# How Well Do Vision Transformers (VTs) Transfer To The Non-Natural Image Domain? An Empirical Study Involving Art Classification

Vincent Tonkes[1] and Matthia Sabatelli[1]

Department of Artificial Intelligence and Cognitive Engineering, University of Groningen, 9712 CP Groningen, The Netherlands
{m.sabatelli}@rug.nl

**Abstract.** Vision Transformers (VTs) are becoming a valuable alternative to Convolutional Neural Networks (CNNs) when it comes to problems involving high-dimensional and spatially organized inputs such as images. However, their Transfer Learning (TL) properties are not yet well studied, and it is not fully known whether these neural architectures can transfer across different domains as well as CNNs. In this paper we study whether VTs that are pre-trained on the popular ImageNet dataset learn representations that are transferable to the non-natural image domain. To do so we consider three well-studied art classification problems and use them as a surrogate for studying the TL potential of four popular VTs. Their performance is extensively compared against that of four common CNNs across several TL experiments. Our results show that VTs exhibit strong generalization properties and that these networks are more powerful feature extractors than CNNs.

**Keywords:** Vision Transformers, Convolutional Neural Networks, Transfer Learning, Art Classification

## 1 Introduction

Since the introduction of AlexNet, roughly a decade ago, Convolutional Neural Networks (CNNs) have played a significant role in Computer Vision (CV) [15]. Such neural networks are particularly well-tailored for vision-related tasks, given that they incorporate several inductive biases that help them deal with high dimensional, rich input representations. As a result, CNNs have found applications across a large variety of domains that are not *per-se* restricted to the realm of natural images. Among such domains, the Digital Humanities (DH) field is of particular interest. Thanks to a long tradition of works that aimed to integrate advances stemming from technical disciplines into the Humanities, they have been serving as a challenging real-world test-bed regarding the applicability of CV algorithms. It naturally follows that over the last years, several works have studied the potential of CNNs within the DH (see [11] for a survey about the topic), resulting in a significant number of successful applications that range from the classification of artworks [36,30,43,23] to the detection of

objects within paintings [12,29], automatic style classification [7] and even art understanding [2].

A major breakthrough within the CV community has recently been achieved by the Vision Transformer (VT) [9], a novel neural architecture that has gained state-of-the-art performance on many standard learning benchmarks including the popular ImageNet dataset [8]. Exciting as this may be, we believe that a plain VT is not likely to become as valuable and powerful as a CNN as long as it does not exhibit the strong generalization properties that, over the years, have allowed CNNs to be applied across almost all domains of science [39,25,1]. Therefore, this paper investigates what VTs offer within the DH by studying this family of neural networks from a Transfer Learning (TL) perspective. Building on top of the significant efforts that the DH have been putting into digitizing artistic collections from all over the world [33], we define a set of art classification problems that allow us to study whether pre-trained VTs can be used outside the domain of natural images. We compare their performance to that of CNNs, which are well-known to perform well in this setting, and present to the best of our knowledge the very first thorough empirical analysis that describes the performance of VTs outside the domain of natural images and within the domain of art specifically.

## 2    Preliminaries

We start by introducing some preliminary background that will be used throughout the rest of this work. We give an introduction about supervised learning and transfer learning (Sec. 2.1), and then move towards presenting some works that have studied Convolutional Neural Networks and Vision Transformers from a transfer learning perspective.

### 2.1   Transfer Learning

With Transfer Learning (TL), we typically denote the ability that machine learning models have to retain and reuse already learned knowledge when facing new, possibly related tasks [26,46]. While TL can present itself within the entire machine learning realm [4,37,45,42], in this paper we consider the supervised learning setting only, a learning paradigm that is typically defined by an input space $\mathcal{X}$, an output space $\mathcal{Y}$, a joint probability distribution $P(X,Y)$ and a loss function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$. The goal of a supervised learning algorithm is to learn a function $f : \mathcal{X} \to \mathcal{Y}$ that minimizes the expectation over $P(X,Y)$ of $\ell$ known as the expected risk. Classically, the only information that is available for minimizing the expected risk is a learning set $\mathcal{L}$ that provides the learning algorithm with $N$ pairs of input vectors and output values $(\mathbf{x}_1, y_1), ..., (\mathbf{x}_N, y_N)$ where $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \mathcal{Y}$ are i.i.d. drawn from $P(X,Y)$. Such learning set can then be used for computing the empirical risk, an estimate of the expected risk, that can be used for finding a good approximation of the optimal function $f^*$ that minimizes the aforementioned expectation. When it comes to TL, however,

we assume that next to the information contained within $\mathcal{L}$, the learning algorithm also has access to an additional learning set defined as $\mathcal{L}'$. Such a learning set can then be used alongside $\mathcal{L}$ for finding a function that better minimizes the loss function $\ell$. In this work, we consider $\mathcal{L}'$ to be the ImageNet dataset, and $f$ to come in the form of either a pre-trained Convolutional Neural Networks or a pre-trained Vision Transformers, two types of neural networks that over the years have demonstrated exceptional abilities in tackling problems modeled by high dimensional and spatially organized inputs such as images, videos, and text.

