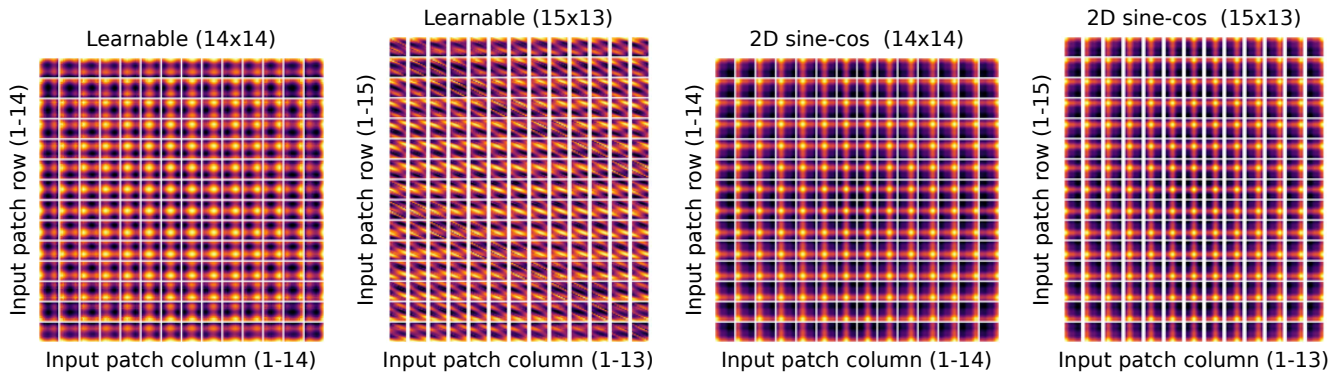| | ViT | | | | CNN | | | |
|---|---|---|---|---|---|---|---|---|
| | Teacher | | Student | | Teacher | | Student | |
| $k$ | $s=14$ | $s=8$ | $s=14$ | $s=8$ | $s=14$ | $s=8$ | $s=14$ | $s=8$ |
| 1 | $71.7 \pm 6.4$ | $\mathbf{78.5 \pm 6.4}$ | $73.6 \pm 5.1$ | $77.4 \pm 5.1$ | $63.6 \pm 5.1$ | $64.4 \pm 1.9$ | $63.8 \pm 3.2$ | $63.9 \pm 2.2$ |
| 5 | $71.5 \pm 1.7$ | $\mathbf{81.7 \pm 3.2}$ | $72.8 \pm 1.9$ | $81.0 \pm 4.0$ | $65.1 \pm 3.3$ | $64.7 \pm 2.1$ | $64.1 \pm 4.6$ | $65.4 \pm 2.7$ |
| 10 | $71.9 \pm 2.4$ | $77.8 \pm 2.8$ | $72.5 \pm 2.5$ | $\mathbf{77.9 \pm 3.4}$ | $62.0 \pm 3.8$ | $58.9 \pm 2.9$ | $61.5 \pm 5.8$ | $61.1 \pm 2.5$ |
| 20 | $71.3 \pm 4.0$ | $76.3 \pm 3.0$ | $72.5 \pm 3.0$ | $\mathbf{76.5 \pm 4.0}$ | $64.0 \pm 6.7$ | $55.5 \pm 1.8$ | $61.0 \pm 9.2$ | $55.4 \pm 2.4$ |
| 50 | $71.2 \pm 4.0$ | $\mathbf{74.7 \pm 4.7}$ | $70.9 \pm 3.3$ | $74.3 \pm 5.7$ | $59.3 \pm 5.3$ | $56.1 \pm 3.2$ | $58.1 \pm 6.2$ | $54.6 \pm 3.9$ |
| 100 | $71.7 \pm 4.1$ | $\mathbf{74.0 \pm 5.5}$ | $71.4 \pm 3.8$ | $73.6 \pm 5.9$ | $57.4 \pm 3.4$ | $50.6 \pm 5.1$ | $56.2 \pm 3.0$ | $48.7 \pm 4.7$ |



**Figure 5: The cosine similarity of a learnable and 2D sine-cos positional encoding is compared**. The learnable positional encoding introduces undesirable artifacts when the aspect ratio changes (*Learnable (15×13)*).

