# EfficientTrain: Exploring Generalized Curriculum Learning for Training Visual Backbones

Yulin Wang[1*]    Yang Yue[1*]    Rui Lu[1]    Tianjiao Liu[2]    Zhao Zhong[2]

Shiji Song[1]    Gao Huang[1,3 ✉]

[1]Department of Automation, BNRist, Tsinghua University    [2]Huawei Technologies Ltd.
[3]Beijing Academy of Artificial Intelligence (BAAI)

{wang-yl19, yueyang22, r-lu21}@mails.tsinghua.edu.cn,
{liutianjiao3, zorro.zhongzhao}@huawei.com, {shijis, gaohuang}@tsinghua.edu.cn

## Abstract

*The superior performance of modern deep networks usually comes at the price of a costly training procedure. In this paper, we present a novel curriculum learning approach for the efficient training of visual backbones (e.g., vision Transformers). The proposed method is inspired by the phenomenon that deep networks mainly learn to recognize some 'easier-to-learn' discriminative patterns within each example at earlier stages of training, e.g., the lower-frequency components of images and the original information before data augmentation. Driven by this observation, we propose a curriculum where the model always leverages all the training data at each epoch, while the curriculum starts with only exposing the 'easier-to-learn' patterns of each example, and introduces gradually more difficult patterns. To implement this idea, we 1) introduce a cropping operation in the Fourier spectrum of the inputs, which enables the model to learn from only the lower-frequency components efficiently, and 2) demonstrate that exposing the features of original images amounts to adopting weaker data augmentation. Our resulting algorithm, EfficientTrain, is simple, general, yet surprisingly effective. For example, it reduces the training time of a wide variety of popular models (e.g., ConvNeXts, DeiT, PVT, and Swin/CSWin Transformers) by more than $1.5\times$ on ImageNet-1K/22K without sacrificing the accuracy. It is effective for self-supervised learning (i.e., MAE) as well. Code is available at https://github.com/LeapLabTHU/EfficientTrain.*

## 1. Introduction

The success of modern visual backbones is largely fueled by the interest in exploring big models on large-scale benchmark datasets [15, 27, 28, 44]. In particular, the recent intro-
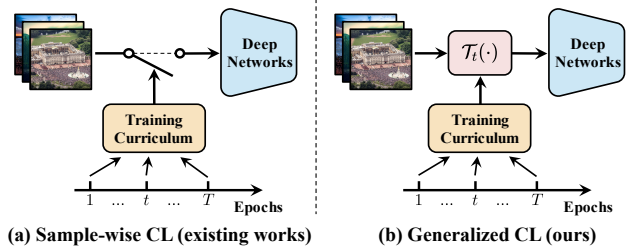


**(a) Sample-wise CL (existing works)**          **(b) Generalized CL (ours)**

Figure 1. **(a) Sample-wise curriculum learning (CL):** making a *discrete* decision on whether each example should be leveraged to train the model. **(b) Generalized CL:** we consider a *continuous* function $\mathcal{T}_t(\cdot)$, which only exposes the *'easier-to-learn'* patterns within each example at the beginning of training (*e.g.*, *lower-frequency* components; see: Section 4), while gradually introducing relatively *more difficult* patterns as learning progresses.

duction of vision Transformers (ViTs) scales up the number of model parameters to more than 1.8 billion, with the training data growing to 3 billion samples [15, 89]. Although a state-of-the-art accuracy has been attained, this huge-model and high-data regime yields a time-consuming and costly training process. For example, it takes 2,500 TPUv3-core-days to train ViT-H/14 on JFT-300M [15], which may be unaffordable for the practitioners from both academia and industry. The power consumption will lead to considerable carbon emissions as well [39, 64]. For both economical and environmental reasons, a surge in the demand for saving the training cost of modern deep networks has been witnessed.