## 2.2   Related Works

While to this date, countless examples have studied the TL potential of CNNs [25,40,34,24], the same cannot yet be said for VTs. In fact, papers that have so far investigated their generalization properties are much rarer. Yet, some works have attempted to compare the TL potential of VTs to that of CNNs, albeit strictly outside the artistic domain. E.g., in [21] the authors consider the domain of medical imaging and show that if CNNs are trained from scratch, then these models outperform VTs; however, if either an off-the-shelf feature extraction approach is used, or a self-supervised learning training strategy is followed, then VTs significantly outperform their CNNs counterparts. Similar results were also observed in [44] where the authors show that both in a single-task learning setting and in a multi-task learning one, transformer-based backbones outperformed regular CNNs on 13 tasks out of 15. Positive TL results were also observed in [16], where the TL potential of VTs is studied for object detection tasks, in [41] where the task of facial expression recognition is considered, and in [10], where similarly to [21], the authors successfully transfer VTs to the medical imaging domain. The generalization properties of VTs have also been studied outside of the supervised learning framework: E.g. in [5] and [6], a self-supervised learning setting is considered. In [5] the authors show that VTs can learn features that are particularly general and well suited for TL, which is a result that is confirmed in [6], where the authors show that VTs pre-trained in a self-supervised learning manner transfer better than the ones trained in a pure supervised learning fashion. While all these works are certainly valuable, it is worth noting that, except for [21] and [10], all other studies have performed TL strictly within the domain of natural images. Therefore, the following research question

*"How transferable are the representations of pre-trained VTs when it comes to the non-natural image domain?"*

remains open. We will now present our work that helps us answer this question by considering **image datasets from the artistic domain.** Three art classification tasks are described, which serve as a surrogate for identifying which family of pre-trained networks, among CNNs and VTs, transfers better to the non-natural realm.
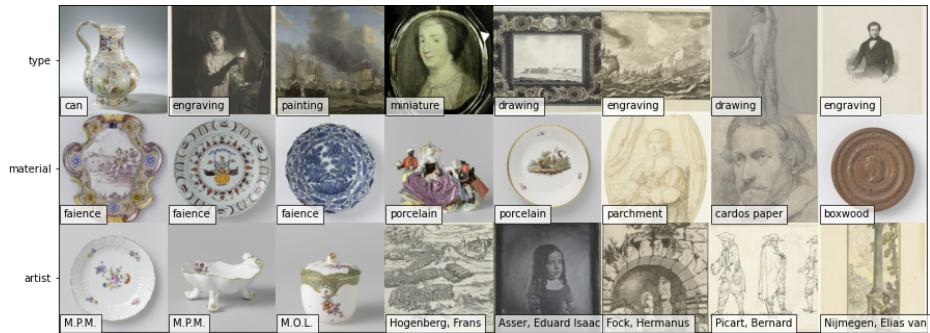
Fig. 1: Samples from the three distinct classification tasks used throughout this paper.

## 3    Methods

Our experimental setup is primarily inspired by the work presented in Sabatelli et al. [30], where the authors report a thorough empirical analysis that studies the TL properties of pre-trained CNNs that are transferred to different art collections. More specifically, the authors investigate whether popular neural architectures such as `VGG19` [32] and `ResNet50` [13], which come as pre-trained on the ImageNet1k [8] dataset, can get successfully transferred to the artistic domain. The authors consider three different classification problems and two different TL approaches: an off-the-shelf (OTS) feature extraction approach where pre-trained models are used as pure feature extractors and a fine-tuning approach (FT) where the pre-trained networks are allowed to adapt all of their pre-trained parameters throughout the training process. Their study suggests that all the considered CNNs can successfully get transferred to the artistic domain and that a fine-tuning training strategy yields substantially better performance than the OTS one. In this work we investigate whether these conclusions also hold for transformer-based architectures, and if so, whether this family of models outperforms that of CNNs.

### 3.1    Data

Similar to Sabatelli et al. [30] we also use data stemming from the *Rijksmuseum Challenge* dataset [22] (see Fig. 1 for an illustration). This dataset consists of a large collection of digitized artworks that come together with `xml`-formatted metadata which can be used for defining a set of supervised learning problems [22,33]. Following [22] and [30] we focus on three, well-known, classification problems namely: (1) *Type classification*, where the goal is to train a model such that it is able to distinguish classes such as 'painting', 'sculpture', 'drawing', etc.; (2) *Material classification*, where the classification problem is defined by labels such as 'paper', 'porcelain', 'silver', etc.; and finally, (3) *Artist classification*, where, naturally, the model has to predict who the creator of a specific artwork is. While

Table 1: Overview of the used datasets. Values between brackets show the situation before balancing operations were performed. 'Sample overlap' gives the average overlap between 2 of the 5 randomly generated sets per task ($i$ and $j$ where $i \neq j$).

| Task | # Samples | # Classes | Sample overlap |
|------|-----------|-----------|----------------|
| Type classification | 9607 (77628) | 30 (801) | 0.686 |
| Material classification | 7788 (96583) | 30 (136) | 0.798 |
| Artist classification | 6530 (38296) | 30 (8592) | 1 |

the full dataset contains 112,039 images, due to computational reasons in the present study we only use a fraction of it as this allows us to run shorter yet more thorough experiments. A smaller dataset also allows us to study an additional research question which we find worth exploring and that [30] did not consider in their study, namely: *"How well do CNNs (and VTs) transfer when the size of the artistic collection is small?"* To this end we decided to select the 30 most occurring classes within their dataset and to set a cap of 1000 randomly sampled instances per class. Table 1 summarizes the datasets used in the present study. Between brackets we report values as they were before balancing operations were performed. For all of our experiments we use 5 fold cross-validation and use 80% of the dataset as training-set, 10% as validation-set and 10% as testing-set.