**Table 6: A weighted $k$ Nearest Neighbors classifier assesses the learned features' discriminability (weighted F1-score) on the low-resolution BRACS dataset**. The performances of CNN and ViT-based architectures are compared, and similarly for two down-scaling ratios ($s = 8$ or $s = 14$). The $k$-NN classifier is trained on the merged train/valid set and evaluated on the test set (see [30]), hence the high performances.

| | ViT | | | | CNN | | | |
|---|---|---|---|---|---|---|---|---|
| | Teacher | | Student | | Teacher | | Student | |
| $k$ | $s=14$ | $s=8$ | $s=14$ | $s=8$ | $s=14$ | $s=8$ | $s=14$ | $s=8$ |
| 1 | 52.5 | 54.3 | 51.6 | **55.0** | 45.2 | 45.5 | 45.4 | 44.7 |
| 5 | 55.2 | **56.1** | 55.4 | 55.8 | 47.1 | 47.6 | 46.6 | 46.2 |
| 10 | 57.2 | 56.4 | **57.5** | 56.7 | 49.3 | 46.5 | 50.5 | 45.8 |
| 20 | 56.9 | 58.0 | **58.1** | 57.6 | 47.1 | 47.6 | 45.9 | 47.0 |
| 50 | 56.2 | **57.5** | 55.7 | 56.9 | 41.2 | 44.9 | 40.6 | 44.9 |
| 100 | 53.9 | 54.0 | **54.3** | 53.7 | 40.3 | 43.5 | 40.1 | 44.2 |

different dimension than those used during the pre-training (BACH). Notably, there seems to be a significant performance discrepancy between the models using a [CLS] token (CPVT) and those based on a global average pooling (CPVT-GAP). Based on these results, we select the CPVT-GAP approach for the remaining experiments. Note that we referred to [CLS] token throughout this text when referring to a GAP token. Additionally, we have slightly modified the method to be able to extract one self-attention map per transformer head: instead of performing the GAP operation after the very last layer of the transformer encoder, we do it after the $(L - 1)^{th}$ layer and concatenate the resulting token to the sequence, thereby producing a pseudo [CLS] token. Similarly, when $m$ pseudo [CLS] tokens are used, this operation is performed after the $(L - m)^{th}$ layer.

**Selecting the Number of Finely Attended Regions.** The effect of the number of selected regions is depicted in Table 10. One can observe that it does not appear as the most determining factor, particularly that the results are not monotonically increasing, which is unexpected. There are two potential explanations for this behavior. The first is due to the heterogeneity of the BRACS dataset. More precisely, it encompasses images containing less than 50 patches, which implies that the image must first be resized, potentially harming the predictions. The second explanation is that the model used for this ablation is a $\mathrm{ScoreNet}/4/1$ variant, which by design relies less on the high-resolution images

**Table 7: A weighted $k$ Nearest Neighbors classifier assesses the discriminability (weighted F1-score) of the learned features on the tile CAMELYON16 dataset** (see Sec. B.3). The performances of CNN and ViT-based architectures is compared, and similarly for two tile dimensions ($128 \times 128$ and $224 \times 224$) corresponding to down-scaling ratios of $s = 8$ and $s = 14$, respectively. A 4-fold approach with $75\%/25\%$ train/test splits is used.