In this paper, we contribute to this issue by revisiting the idea of curriculum learning [2], which reveals that a model can be trained efficiently by starting with the easier aspects of a given task or certain easier subtasks, and increasing the difficulty level gradually. Most existing works implement this idea by introducing easier-to-harder examples progressively during training [17, 19, 22, 24, 31, 32, 63, 80]. However, obtaining a light-weighted and generalizable difficulty measurer is typically non-trivial [63, 80]. In general, these methods tend to be immature for being applied as a regular efficient training technique.

---

*Equal contribution.          ✉Corresponding author.

1

In contrast to the prior works, this paper seeks a simple yet broadly applicable efficient learning approach that has the potential to be widely implemented. To attain this goal, we consider a generalization of curriculum learning. In specific, we extend the notion of training curricula beyond only differentiating between 'easier' and 'harder' examples, and adopt a more flexible hypothesis, which indicates that the discriminative features of each training sample consist of both 'easier-to-learn' and 'harder-to-learn' patterns. Instead of making a *discrete* decision on whether each example should appear in the training set, we argue that it would be more proper to establish a *continuous* function that adaptively extracts the simpler and more learnable discriminative patterns within every example. In other words, a curriculum may always leverage all examples at any stage of learning, but it should eliminate the relatively more difficult or complex patterns within inputs at earlier learning stages. An illustration of our idea is shown in Figure 1.

Driven by our hypothesis, a straightforward yet surprisingly effective algorithm is derived. We first demonstrate that the 'easier-to-learn' patterns incorporate the lower-frequency components of the images. We further show that a lossless extraction of these components can be achieved by introducing a cropping operation in the frequency domain. This operation not only retains exactly all the lower-frequency information, but also yields a smaller input size for the model to be trained. By triggering this operation at earlier training stages, the overall computational cost for training can be considerably reduced while the final performance of the model will not be sacrificed. Moreover, we observe that the original information before heavy data augmentation is more learnable, and hence starting the training with weaker augmentation techniques is beneficial. These theoretical and experimental insights are integrated into a unified curriculum, which we refer to as '*EfficientTrain*'. One of the most appealing advantages of EfficientTrain may be its simplicity and generalizability. Our method can be conveniently applied to most deep networks *without any modification or hyper-parameter tuning*, but significantly improves their training efficiency. Empirically, for the supervised learning on ImageNet-1K/22K [12], EfficientTrain reduces the wall-time training cost of a wide variety of state-of-the-art visual backbones (*e.g.*, ResNet [27], ConvNeXt [45], DeiT [69], PVT [79], Swin [44] / CSWin [14] Transformer) by *more than* $1.5\times$, while achieving competitive or better performance compared with the baselines. Importantly, EfficientTrain is also effective on top of a representative self-supervised learning algorithm (*i.e.*, MAE [25]).

## 2. Related Work

**Curriculum learning** is a training paradigm inspired by the organized learning order of examples in human curricula [2, 16, 37]. This idea has been widely explored in the context of training deep networks from easier data to harder data [17, 22, 24, 38, 55, 63, 80]. Typically, a pre-defined [2, 7, 73, 83] or automatically-learned [17, 22, 24, 30, 32, 38, 49, 58, 74, 84, 90] difficulty measurer is deployed to differentiate between easier and harder samples, while a scheduler [2, 22, 55, 80] is defined to determine when and how to introduce harder training data. Our method is also based on the 'starting small' spirit [16], but we always leverage all the training data simultaneously. Our work is also related to curriculum by smoothing [62] and curriculum dropout [51], which do not perform example selection as well. However, our method is orthogonal to them since we propose to reduce the training cost by modifying the model inputs, while they focus on regularizing the deep features during training (*e.g.*, via anti-aliasing smoothing or feature dropout).

**Progressive training.** It has been shown that deep networks can be trained efficiently with an increasing model size during training, *e.g.*, a growing number of layers [33, 39, 61, 75], a growing network width [6], or a dynamically changing topology of network connections [81, 82]. These methods are mainly motivated by that smaller models are more efficient to train at earlier learning stages. This idea is also explored in natural language processing [20, 92], recommendation systems [77] and graph ConvNets [87].