## 3.2   Neural Architectures

In total, we consider eight different neural architectures, of which four are CNN-based networks while the remaining four are VT-based models. All models are pre-trained on the ImageNet1K dataset. When it comes to CNNs we consider `ResNet50` [13] and `VGG19` [32], as the first one has widely been adopted by researchers working at the intersection of computer vision and digital heritage [43,2], whereas the latter was among the best performing models considered by Sabatelli et al. [30]. Next to these two architectures we also consider two additional, arguably more recent networks, namely `ConvNext` [19], which is a purely CNN-based model that is inspired by VTs' recent successes, and that therefore fits well within the scope of this work, and `EfficientNetV2` [35], which is a network that is well known for its computational efficiency and potentially faster training times. Regarding the VTs, we use models that have $16 \times 16$ patch sizes. As a result, we start by considering the first original VT model presented in [9] which we refer to as `ViT`. We then consider the `Swin` architecture [18] as it showed promising TL performance in [44] and the `BeiT` [3] and the `DeiT` [38] transformers. The usage of `DeiT` is motivated by Matsoukas et al. work [21] reviewed in Section 2.2 where the authors compared it to `ResNet50`. Yet, note that in this work, we use the version of the model known as the 'base' version that does not take advantage of distillation learning in the FT stage. All VT-based models have $\approx 86$ million trainable parameters, and so does `ConvNext`. `ResNet50` and `EfficientNetV2` are, however much smaller networks as they come with $\approx 25.6$

and $\approx$ 13.6 million trainable parameters, respectively. Lastly, `VGG19` is by far the largest model of all with its 143.7 million trainable parameters.

### 3.3   Training Procedure

For all experiments, images are resized to a $224 \times 224$ resolution by first scaling them to the desired size along the shortest axis (retaining aspect ratio), and then taking a center crop along the longer axis. In addition, all images are normalized to the RGB mean and standard deviation of the images presented within the ImageNet1K dataset ([0.485, 0.456, 0.406] and [0.229, 0.224, 0.225]). For all models we replace the final linear classification layer with a new layer with as many output nodes as there are classes to classify ($C$) within the dataset and optimize the parameters of the model $\theta$ such that the categorical cross-entropy loss function

$$\mathscr{L}(\theta) = -\mathbb{E}_{(\mathbf{x},y) \sim P(X,Y)} \sum_{i=1}^{C} 1(y = i) \log p_{\mathrm{model}} f_i(\mathbf{x}; \theta), \qquad (1)$$

is minimized.

Training is regularized through the early stopping method, which interrupts training if for 10 epochs in a row no improvement on the validation loss is observed. The model with the lowest validation loss is then benchmarked on the final testing set. For our OTS experiments, we use the Adam optimizer [14] with standard `PyTorch` [27] parameters (`lr=1e-3`, $\beta_1 = 0.9$, $\beta_2 = 0.999$), and use a batch size of 256. For the FT experiments, hyperparameters are partially inspired by [21,44]. We again use the Adam optimizer, this time initialized with `lr=1e-4` which gets reduced by a factor of 10 after three epochs without improvements on the validation loss; inspired by [20], the batch size is now reduced to 32; whereas label smoothing (0.1) and dropout (p=0.2) are used for regularizing training even further. Finally, input images are augmented with random horizontal flips and rotations in a $\pm 10°$ range.

### 3.4   Hardware and Software

All experiments are conducted on a single compute node containing one 32 GB `Nvidia V100 GPU`. The FT experiments take advantage of the V100's mixed precision acceleration. An exception is made for the *Type Classification* experiment, as this one is also used to compare OTS TL and FT in terms of time/accuracy trade-offs (see Sec. 6 for further details). The `PyTorch` machine learning framework [27] is used for all experiments, and many pre-trained models are taken from its `Torchvision` library. Exceptions are made for `EfficientNetV2`, `Swin`, `DeiT` and `Beit`, which are taken from the `Timm` library[1]. We release all source code and data on the following GitHub link[2].

---

[1] `https://timm.fast.ai/`.
[2] `https://github.com/IndoorAdventurer/ViTTransferLearningForArtClassification`

## 4　Results

We now present the main findings of our study and report results for the three aforementioned classification problems and for all previously introduced architectures. All networks are either trained with an OTS training scheme (Sec. 4.1) or with a FT one (Sec. 4.2). Results come in the form of line plots and tables: the former visualize the performance of all models in terms of accuracy on the validation set, whereas the latter report the final performance that the best validation models obtained on the separate testing-set. All line plots report the average accuracy obtained across five different experiments, while the shaded areas correspond to the standard error of the mean ($\pm s \div \sqrt{N}$, where $N = 5$). The dashed lines represent CNNs, whereas continuous lines depict VTs. Plots end when an early stop occurred for the first of the five trials. Regarding the tables, a green-shaded cell marks the best overall performing model, while yellow and red cells depict the second best and worst performing networks, respectively. We quantitatively assess the performance of the models with two separate metrics: the accuracy and the balanced accuracy, where the latter is defined as the average recall over all classes. Note that compared to the plain accuracy metric, the balanced accuracy allows us to penalize type-II errors more when it comes to the less occurring classes within the dataset.

### 4.1　Off-The-Shelf Learning

For this set of experiments, we start by noting that overall, both the CNNs and the VTs can perform relatively well on all three different classification tasks. When it comes to (a) *Type Classification* we can see that all models achieve a final accuracy between $\approx 80\%$ and $\approx 90\%$, whereas on (b) *Material Classification* the performance deviates between $\approx 75\%$ and $85\%$, and on (c) *Artist Classification* we report accuracies between $75\%$ and $90\%$. Overall, however, as highlighted by the green cells in Table 2, the best performing model on all classification tasks is the VT `Swin`, which confirms the good TL potential that this architecture has and that was already observed in [44]. Yet, we can also observe that the second-best performing model is not a VT, but the `ConvNext` CNN. Despite performing almost equally well, it is important to mention though that `ConvNext` required more training epochs to converge when compared to `Swin`, as can clearly be seen in all three plots reported in Fig. 2. We also note that on the *Type Classification* task the worst performing model is the `Beit` transformer ($82.26\%$), but when it comes to the classification of the materials and artists the worst performing model becomes `EfficientNetV2` with final accuracies of $75.96\%$ and $73.92\%$ respectively. Among the different VTs the `Beit` transformer also appears to be the architecture that requires the longest training. In fact, as can be seen in Fig. 3, this network does not exhibit substantial "Jumpstart-Improvements" as the other VTs (learning starts much lower in all plots as is depicted by the green lines).