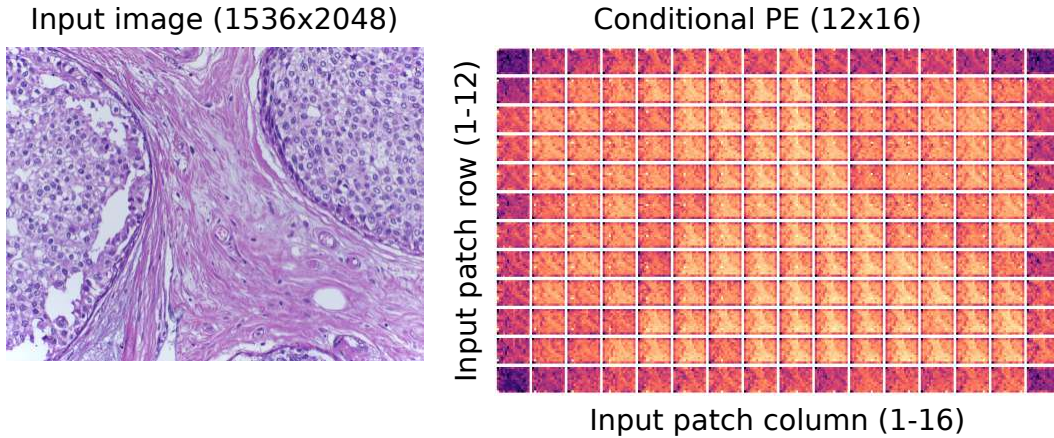| | ViT | | | | CNN | | | |
| | Teacher | | Student | | Teacher | | Student | |
| $k$ | $s = 14$ | $s = 8$ | $s = 14$ | $s = 8$ | $s = 14$ | $s = 8$ | $s = 14$ | $s = 8$ |
|---|---|---|---|---|---|---|---|---|
| 1 | $89.7 \pm 0.6$ | $89.1 \pm 0.4$ | $\mathbf{89.6 \pm 0.6}$ | $89.1 \pm 0.4$ | $87.0 \pm 0.8$ | $85.8 \pm 1.1$ | $87.2 \pm 0.9$ | $85.8 \pm 0.8$ |
| 5 | $\mathbf{91.7 \pm 0.4}$ | $91.1 \pm 0.5$ | $91.6 \pm 0.3$ | $91.2 \pm 0.6$ | $89.9 \pm 1.7$ | $88.8 \pm 1.8$ | $89.8 \pm 1.7$ | $88.7 \pm 1.7$ |
| 10 | $91.9 \pm 0.5$ | $91.4 \pm 0.5$ | $\mathbf{91.9 \pm 0.5}$ | $91.5 \pm 0.4$ | $90.3 \pm 1.0$ | $89.0 \pm 0.5$ | $90.2 \pm 1.1$ | $89.0 \pm 0.6$ |
| 20 | $\mathbf{91.6 \pm 0.6}$ | $91.2 \pm 0.3$ | $91.6 \pm 0.4$ | $91.2 \pm 0.4$ | $90.0 \pm 1.1$ | $89.0 \pm 0.5$ | $89.8 \pm 1.1$ | $88.9 \pm 0.6$ |
| 50 | $\mathbf{91.4 \pm 0.9}$ | $90.7 \pm 0.6$ | $91.3 \pm 1.0$ | $90.7 \pm 0.5$ | $88.8 \pm 1.1$ | $88.6 \pm 0.8$ | $88.9 \pm 1.1$ | $88.5 \pm 0.8$ |
| 100 | $\mathbf{90.9 \pm 1.1}$ | $90.1 \pm 0.6$ | $\mathbf{90.9 \pm 1.1}$ | $90.0 \pm 0.5$ | $88.2 \pm 1.0$ | $87.6 \pm 1.0$ | $88.1 \pm 1.0$ | $87.6 \pm 0.9$ |



**Figure 6: The conditional position encoding [8] of a non-squared input image** is represented. The PE is image-dependent and captures the local interactions between tokens.

than its $\mathrm{ScoreNet}/4/3$ counterpart. The respective properties of these two variants are subject to Sec. C.1.

## C.1. ScoreNet Under the Magnifying Glass

**Just a Glorified Low-resolution ViT?** We explore the usage of high-resolution images for predictions. For that purpose, at test time, we mask $75\%$ of the selected high-resolution regions and report the obtained results in Table 11. As expected, we observe that the $\mathrm{ScoreNet}/4/3$ variant uses the high-resolution content more. Furthermore, these results shed light on how the high-resolution information is not equally relevant for each class. An interesting observation is that for each variant of $\mathrm{ScoreNet}$, the higher the performance of a given model is, the more it is affected by the removal of the high-resolution information (see Table 12).