A similar work to us is progressive learning (PL) [67], which down-samples the images to save the training cost. Nevertheless, our work differs from PL in several important aspects: 1) EfficientTrain is drawn from a distinctly different motivation of generalized curriculum learning, based on which we present a novel frequency-inspired analysis; 2) we introduce a cropping operation in the frequency domain, which is not only theoretically different from the down-sampling operation in PL (see: Proposition 1), but also outperforms it empirically (see: Table 12); 3) from the perspective of system-level comparison, we design an EfficientTrain curriculum, achieving a significantly higher training efficiency than PL on a variety of state-of-the-art models (see: Tables 8). In addition, FixRes [72] shows that a smaller training resolution may improve the accuracy by fixing the discrepancy between the scale of training and test inputs. However, we do not borrow gains from FixRes as we adopt a standard resolution at the final stages of training. Our method is actually orthogonal to FixRes (see: Table 8).

**Frequency-based analysis of deep networks.** Our observation that deep networks tend to capture the low-frequency components first is inline with [76], but the discussions in [76] focus on the robustness of ConvNets and are mainly based on some small models and tiny datasets. Towards this direction, several existing works also explore decomposing the inputs of models in the frequency domain [46, 53, 86] in order to understand or improve the robustness of the networks. In contrast, our aim is to improve the training efficiency of modern deep networks.

## 3. A Generalization of Curriculum Learning

As uncovered in previous research, machine learning algorithms generally benefit from a 'starting small' strategy, *i.e.*, to first learn certain easier aspects of a task, and increase the level of difficulty progressively [2,16,37]. The dominant implementation of this idea, curriculum learning, proposes to introduce gradually more difficult examples during training [24,63,80]. In specific, a curriculum is defined on top of the training process to determine whether or not each sample should be leveraged at a given epoch (see: Figure 1 (a)).

**On the limitations of sample-wise curriculum learning.** Although curriculum learning has been widely explored from the lens of the sample-wise regime, its extensive usage is usually limited by two major issues. *First*, differentiating between 'easier' and 'harder' training data is non-trivial. It typically requires deploying additional deep networks as a 'teacher' or exploiting specialized automatic learning approaches [17, 22, 24, 32, 38]. The resulting implementation complexity and the increased overall computational cost are both noteworthy weaknesses in terms of improving the training efficiency. *Second*, it is challenging to attain a principled approach that specifies which examples should be attended to at the earlier stages of learning. As a matter of fact, the 'easy to hard' strategy is not always helpful [80]. The hard-to-learn samples can be more informative and may be beneficial to be emphasized in many cases [1,18,21,47,60,97], sometimes even leading to a 'hard to easy' anti-curriculum [3, 43, 54, 78, 93, 96].

Our work is motivated by the above two issues. In the following, we start by proposing a generalized assumption for curriculum learning, aiming to address the second issue. Then we further demonstrate that an implementation of our idea naturally addresses the first issue.

**Generalized curriculum learning.** We argue that simply measuring the easiness of training samples tends to be ambiguous and may be insufficient to reflect the effects of a sample on the learning process. As aforementioned, even the difficult examples may provide beneficial information for guiding the training, and they do not necessarily need to be introduced after easier examples. To this end, we hypothesize that every training sample, either 'easier' or 'harder', contains both *easier-to-learn* or *more accessible* patterns, as well as certain *difficult* discriminative information which may be challenging for the deep networks to capture. Ideally, a curriculum should be a continuous function on top of the training process, which starts with a focus on the 'easiest' patterns of the inputs, while the 'harder-to-learn' patterns are gradually introduced as learning progresses.