　　Several other interesting conclusions can be made from this experiment: we observed that the `VGG19` network yielded worse performance than `ResNet50`,

a result which is not in line with what was observed by Sabatelli et al. [30] where `VGG19` was the best performing network when used with an OTS training scheme. Also, differently from [30], the most challenging classification task in our experiments appeared to be that of *Material Classification* as it resulted in the overall lowest accuracies. These results seem to suggest that even though both studies considered images stemming from the same artistic collection, the training dataset size can significantly affect the TL performance of the different models, a result which is in line with [28].

In general, VTs seem to be very well suited for an OTS transfer learning approach, especially regarding the `Swin` and `DeiT` architectures which performed well across all tasks. Equally interesting is the performance obtained by `ConvNext` which is by far the most promising CNN-based architecture. Yet, on average, the performance of the VTs is higher than the one obtained by the CNNs: on *Type Classification* the first ones perform on average $\approx 86.5\%$ while the CNNs reach an average classification rate of $\approx 85.51\%$, whereas on *Material* and *Artist Classification* VTs achieve average accuracies of $\approx 81.8\%$ and $\approx 85.8\%$ respectively, whereas CNNs $\approx 79.5\%$ $\approx 81.2\%$. These results suggest that this family of methods is better suited for OTS TL when it comes to art classification problems.
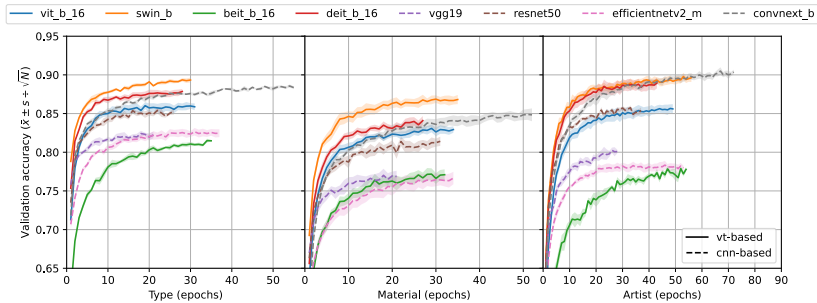


Fig. 2: The validation accuracy obtained by all architectures when trained with an off-the-shelf (OTS) feature extraction approach. We can see that the best performing models are the VT `Swin` and the CNN `ConvNext`.

### 4.2   Fine-Tuning

When it comes to the fine-tuning experiments, we observe, in part, consistent results with what we have presented in the previous section. We can again note that the lower classification rates have been obtained when classifying the material of the different artworks, while the best performance is again achieved when classifying the artists of the heritage objects. The `Swin` VT remains the network that overall performs best, while the `ConvNext` CNN remains the overall second

Table 2: The performance of all models trained with an OTS approach on the final testing sets. We can see that the best overall model is the `Swin` transformer (green cells) followed by the `ConvNext` CNN architecture (yellow cells). The worst performing models are `Beit` when it comes to type classification and `EfficientNetV2` when it comes to material and artist classification.

| Model | Type | | Material | | Artist | |
|---|---|---|---|---|---|---|
| | Accuracy | Bal. accuracy | Accuracy | Bal. accuracy | Accuracy | Bal. accuracy |
| vit_b_16 | 86.06% ±1.06% | 84.13% ±1.57% | 81.78% ±0.48% | 67.38% ±1.37% | 84.89% ±0.46% | 81.42% ±0.42% |
| swin_b | 89.43% ±0.93% | 87.47% ±1.02% | 85.87% ±0.35% | 71.19% ±1.60% | 90.40% ±0.65% | 88.64% ±0.78% |
| beit_b_16 | 82.26% ±0.72% | 77.75% ±0.27% | 76.87% ±0.96% | 60.16% ±1.56% | 79.70% ±0.69% | 75.35% ±1.09% |
| deit_b_16 | 88.18% ±0.66% | 85.36% ±0.54% | 82.80% ±1.12% | 66.46% ±1.03% | 88.13% ±0.76% | 85.62% ±0.87% |
| vgg19 | 83.93% ±0.72% | 83.35% ±0.81% | 76.87% ±0.44% | 61.39% ±1.47% | 82.01% ±0.66% | 78.10% ±0.77% |
| resnet50 | 85.51% ±0.64% | 82.33% ±1.85% | 80.99% ±0.82% | 65.51% ±0.93% | 87.71% ±1.06% | 85.12% ±1.34% |
| eff. netv2_m | 83.41% ±0.76% | 82.05% ±1.25% | 75.96% ±1.24% | 59.15% ±1.24% | 78.62% ±1.07% | 73.92% ±0.96% |
| convnext_b | 89.19% ±0.64% | 86.95% ±1.38% | 84.14% ±0.92% | 69.10% ±1.05% | 90.13% ±0.94% | 87.84% ±1.07% |

best performing model, even though on *Artist Classification* it is slightly outperformed by `ResNet50`. Unlike our OTS experiments, however, this time, we note that VTs are on average the worst-performing networks. When it comes to the *Type Classification* problem, the `ViT` model achieves the lowest accuracy (although note that if the balance accuracy metric is used, then the worst performing network becomes `VGG19`). Further, when considering the classification of materials and artists, the lowest accuracies are instead achieved by the `Beit` model. While it is true that overall the performance of both CNNs and VTs significantly improves through fine-tuning, it is also true that, perhaps surprisingly, such a training approach reduces the differences in terms of performance between these two families of networks. This result was not observed in [44], where VTs outperformed CNNs even in a fine-tuning training regime. While these results indicate that an FT training strategy does not seem to favor either VTs or CNNs, it is still worth noting that this training approach is beneficial. In fact, no matter whether VTs or CNNs are considered, the worst-performing fine-tuned model still performs better than the best performing OTS network. The only exception to this is *Material Classification*, where `BeiT` obtained a testing accuracy of 85.75% after fine-tuning which is slighlty lower than the one obtained by the `Swin` model trained with an OTS training scheme (85.87%).