**Table 12: The performance drop incurred by the high-resolution masking operation of individual models is monitored**. The models that rely the most on the high-resolution content are the ones that perform the best.

| ScoreNet/4/1 | | | ScoreNet/4/3 | | |
|---|---|---|---|---|---|
| 63.3 | $\xrightarrow{-0.6}$ | $62.7 \pm 0.2$ | 63.3 | $\xrightarrow{-2.8}$ | $60.5 \pm 0.1$ |
| 63.8 | $\xrightarrow{-2.2}$ | $61.6 \pm 0.1$ | 64.8 | $\xrightarrow{-5.2}$ | $59.6 \pm 0.3$ |
| 64.9 | $\xrightarrow{-2.2}$ | $62.7 \pm 0.3$ | 65.0 | $\xrightarrow{-6.4}$ | $58.6 \pm 0.3$ |

Despite that, we expected a more considerable drop in performance from this masking operation, which raises the question; *is* $\mathrm{ScoreNet}$ *nothing but a glorified low-resolution ViT?* To answer that question, we train the same ViT as the one used in the recommendation stage and the same setting, but basing the predictions on the scorer's [CLS] tokens and hence without the feedback from the high-resolution stage. Table 13 clearly shows a gap of almost $10\%$ compared to

**Table 8: A weighted $k$ Nearest Neighbors classifier assesses the learned features' discriminability (weighted F1-score) on the low-resolution BACH and BRACS datasets**. A fixed and absolute PE (2D sine-cos)'s performances are compared to a learnable and conditional PE (CPVT and CPVT-GAP). The $k$-NN classifier is trained on the merged train/valid set and evaluated on the test set (BRACS), and a 4-fold approach with $75\%/25\%$ train/test splits is used for BACH dataset.

| | BACH | | | | | | BRACS | | | | | |
| | 2D sine-cos | | CPVT | | CPVT-GAP | | 2D sine-cos | | CPVT | | CPVT-GAP | |
| $k$ | Teacher | Student | Teacher | Student | Teacher | Student | Teacher | Student | Teacher | Student | Teacher | Student |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | $76.0 \pm 3.4$ | $75.0 \pm 4.0$ | $76.6 \pm 2.9$ | $77.7 \pm 2.0$ | $\mathbf{78.5 \pm 6.4}$ | $77.4 \pm 5.1$ | $42.2$ | $42.3$ | $49.6$ | $49.2$ | $54.3$ | $\mathbf{55.0}$ |
| 5 | $74.6 \pm 4.2$ | $75.6 \pm 4.2$ | $76.8 \pm 3.0$ | $76.3 \pm 3.7$ | $\mathbf{81.7 \pm 3.2}$ | $81.0 \pm 4.0$ | $45.3$ | $45.7$ | $53.3$ | $53.2$ | $\mathbf{56.1}$ | $55.8$ |
| 10 | $76.3 \pm 4.1$ | $75.6 \pm 4.6$ | $76.3 \pm 5.0$ | $76.0 \pm 5.2$ | $77.8 \pm 2.8$ | $\mathbf{77.9 \pm 3.4}$ | $47.2$ | $46.3$ | $54.3$ | $54.5$ | $56.4$ | $\mathbf{56.7}$ |
| 20 | $73.9 \pm 3.5$ | $73.9 \pm 3.5$ | $75.7 \pm 5.3$ | $72.9 \pm 5.8$ | $76.3 \pm 3.0$ | $\mathbf{76.5 \pm 4.0}$ | $48.2$ | $47.6$ | $53.3$ | $51.5$ | $\mathbf{58.0}$ | $57.6$ |
| 50 | $73.5 \pm 4.3$ | $73.0 \pm 4.1$ | $74.2 \pm 5.1$ | $73.4 \pm 6.5$ | $\mathbf{74.7 \pm 4.7}$ | $74.3 \pm 5.7$ | $47.0$ | $47.3$ | $50.8$ | $49.7$ | $\mathbf{57.5}$ | $56.9$ |
| 100 | $72.8 \pm 3.7$ | $73.0 \pm 3.1$ | $73.6 \pm 5.8$ | $71.4 \pm 7.4$ | $\mathbf{74.0 \pm 5.5}$ | $73.6 \pm 5.9$ | $45.5$ | $45.0$ | $48.4$ | $48.1$ | $\mathbf{54.0}$ | $53.7$ |