A formal illustration is shown in Figure 1 (b). Any input data $X$ will be processed by a transformation function $\mathcal{T}_t(\cdot)$ conditioned on the training epoch $t$ ($t \leq T$) before being fed into the model, where $\mathcal{T}_t(\cdot)$ is designed to dynamically filter out the excessively difficult and less learnable patterns

within the training data. We always let $\mathcal{T}_T(X) = X$. Notably, our approach can be seen as a generalized form of the sample-wise curriculum learning. It reduces to example-selection by setting $\mathcal{T}_t(X) \in \{\emptyset, X\}$.

**Overview.** In the rest of this paper, we will demonstrate that an algorithm drawn from our hypothesis dramatically improves the implementation efficiency and generalization ability of curriculum learning. We will show that a zero-cost criterion pre-defined by humans is able to effectively measure the difficulty level of different patterns within images. Based on such simple criteria, even a surprisingly straightforward implementation of introducing 'easier-to-harder' patterns yields significant and consistent improvements on the training efficiency of modern visual backbones.

## 4. The EfficientTrain Approach

To obtain a training curriculum following our aforementioned hypothesis, we need to solve two challenges: 1) identifying the patterns that are easier or more difficult for the networks to capture; 2) designing a schedule to dynamically filter out the relatively more difficult patterns within inputs. This section will demonstrate how both issues are addressed, leading to our simple but effective EfficientTrain algorithm. The implementation details of the experiments in this section are deferred to Appendix A.

### 4.1. Frequency-inspired Curricula

Image-based data can naturally be decomposed in the frequency domain. The lower-frequency components describe the smoothly changing contents, while the higher-frequency counterparts correspond to the finer details [5, 8, 48, 65]. In this subsection, an analysis will be presented to reveal that the 'easier-to-learn' patterns we are looking for at least include the lower-frequency parts of images.

**Experiments with the low-pass filtered input data.** We first consider an ablation study, where the low-pass filtering is performed on the data we use. As shown in Figure 2 (a), we map the images to the Fourier spectrum with the lowest frequency at the centre, set all the components outside a centred circle (radius: $r$) to zero, and map the spectrum back to the pixel space. Figure 2 (b) illustrates the effects of $r$. The curves of accuracy v.s. training epochs on top of the low-pass filtered data are presented in Figure 3. Here both training and validation data is processed by the filter to ensure the compatibility with the i.i.d. assumption.

**Lower-frequency components are captured first.** The models in Figure 3 are imposed to leverage only the lower-frequency components of the inputs. However, an appealing phenomenon arises that their training process is approximately identical to the original baseline at the beginning of the training. The baseline outperforms finally, but this tendency starts midway in the training process, instead of from the very beginning. Moreover, consider increasing the value of the filter bandwidth $r$, which progressively pre-

**(a) Low-pass Filtering**
**(DFT: discrete Fourier transform)**
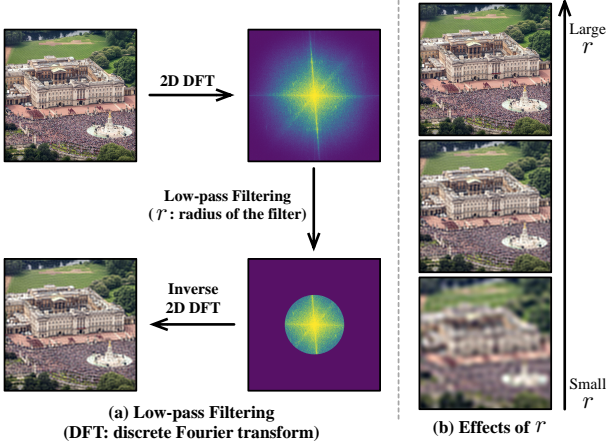
**(b) Effects of $r$**

Figure 2. **Low-pass filtering.** Following [76], we adopt a circular filter.

serves more information within the images from the lowest frequency. The corresponding separation point of the baseline and the training process with the low-pass filtered data moves towards the end of training. A similar phenomenon occurs as well if the low-pass filtering is performed only on either the training or the validation data (see: Figure 4). To explain these observations, we postulate that *the models tend to first learn to capture the lower-frequency information of the inputs, while the relatively higher-frequency components are gradually exploited on the basis of it*.