## 5  Discussion

Our results show that VTs possess strong transfer learning capabilities and that this family of models learns representations that can generalize to the artistic realm. Specifically, as demonstrated by our OTS experiments, when these architectures are used as pure feature extractors, their performance is, on average, substantially better than the one of CNNs. To the best of our knowledge, this is the first work that shows that this is the case for non-natural and artistic images. Next to be attractive to the computer vision community, as some light is shed on the generalization properties of VTs, we believe that these results are
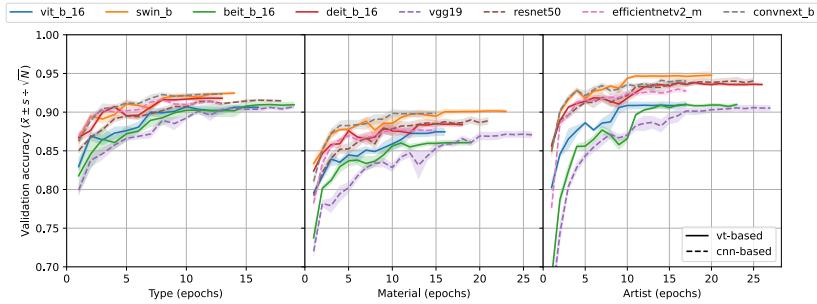
Fig. 3: The validation accuracy learning curves obtained by all models when a fine-tuning FT training approach is followed. We can again see that the `Swin` and `ConvNext` architectures are among the best performing networks.

Table 3: The testing set performance of all models trained with a FT approach. We can see that compared to the results presented in Table 2, all models perform substantially better, yet differences in terms of performance between the VTs and the CNNs now seem to be smaller.

| Model | Type | | Material | | Artist | |
|---|---|---|---|---|---|---|
| | Accuracy | Bal. accuracy | Accuracy | Bal. accuracy | Accuracy | Bal. accuracy |
| **vit_b_16** | 90.11% ±0.35% | 87.40% ±0.29% | 87.42% ±0.45% | 73.46% ±1.44% | 92.05% ±0.44% | 89.77% ±0.38% |
| **swin_b** | 92.17% ±0.98% | 89.71% ±1.03% | 89.35% ±0.68% | 77.16% ±2.98% | 95.05% ±0.47% | 93.94% ±0.81% |
| **beit_b_16** | 90.81% ±0.41% | 87.95% ±0.69% | 85.74% ±0.37% | 72.12% ±1.38% | 91.27% ±1.13% | 88.83% ±1.63% |
| **deit_b_16** | 91.78% ±0.64% | 89.22% ±0.90% | 87.85% ±1.12% | 74.42% ±1.99% | 93.37% ±1.15% | 91.67% ±1.55% |
| **vgg19** | 90.54% ±0.37% | 87.05% ±1.03% | 85.74% ±1.40% | 72.43% ±3.03% | 92.20% ±0.49% | 90.18% ±0.72% |
| **resnet50** | 91.78% ±0.44% | 88.24% ±0.59% | 88.69% ±0.99% | 77.97% ±2.25% | 94.72% ±0.74% | 93.41% ±1.05% |
| **eff. netv2_m** | 90.87% ±0.67% | 88.34% ±1.37% | 87.55% ±1.15% | 75.31% ±1.60% | 92.65% ±0.54% | 90.84% ±0.51% |
| **convnext_b** | 92.15% ±0.40% | 89.82% ±1.18% | 88.79% ±1.07% | 78.40% ±1.26% | 94.60% ±0.54% | 93.13% ±0.61% |

also of particular interest for practitioners working in the digital humanities with limited access to computing power. As the resources for pursuing an FT training approach might not always be available, it is interesting to know that between CNNs and VTs, the latter models are the ones that yield the best results in the OTS training regime.

Equally interesting and novel are the results obtained through fine-tuning, where the performance gap between VTs and CNNs gets greatly reduced, with the latter models performing only slightly worse than the former. In line with the work presented in [30], which considered CNNs exclusively, we clearly show that this TL approach also substantially improves the performance of VTs.

## 6    Additional Studies

We now present four additional studies that we hope can help practitioners that work at the intersection of computer vision and the digital humanities. Specifically we aim to shed some further light into the classification performance

of both CNNs and VTs, while also providing some practical insights that consider the training times of both families of models.

### 6.1  Saliency Maps

We start by performing a qualitative analysis that is based on the visual investigation of saliency maps. For this set of studies we consider the `ConvNext` CNN and the `Swin` and `DeiT` VTs. When it comes to the former, saliency maps are computed through the popular GradCam method [31], whereas attention-rollout is used when it comes to the `DeiT` architecture. `Swin`'s saliency maps are also computed with a method similar to GradCam, with the main difference being that instead of taking an average of the gradients per channel as weights, we directly multiply gradients by their activation and visualize patch means of this product. Motivated by the nice performance of VTs in the OTS transfer learning setting, we start by investigating how the representation of an image changes with respect to the depth of the network. While it is well known that CNNs build up a hierarchical representation of the input throughout the network, similar studies involving transformer-based architectures are rarer. In Fig. 4 we present some examples that show that also within VTs, the deeper the network becomes, the more the network starts focusing on lower level information within the image.