**Table 9: A weighted $k$ Nearest Neighbors classifier assesses the discriminability (weighted F1-score) of the learned features** on the tile CAMELYON16 dataset (see Sec. B.3). A fixed and absolute PE (2D sine-cos)'s performances are compared to a learnable and conditional PE (CPVT and CPVT-GAP). A 4-fold approach with $75\%/25\%$ train/test splits is used.

| | 2D sine-cos | | CPVT | | CPVT-GAP | |
| $k$ | Teacher | Student | Teacher | Student | Teacher | Student |
|---|---|---|---|---|---|---|
| 1 | $88.2 \pm 0.9$ | $88.4 \pm 0.6$ | $88.8 \pm 0.4$ | $88.8 \pm 0.4$ | $\mathbf{89.1 \pm 0.4}$ | $89.1 \pm 0.4$ |
| 5 | $91.1 \pm 0.9$ | $91.0 \pm 1.0$ | $90.9 \pm 0.5$ | $90.9 \pm 0.6$ | $91.1 \pm 0.5$ | $\mathbf{91.2 \pm 0.6}$ |
| 10 | $91.2 \pm 0.7$ | $91.2 \pm 0.5$ | $91.1 \pm 0.4$ | $91.1 \pm 0.3$ | $91.4 \pm 0.5$ | $\mathbf{91.5 \pm 0.4}$ |
| 20 | $91.1 \pm 0.7$ | $91.1 \pm 0.6$ | $91.3 \pm 0.5$ | $\mathbf{91.4 \pm 0.6}$ | $91.2 \pm 0.3$ | $91.2 \pm 0.4$ |
| 50 | $90.5 \pm 0.7$ | $90.6 \pm 0.7$ | $\mathbf{90.8 \pm 0.6}$ | $\mathbf{90.8 \pm 0.5}$ | $90.7 \pm 0.6$ | $90.7 \pm 0.5$ |
| 100 | $\mathbf{90.2 \pm 1.8}$ | $\mathbf{90.2 \pm 0.8}$ | $\mathbf{90.2 \pm 0.6}$ | $\mathbf{90.2 \pm 0.6}$ | $90.1 \pm 0.6$ | $90.0 \pm 0.5$ |

ScoreNet's results and, more interestingly, a gap of more than $5\%$ when compared to the same ViT, but trained with the high-resolution feedback. The above results indicate that **high-resolution information distillation occurs during the training of ScoreNet**.

# D. Computational Cost

Vision transformers heavily rely on the attention mechanism to learn a high-level representation from low-level regions. The underlying assumption is that the different sub-regions of the image are not equally important for the overall representation. Despite this key observation, the computation cost dedicated to a sub-region is independent of its contribution to the high-level representation, which is inefficient and undesirable. Our ScoreNet attention mechanism overcomes this drawback by learning to attribute more computational resources to regions of high interest. Let us consider a high-resolution input image $x_h \in \mathbb{R}^{C \times H \times W}$, a low-resolution version of the image $x_l \in \mathbb{R}^{C \times h \times w}$ is obtained by applying a down-scaling factor $s$, as $h = H/s$ and $w = W/s$. The low-resolution image is fed to a scorer model (recommendation stage), which recommends the regions where to apply fine-grained attention. If this operation is implemented by a ViT, its computational cost is