**More evidence to support our assumption** can be found in Appendix C.1. As an overview, we highlight two observations. *First*, as an opposite case to Figure 3, learning on the *high-pass* filtered data yields a significantly reduced accuracy compared to the baseline from the very beginning, while the gap is not closed or shrunk during the whole training process. Besides, this accuracy degradation becomes more serious when more lower-frequency components are eliminated. *Second*, we train a model using original images, and evaluate all the intermediate checkpoints on both low-pass and high-pass filtered validation sets. The bandwidths of the low/high-pass filters are adjusted to make the finally trained model have the same accuracy on the two validation sets. We find that the accuracy on the low-pass filtered validation set grows much faster at the beginning stages of training, although the two final accuracies are equal.

**Frequency-based training curricula.** Returning to our hypothesis in Section 3, our analysis thus far has validated that the lower-frequency components are naturally captured earlier. Hence, it would be straightforward to consider them as a type of the 'easier-to-learn' patterns. This begs a question: can we design a curriculum, which starts with providing only the lower-frequency information for the model to be trained, while gradually introducing the components with higher-frequency? We investigate this idea in Table 1, where we perform low-pass filtering on the training data only in a given number of the beginning epochs. The rest of the training process remains unchanged.
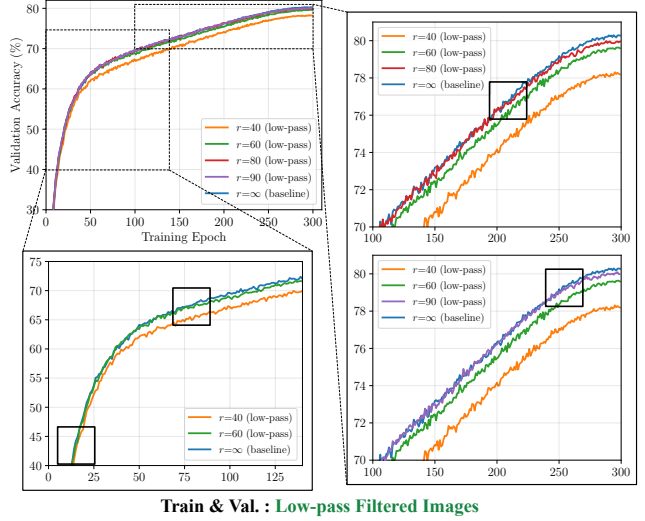


**Train & Val. : Low-pass Filtered Images**

Figure 3. **Ablation studies with low-pass filtering** ($r$: bandwidth, see: Figure 2). We ablate the higher-frequency components of the inputs of a DeiT-Small [69], and present the curves of validation accuracy v.s. training epochs on ImageNet-1K. The low-pass filtering is performed on both training and validation data to ensure the compatibility with the i.i.d. assumption. We highlight the separation points of the curves with black boxes.



**(a) Train: Original Images;**
**Val. : Low-pass Filtered Images;**

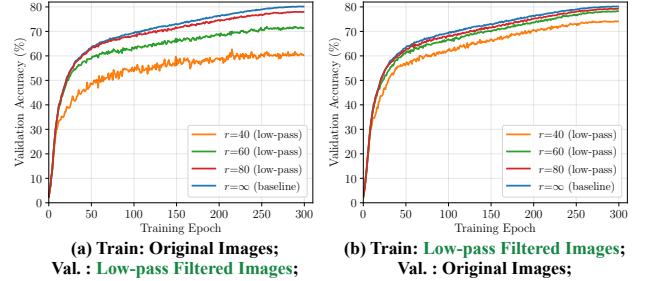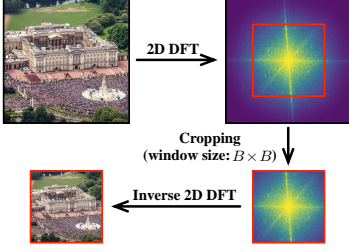**(b) Train: Low-pass Filtered Images;**
**Val. : Original Images;**