In Fig. 5 we show how different network architectures and TL approaches result in different saliency maps. In the first image of Fig. 5 we can observe that the `ConvNext` architecture miss-classifies a 'Dish' as a 'Plate' when an OTS approach is used. Note that this is a mistake not being made by the transformer-based architectures. However, we observe that after the fine-tuning stage, the CNN can classify the image's *Type* correctly after having shifted its attention toward the bottom and center of the dish rather than its top. Interestingly, none of the transformer-based architecture focuses on the same image regions, independently of whether an OTS or a FT approach is used. We can see that most saliency maps are clustered within the center of the image, both when an OTS training strategy is used and when an FT approach is adopted. When looking at the second image of Fig. 5, we see that all networks, independently from the adopted TL approach, correctly classify the image as a *'Picture'*. Yet the saliency maps of the `ConvNext` change much more across TL approaches, and we again see that the transformer-based architectures focus more on the central regions of the image rather than on the borders. Similar behavior can also be observed in the last image of Fig. 5.

### 6.2  Dealing With Small Artistic Collections

It is not uncommon for heritage institutions to deal with datasets far smaller than those typically used by the computer vision community. While it is true that the number of samples used throughout this study is far smaller than the one used within the naturalistic domain, we still wondered whether the results reported in Sec. 4 would generalize to even smaller artistic collections. Inspired by [17] we
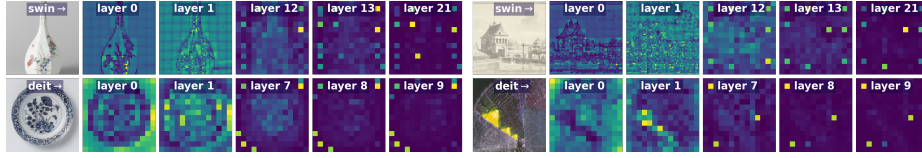
Fig. 4: Saliency maps for different attention layers of successively deeper transformer blocks. The deeper the network, the lower level the representations learned by the Vision Transformer become.
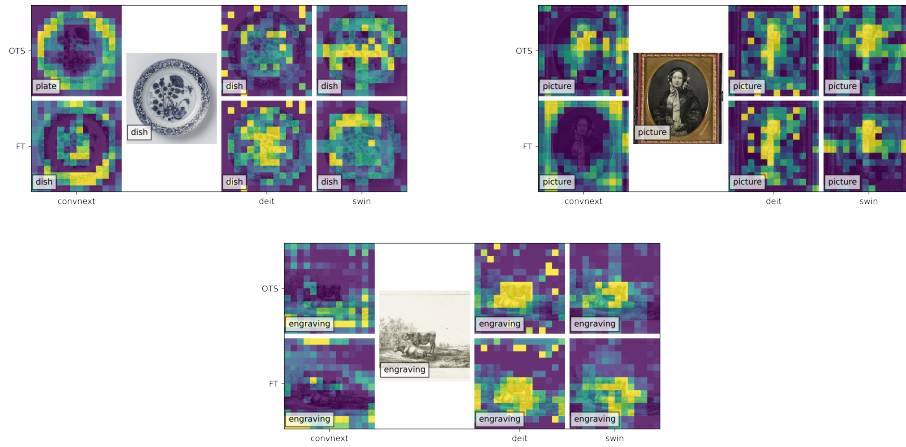


Fig. 5: Saliency maps computed for `ConvNext`, `DeiT` and `Swin` when classifying the *"Type"* of three heritage objects. Saliency maps are computed both for the off-the-shelf experiments as well as for the fine-tuning ones.

have designed the following experiment where we consider the *Type Classification* task. We consider the top 15 most occurring classes within the original dataset presented in Table 1 and scale the number of samples four times by a factor $\sqrt[4]{\frac{1}{10}} \approx 0.56$. This results in five datasets which are respectively 100%, 56%, 32%, 18% and 10% the size of the original dataset shown in Table 1. Note that we ensure that the distribution of samples per class remains the same across all datasets. We then trained all models following the exact experimental protocol described in Sec. 3. Our results are reported in Fig. 6 where we can observe how the testing accuracy decreases when smaller portions of the full dataset are taken. Note that the x-axes show a logarithmic scale, with the rightmost value being roughly 10% the size of the leftmost one. For both an OTS approach (left plot of Fig. 6), as well as a FT one (right plot of Fig. 6), we show that the findings discussed in Sec. 4.1 and Sec. 4.2 generalize to smaller datasets.
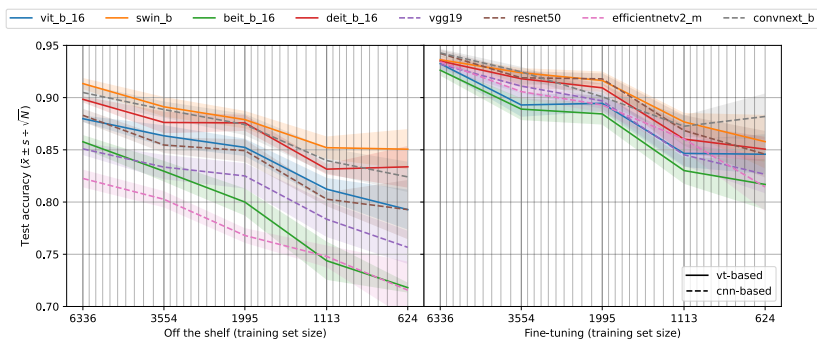
Fig. 6: The testing accuracy obtained on the *Type Classification* problem when gradually reducing the size of the dataset. We can observe that different dataset sizes do not affect the best performing models, which as presented in Sec. 4 remain the `Swin` VT and the `ConvNext` CNN.