$\mathcal{O}\left( \left( \frac{h}{P_l} \cdot \frac{w}{P_l} \right)^2 \right)$ with $P_l$ is the dimension of the patches in low-resolution. Using a ViT as the scorer model, there is a one-to-one mapping between the low-resolution patches and the regions the model can process with fine-grained attention; as a consequence, the dimension of the regions is $P_h = s \cdot P_l$. Attending to such regions with a patch size, $P_a$, has a computational cost of $\mathcal{O}\left( \left( \frac{P_h}{P_a} \cdot \frac{P_h}{P_a} \right)^2 \right)$ and the model processes $k$ of them, hence $\mathcal{O}\left( k \cdot \left( \frac{P_h}{P_a} \cdot \frac{P_h}{P_a} \right)^2 \right)$. Finally, a coarse attention mechanism is applied to endow the locally attended regions with contextual information. This final step costs $\mathcal{O}\left( k^2 \right)$. On the other hand, a vanilla ViT would attend uniformly across the whole image with a cost of $\mathcal{O}\left( \left( \frac{H}{P_a} \cdot \frac{W}{P_a} \right)^2 \right)$. Importantly, we observe that only the recommendation stage's cost depends on the input size; consequently, if this step is implemented as a ViT and with a down-scaling ratio $s \in [8, 14]$, the asymptotic cost is reduced by approximately two orders of magnitude, as we typically used $P_a = P_l$ in practice. At last, one can observe that the asymptotic cost can be made linear w.r.t. the input dimension by adopting a convolution-based architecture for the recommendation stage.

# E. ScoreMix Investigation & Examples

The underlying assumption of the "cut-and-paste"-based augmentation methods is that the trained model can assimilate the pasted region to the representation of the target image. In the case of ScoreNet, it translates to attending to the pasted area in a low or high-resolution image. Fig. 7 depicts an example of ScoreNet being able to detect and localize the pasted regions even when the pasted region is small and hard to distinguish. We further observe that a local change in the image directly affects the global representations as the representation of each token is adapted to accommodate the local change in information. This behavior would typically

**Table 10: The number of finely attended regions is selected** by independently training our pipeline 5 times on 10% of the BRACS dataset with a varying number of proposal regions. The number of training epochs is fixed and is the same for all experiments. The models are trained with standard data augmentation methods, i.e., none of ScoreMix, SaliencyMix, or CutMix.

| # Regions | Normal | Benign | UDH | ADH | FEA | DCIS | Invasive | Weighted F1 |
|---|---|---|---|---|---|---|---|---|
| $k = 5$ | **53.7 ± 5.2** | **44.0 ± 5.1** | 29.7 ± 5.3 | 28.8 ± 6.8 | 69.3 ± 4.2 | 56.9 ± 6.5 | **86.9 ± 3.2** | 54.2 ± 1.8 |
| $k = 10$ | 52.1 ± 6.2 | **44.0 ± 3.9** | 31.0 ± 5.3 | 28.6 ± 4.3 | 69.8 ± 3.6 | 56.4 ± 3.9 | 85.9 ± 1.4 | 54.0 ± 0.8 |
| $k = 20$ | 52.2 ± 3.4 | 42.2 ± 5.6 | 29.6 ± 7.5 | **31.9 ± 5.3** | **71.9 ± 2.3** | **57.5 ± 3.6** | **86.9 ± 2.5** | **54.7 ± 0.8** |
| $k = 50$ | 51.5 ± 5.4 | 42.8 ± 4.7 | 30.0 ± 6.8 | 25.9 ± 7.1 | 70.5 ± 4.0 | 55.8 ± 5.2 | 85.7 ± 0.9 | 53.3 ± 2.5 |

**Table 11: At test time, 75% of the selected high-resolution regions are randomly masked.** ScoreNet/4/1 and ScoreNet/4/3 do not equally rely on the high-resolution content.

| Masking | Normal | Benign | UDH | ADH | FEA | DCIS | Invasive | Weighted F1 |
|---|---|---|---|---|---|---|---|---|
| ScoreNet/4/1 | 64.6 ± 2.2 | 52.6 ± 2.8 | 48.4 ± 2.2 | 47.4 ± 2.4 | 77.9 ± 0.7 | 59.3 ± 1.1 | 90.6 ± 1.5 | 64.1 ± 0.7 |
| Masked ScoreNet/4/1 | 61.1 ± 2.7 | 50.8 ± 1.4 | 45.9 ± 2.2 | 41.0 ± 3.5 | 78.8 ± 0.5 | 59.9 ± 3.3 | 90.6 ± 1.1 | 62.4 ± 0.6 |
| ScoreNet/4/3 | 64.3 ± 1.5 | 54.0 ± 2.2 | 45.3 ± 3.4 | 46.7 ± 1.0 | 78.1 ± 2.8 | 62.9 ± 2.0 | 91.0 ± 1.4 | 64.4 ± 0.9 |
| Masked ScoreNet/4/3 | 64.9 ± 2.4 | 51.7 ± 0.5 | 44.4 ± 4.0 | 22.0 ± 6.2 | 77.6 ± 1.0 | 60.8 ± 1.6 | 87.2 ± 1.3 | 59.6 ± 0.7 |