Figure 4. **Performing low-pass filtering only on the training or the validation data** (other setups follow from Figure 3). Although these results do not strictly satisfy the i.i.d. assumption, they provide consistent observations with Figure 3, which may serve as effective supplementary evidence.

| Curricula (ep: epoch) | | Final Top-1 Accuracy | | |
| Low-pass Filtered Training Data | Original Training Data | ($r$: filter bandwidth) | | |
| | | $r$=40 | $r$=60 | $r$=80 |
|---|---|---|---|---|
| $1^{st}$ – $300^{th}$ ep | – | 78.3% | 79.6% | 80.0% |
| $1^{st}$ – $225^{th}$ ep | $226^{th}$ – $300^{th}$ ep | 79.4% | <u>80.2%</u> | <u>80.5%</u> |
| $1^{st}$ – $150^{th}$ ep | $151^{th}$ – $300^{th}$ ep | <u>80.1%</u> | 80.2% | 80.6% |
| $1^{st}$ – $75$ ep | $76^{th}$ – $300^{th}$ ep | 80.3% | 80.4% | 80.6% |
| – | $1^{st}$ – $300^{th}$ ep | *80.3% (baseline)* | | |

Table 1. **Results with the straightforward frequency-based training curricula** (DeiT-Small [69] on ImageNet-1K). Observation: one can eliminate the higher-frequency components of the inputs in 50-75% of the training process without sacrificing the final accuracy (see: <u>underlined</u> data).

**Observations from Table 1.** At the first glance, the results in Table 1 may seem less dramatic, *i.e.*, by processing the images with a properly-configured low-pass filter at earlier epochs, the accuracy is improved moderately (*e.g.*, 0.2-0.3%). However, an important observation is noteworthy: the performance of the model is capable of being largely preserved even with a relatively aggressive filtering (*e.g.*, $r = 40, 60$) performed in 50-75% of the training process. This phenomenon turns our attention to the perspective of

Figure 5. **Low-frequency cropping in the frequency domain** ($B^2$: bandwidth).

Table 2. **Results on ImageNet-1K with the low-pass filtering in Table 1 replaced by the low-frequency cropping**, which *yields competitive accuracy with a significantly reduced training cost* (see: underlined data).

| Curricula (ep: epoch) | | Final Top-1 Accuracy/Computational Cost for Training (relative) | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Low-frequency Cropping ($B^2$) | Original Training Data ($B=224$) | DeiT-Small [69] | | | | Swin-Tiny [44] | | | |
| | | $B=96$ | $B=128$ | $B=160$ | $B=192$ | $B=96$ | $B=128$ | $B=160$ | $B=192$ |
| 1st – 300th ep | – | 70.5%/0.18 | 75.3%/0.31 | 77.9%/0.49 | 79.1%/0.72 | 73.3%/0.18 | 76.8%/0.32 | 78.9%/0.50 | 80.5%/0.73 |
| 1st – 225th ep | 226th – 300th ep | 78.7%/0.38 | 79.6%/0.48 | 80.0%/0.62 | <u>80.3%/0.79</u> | 80.0%/0.38 | 80.5%/0.49 | 81.0%/0.63 | <u>81.2%/0.80</u> |
| 1st – 150th ep | 151th – 300th ep | 79.2%/0.59 | 79.8%/0.66 | <u>80.3%/0.75</u> | 80.4%/0.86 | 80.9%/0.59 | 80.9%/0.66 | <u>81.2%/0.75</u> | 81.3%/0.86 |
| 1st – 75th ep | 76th – 300th ep | 79.6%/0.79 | <u>80.2%/0.83</u> | 80.4%/0.87 | 80.3%/0.93 | <u>81.2%/0.79</u> | <u>81.2%/0.83</u> | 81.3%/0.88 | 81.3%/0.93 |
| – | 1st – 300th ep | *80.3%/1.00 (baseline)* | | | | *81.3%/1.00 (baseline)* | | | |

training efficiency. At earlier learning stages, it is feasible to train the model with only the lower-frequency components. These components incorporate only a selected subset of all the information within the original input images. Therefore, *can we enable the model to learn from them efficiently with less computational cost than processing the original inputs?*