## 6.3   Training Times

We now report some final results that describe the training times that are required by CNNs and VTs when trained with the aforementioned TL strategies. These results have also been obtained on the *Type* classification task but with the V100s mixed precision capabilities disabled. In Fig. 7 we show the classification accuracy as reported in Table 2 and 3 and plot it against the average duration of one training epoch. This allows us to understand the time/accuracy trade-offs that characterize all neural networks. As shown by the blue dots, we can observe that when it comes to VTs trained in an OTS setting, all transformer-based architectures require approximately the same number of seconds to successfully go through one training epoch ($\approx 40$). This is, however, not true for the CNNs, which on average, require less time and for which there is a larger difference between the fastest OTS model (`ResNet50` requiring $\approx 20$ seconds) and the slowest model (`ConvNext`). While, as discussed in Sec. 4, VTs are more powerful feature extractors than CNNs, it is worth noting that the gain in performance these models have to offer comes at a computational cost. Note, however, that this is not true anymore when it comes to an FT training regime, as all models (red and purple dots) now perform almost equally well. Yet it is interesting to point out that the computational costs required by the VTs stay approximately the same across all architectures, which is not the case for the CNN-based networks as there is a clear difference between the fastest fine-tuned network (`ResNet50`) and the slowest (`ConvNext`). We believe that the design of a transformer-based architecture that, if fine-tuned, results in the same training times as `ResNet50` provides an interesting avenue for future work.
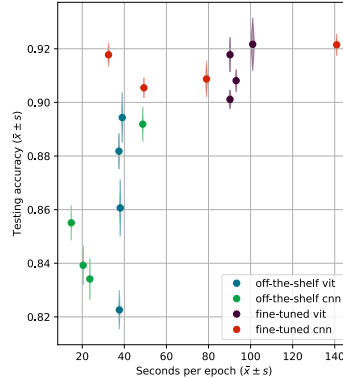
Fig. 7: Time/accuracy trade-offs for VTs and CNNs trained either in an off-the-shelf setting or with a fine-tuning strategy.

## 7   Conclusion

This work examined how well VTs can transfer knowledge to the non-natural image domain. To this end, we compared four popular VT architectures with common CNNs, in terms of how well they transfer from ImageNet1k to classification tasks presented in the *Rijksmuseum Challenge* dataset. We have shown that when fine-tuned VTs performed on par with CNNs and that they performed even better than their CNN counterparts when used as feature extractors. We believe that our study proves that VTs can become a valuable alternative to CNN-based architectures in the non-natural image domain and in the realm of artistic images specifically. Especially `Swin` and `DeiT` showed promising results throughout this study and, therefore, we aim to investigate their potential within the Digital Humanities in the future. To this end, inspired by [33] we plan on using them in a multi-task learning setting; as backbone feature extractors when it comes to the object detection within artworks [12,29], and finally, in a self and semi-supervised learning setting. We also plan on performing a similar analysis for the `ConvNext` CNN, as this architecture showed very good performance as well.

To conclude, we believe that the study presented in this work opens the door for a more fundamental CV question that deserves attention: *"What makes transformer-based architectures such powerful feature extractors?"*. We plan on investigating whether the results obtained within the domain of DH will also generalize to other non-natural image datasets with the hope of partially answering this question.

# References

1. Ackermann, S., Schawinski, K., Zhang, C., Weigel, A.K., Turp, M.D.: Using transfer learning to detect galaxy mergers. Monthly Notices of the Royal Astronomical Society **479**(1), 415–425 (2018)
2. Bai, Z., Nakashima, Y., Garcia, N.: Explain me the painting: Multi-topic knowledgeable art description generation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5422–5432 (2021)
3. Bao, H., Dong, L., Piao, S., Wei, F.: Beit: Bert pre-training of image transformers. In: ICLR 2022 (2022)
4. Bengio, Y.: Deep learning of representations for unsupervised and transfer learning. In: Proceedings of ICML workshop on unsupervised and transfer learning. pp. 17–36. JMLR Workshop and Conference Proceedings (2012)
5. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9650–9660 (2021)
6. Chen, X., Xie, S., He, K.: An empirical study of training self-supervised vision transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9640–9649 (2021)
7. Chu, W.T., Wu, Y.L.: Image style classification based on learnt deep correlation features. IEEE Transactions on Multimedia **20**(9), 2491–2502 (2018)
8. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009)
9. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)
10. Duong, L.T., Le, N.H., Tran, T.B., Ngo, V.M., Nguyen, P.T.: Detection of tuberculosis from chest x-ray images: boosting the performance with vision transformer and transfer learning. Expert Systems with Applications **184**, 115519 (2021)
11. Fiorucci, M., Khoroshiltseva, M., Pontil, M., Traviglia, A., Del Bue, A., James, S.: Machine learning for cultural heritage: A survey. Pattern Recognition Letters **133**, 102–108 (2020)
12. Gonthier, N., Gousseau, Y., Ladjal, S., Bonfait, O.: Weakly supervised object detection in artworks. In: Proceedings of the European Conference on Computer Vision (ECCV) Workshops. pp. 0–0 (2018)
13. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
14. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
15. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems. vol. 25. Curran Associates, Inc. (2012)
16. Li, Y., Xie, S., Chen, X., Dollar, P., He, K., Girshick, R.: Benchmarking detection transfer learning with vision transformers. arXiv preprint arXiv:2111.11429 (2021)
17. Liu, Y., Sangineto, E., Bi, W., Sebe, N., Lepri, B., Nadai, M.: Efficient training of visual transformers with small datasets. In: Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., Vaughan, J.W. (eds.) Advances

in Neural Information Processing Systems. vol. 34, pp. 23818–23830. Curran Associates, Inc. (2021), https://proceedings.neurips.cc/paper/2021/file/c81e155d85dae5430a8cee6f2242e82c-Paper.pdf

18. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10012–10022 (2021)

19. Liu, Z., Mao, H., Wu, C.Y., Feichtenhofer, C., Darrell, T., Xie, S.: A convnet for the 2020s. arXiv preprint arXiv:2201.03545 (2022)

20. Masters, D., Luschi, C.: Revisiting small batch training for deep neural networks. arXiv preprint arXiv:1804.07612 (2018)

21. Matsoukas, C., Haslum, J.F., Söderberg, M., Smith, K.: Is it time to replace cnns with transformers for medical images? arXiv preprint arXiv:2108.09038 (2021)

22. Mensink, T., van Gemert, J.: The rijksmuseum challenge: Museum-centered visual recognition. In: ACM International Conference on Multimedia Retrieval (ICMR) (2014)

23. Milani, F., Fraternali, P.: A dataset and a convolutional model for iconography classification in paintings. Journal on Computing and Cultural Heritage (JOCCH) **14**(4), 1–18 (2021)

24. Minaee, S., Kafieh, R., Sonka, M., Yazdani, S., Soufi, G.J.: Deep-covid: Predicting covid-19 from chest x-ray images using deep transfer learning. Medical image analysis **65**, 101794 (2020)

25. Mormont, R., Geurts, P., Maree, R.: Comparison of deep transfer learning strategies for digital pathology. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops (June 2018)

26. Pan, S.J., Yang, Q.: A survey on transfer learning. IEEE Transactions on knowledge and data engineering **22**(10), 1345–1359 (2009)

27. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch (2017)

28. Sabatelli, M.: Contributions to Deep Transfer Learning: from Supervised to Reinforcement Learning. Ph.D. thesis, Universitè de Liegè, Liegè, Belgique (2022)

29. Sabatelli, M., Banar, N., Cocriamont, M., Coudyzer, E., Lasaracina, K., Daelemans, W., Geurts, P., Kestemont, M.: Advances in digital music iconography: Benchmarking the detection of musical instruments in unrestricted, non-photorealistic images from the artistic domain. Digital Humanities Quarterly **15**(1) (2021)

30. Sabatelli, M., Kestemont, M., Daelemans, W., Geurts, P.: Deep transfer learning for art classification problems. In: Proceedings of the European Conference on Computer Vision (ECCV) Workshops (2018)

31. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Gradcam: Visual explanations from deep networks via gradient-based localization. In: Proceedings of the IEEE international conference on computer vision. pp. 618–626 (2017)

32. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)

33. Strezoski, G., Worring, M.: Omniart: multi-task deep learning for artistic data analysis. arXiv preprint arXiv:1708.00684 (2017)

34. Talo, M., Baloglu, U.B., Yıldırım, Ö., Acharya, U.R.: Application of deep transfer learning for automated brain abnormality classification using mr images. Cognitive Systems Research **54**, 176–188 (2019)

35. Tan, M., Le, Q.: Efficientnetv2: Smaller models and faster training. In: International Conference on Machine Learning. pp. 10096–10106. PMLR (2021)
36. Tan, W.R., Chan, C.S., Aguirre, H.E., Tanaka, K.: Ceci n'est pas une pipe: A deep convolutional network for fine-art paintings classification. In: 2016 IEEE international conference on image processing (ICIP). pp. 3703–3707. IEEE (2016)
37. Taylor, M.E., Stone, P.: Transfer learning for reinforcement learning domains: A survey. Journal of Machine Learning Research **10**(7) (2009)
38. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jégou, H.: Training data-efficient image transformers & distillation through attention. In: International Conference on Machine Learning. vol. 139, pp. 10347–10357 (2021)
39. Van Den Oord, A., Dieleman, S., Schrauwen, B.: Transfer learning by supervised pre-training for audio-based music classification. In: Conference of the International Society for Music Information Retrieval (ISMIR 2014) (2014)
40. Vandaele, R., Dance, S.L., Ojha, V.: Deep learning for automated river-level monitoring through river-camera images: An approach based on water segmentation and transfer learning. Hydrology and Earth System Sciences **25**(8), 4435–4453 (2021)
41. Xue, F., Wang, Q., Guo, G.: Transfer: Learning relation-aware facial expression representations with transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 3601–3610 (2021)
42. Ying, W., Zhang, Y., Huang, J., Yang, Q.: Transfer learning via learning to transfer. In: International Conference on Machine Learning. pp. 5085–5094. PMLR (2018)
43. Zhong, S.h., Huang, X., Xiao, Z.: Fine-art painting classification via two-channel dual path networks. International Journal of Machine Learning and Cybernetics **11**(1), 137–152 (2020)
44. Zhou, H.Y., Lu, C., Yang, S., Yu, Y.: Convnets vs. transformers: Whose visual representations are more transferable? In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 2230–2238 (2021)
45. Zhu, Z., Lin, K., Zhou, J.: Transfer learning in deep reinforcement learning: A survey. arXiv preprint arXiv:2009.07888 (2020)
46. Zhuang, F., Qi, Z., Duan, K., Xi, D., Zhu, Y., Zhu, H., Xiong, H., He, Q.: A comprehensive survey on transfer learning. Proceedings of the IEEE **109**(1), 43–76 (2020)