not be observed in a CNN-based architecture until the very last layers. Fig. 7 further highlights the ability of ScoreMix to treat images of different dimensions and aspect ratios.

## F. Learning From Uncurated Data.

We gauge the ability of ScoreNet to learn from unlabeled data on the BACH dataset [1], which encompasses both a small set of 400 annotated TRoIs images, and the WSIs containing the aforementioned TRoIs. Our model is first pre-trained using DINO's self-supervised learning scheme [5] on an unlabeled set of $\approx 11k$ images extracted from WSIs and then is evaluated on the labeled image set using standard protocols, namely linear probing and $k$-NN (see Table 14). We also report the non-empty cluster's purity for the clusters learned by DINO. This metric indicates the quality of a cluster containing samples from a single class. Learning from large uncurated images is particularly challenging, as the increased receptive field allows for the representation of more complex tissue interactions. This further deviates from the discriminative pretext task's assumption that the images represent a single centered object. Since the DINO method enforces a local-to-global correspondence between large and smaller image crops, it may enforce similarity between different tissue types. For that purpose, we modify DINO's pretext task so that the student network only processes the highly discriminative patches to match the teacher's representation, allowing the processing of all the high-resolution patches. To ensure that the pretext task does not encourage contextual biases [35], we only employ augmentations that change the image pixels' values, but not their locations, such that the semantic content of the two augmented views is identical. As can be observed in Table 14, this proposed strategy yields significant improvements compared to other baselines.

**Table 14: Comparison with the prior art for learning capabilities from uncurated data** on the BACH dataset using DINO's pre-training. A comparison results between the effectiveness of DINO's standard pretext task (ScoreNet) and the proposed unbiased pretext task (ScoreNet$^\dagger$) are also reported.

| | ScoreNet$^\dagger$ | ScoreNet | TransPath [43] | SwinTransformer [24] |
|---|---|---|---|---|
| $k$-NN | **73.7 ± 1.7** | 65.0 ± 3.7 | 65.2 ± 1.4 | 63.7 ± 4.1 |
| Lin. eval | **73.0 ± 2.9** | 66.0 ± 2.6 | 64.2 ± 4.0 | 62.5 ± 1.7 |
| Purity | **78.3 ± 23.9** | 76.4 ± 24.9 | 74.0 ± 23.3 | 71.8 ± 23.9 |

**Table 13:** **The ViT network of recommendation stage is trained without receiving any feedback from the high-resolution-based predictions.** Its features discriminability is significantly worth than that of the same model but trained jointly with the high-resolution stage.