**Low-frequency cropping in the Fourier spectrum.** As a matter of fact, the above question can be answered by introducing a cropping operation in the frequency domain. Consider mapping an $H \times W$ image $\boldsymbol{X}$ into the frequency domain with the 2D discrete Fourier transform (DFT), obtaining a $H \times W$ Fourier spectrum, where the value in the centre denotes the strength of the component with the lowest frequency. The positions distant from the centre correspond to the higher-frequency components. Herein, we propose to crop a $B \times B$ patch from the centre of the spectrum, where $B$ is the window size. Since the patch is still centrosymmetric, we can map it back to the pixel space with the inverse 2D DFT, obtaining a $B \times B$ new image $\boldsymbol{X}_c$, namely

$$\boldsymbol{X}_c = \mathcal{F}^{-1} \circ \mathcal{C}_{B,B} \circ \mathcal{F}(\boldsymbol{X}) \in \mathbb{R}^{B \times B}, \tag{1}$$
$$\boldsymbol{X} \in \mathbb{R}^{H \times W}, \; B < H, B < W,$$

where $\mathcal{F}$ and $\mathcal{F}^{-1}$ denote 2D DFT and inverse 2D DFT, while $\mathcal{C}_{B,B}$ denotes $B^2$ centre-cropping. An illustration of this procedure is presented in Figure 5.

Notably, $\boldsymbol{X}_c$ attains the goal of extracting exactly all the lower-frequency components while eliminating the rest higher-frequency parts. Hence, feeding $\boldsymbol{X}_c$ into the model at earlier training stages can provide the vast majority of the useful information, such that the final accuracy will be minimally affected or even not affected. In contrast, importantly, due to the reduced input size of $\boldsymbol{X}_c$, the computational cost for a model to process $\boldsymbol{X}_c$ is able to be quadratically saved, yielding a considerably more efficient training process.

Our claims are empirically supported by the results in Table 2, where we simply replace the low-pass filtering in Table 1 with our low-frequency cropping. Even such a straightforward implementation yields favorable results. The training cost can be saved by $\sim$20% while a competitive final accuracy is preserved. These results can be interpreted through the feature similarity. As shown in Figure 6, at the intermediate stages of training, the models trained using the inputs with low-frequency cropping learn similar deep rep-
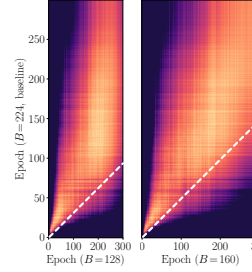


Figure 6. **CKA feature similarity** heatmaps [9, 35] between the DeiT-Small [69] trained using the inputs with low-frequency cropping ($B=128, 160$) and the original inputs ($B=224$). The X and Y axes index training epochs (scaled according to the computational cost of training). Here we feed the same original images into all the models (including the ones trained with $B=128, 160$) and take the features from the final layer. The 45° lines are highlighted in white.

| Curricula (ep: epoch) | | | | Final Top-1 Accuracy | | | | |
|---|---|---|---|---|---|---|---|---|
| Weaker RandAug | RandAug ($m=9$) | Low-frequency ($B=128$) | Original ($B=224$) | ($m$: magnitude of RandAug) | | | | |
| | | | | $m=1$ | $m=3$ | $m=5$ | $m=7$ | $m=9$ |
| 1st – 150th ep | 151th – 300th ep | – | 1st–300th ep | 80.4% | 80.6% | **80.7%** | 80.5% | 80.3% |
| | | 1st–150th ep | 151th–300th ep | **80.2%** | **80.2%** | **80.2%** | 79.9% | 79.8% |

Table 3. **Performance of the data-augmentation-based curricula** (DeiT-Small [69] on ImageNet-1K). We test reducing the magnitude of RandAug in 1st-150th training epochs ($m=9$ refers to the baselines).

resentations to the baseline at a significantly reduced cost, *i.e.*, the bright parts are clearly above the white lines.