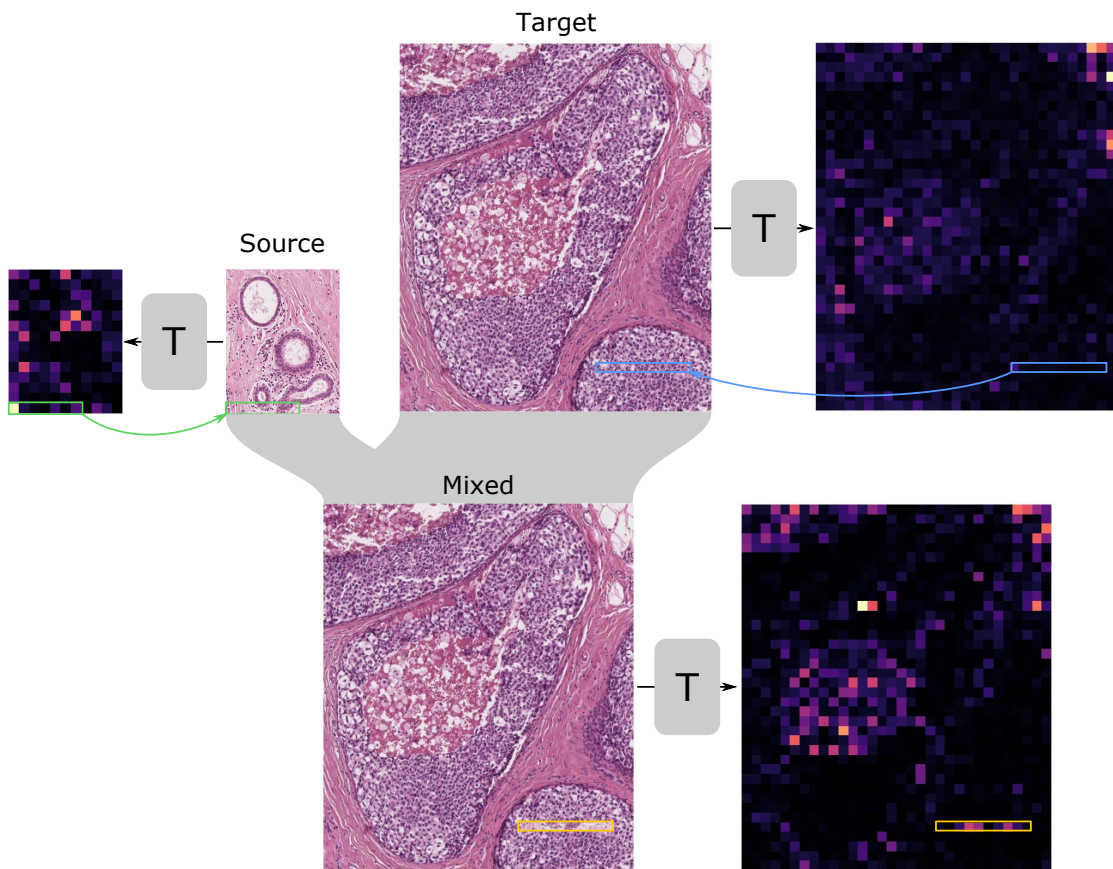| Model | Normal | Benign | UDH | ADH | FEA | DCIS | Invasive | Weighted F1 |
|---|---|---|---|---|---|---|---|---|
| ViT | $53.3 \pm 2.8$ | $42.8 \pm 1.9$ | $37.1 \pm 2.9$ | $32.4 \pm 2.4$ | $77.3 \pm 0.2$ | $51.2 \pm 1.3$ | $85.0 \pm 1.8$ | $55.5 \pm 0.1$ |
| Lin. scorer's [CLS] | $57.5 \pm 4.2$ | $48.8 \pm 5.5$ | $42.7 \pm 3.5$ | $42.7 \pm 7.4$ | $74.3 \pm 5.2$ | $60.5 \pm 2.4$ | $90.6 \pm 0.2$ | $60.9 \pm 3.1$ |
| ScoreNet/4/1 | $\mathbf{64.6 \pm 2.2}$ | $52.6 \pm 2.8$ | $\mathbf{48.4 \pm 2.2}$ | $\mathbf{47.4 \pm 2.4}$ | $77.9 \pm 0.7$ | $59.3 \pm 1.1$ | $90.6 \pm 1.5$ | $64.1 \pm 0.7$ |
| ScoreNet/4/3 | $64.3 \pm 1.5$ | $\mathbf{54.0 \pm 2.2}$ | $45.3 \pm 3.4$ | $46.7 \pm 1.0$ | $\mathbf{78.1 \pm 2.8}$ | $\mathbf{62.9 \pm 2.0}$ | $\mathbf{91.0 \pm 1.4}$ | $\mathbf{64.4 \pm 0.9}$ |



**Figure 7: The learned semantic distribution can detect and localize the newly pasted content.** The green box highlights the region pasted from the source to the target image. The blue box represents the region where the new content is pasted. The yellow box highlights the modified region in the mixed image. $T$ represents the scorer network of ScoreNet.