**Comparisons with image down-sampling.** Our method seeks to only expose the lower-frequency information at earlier training stages. This idea can also be approximately implemented by the pixel-space down-sampling operation, such that the cost of 2D DFT can be saved. However, down-sampling cannot strictly filter out all the higher-frequency components as low-frequency cropping does (see: Proposition 1). We also observe that it degrades the performance empirically (see: Table 12). In addition, we note that 2D DFT can be accomplished with a negligible cost on GPUs.

**Proposition 1.** *Suppose that $\boldsymbol{X}_c = \mathcal{F}^{-1} \circ \mathcal{C}_{B,B} \circ \mathcal{F}(\boldsymbol{X})$, and that $\boldsymbol{X}_d = \mathcal{D}_{B,B}(\boldsymbol{X})$, where $B \times B$ down-sampling $\mathcal{D}_{B,B}(\cdot)$ is realized by a common interpolation algorithm (e.g., nearest, bilinear or bicubic). Then we have two properties.*
*● 1) $\boldsymbol{X}_c$ is only determined by the lower-frequency spectrum of $\boldsymbol{X}$ (i.e., $\mathcal{C}_{B,B} \circ \mathcal{F}(\boldsymbol{X})$). In addition, the mapping to $\boldsymbol{X}_c$ is reversible. We can always recover $\mathcal{C}_{B,B} \circ \mathcal{F}(\boldsymbol{X})$ from $\boldsymbol{X}_c$.*
*● 2) $\boldsymbol{X}_d$ has a non-zero dependency on the higher-frequency spectrum of $\boldsymbol{X}$ (i.e., the regions outside $\mathcal{C}_{B,B} \circ \mathcal{F}(\boldsymbol{X})$).*

***Proof.*** See: Appendix D. □

### 4.2. Data Augmentation

Moderns deep networks (*e.g.*, vision Transformers [14, 44, 69] and ConvNets [10, 45, 94]) are typically trained with strong and delicate data augmentation techniques. We ar-

| | Model | Input Size (inference) | #Param. | #FLOPs | Top-1 Accuracy (300 epochs) | | | Training Speedup | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | Original Paper | Baseline (ours) | **EfficientTrain** | Computation | **Wall-time** |
| *ConvNets* | ResNet-50 [27] | $224^2$ | 26M | 4.1G | – | 78.8% | **79.4%** | 1.53× | **1.44×** |
| | ConvNeXt-Tiny [45] | $224^2$ | 29M | 4.5G | 82.1% | 82.1% | **82.2%** | 1.53× | **1.49×** |
| | ConvNeXt-Small [45] | $224^2$ | 50M | 8.7G | 83.1% | 83.1% | **83.2%** | 1.53× | **1.50×** |
| | ConvNeXt-Base [45] | $224^2$ | 89M | 15.4G | 83.8% | **83.8%** | 83.7% | 1.53× | **1.48×** |
| *Isotropic ViTs* | DeiT-Tiny [69] | $224^2$ | 5M | 1.3G | 72.2% | 72.5% | **73.3%** | 1.59× | **1.55×** |
| | DeiT-Small [69] | $224^2$ | 22M | 4.6G | 79.9% | 80.3% | **80.3%** | 1.56× | **1.51×** |
| *Multi-stage ViTs* | PVT-Tiny [79] | $224^2$ | 13M | 1.9G | 75.1% | 75.5% | **75.5%** | 1.55× | **1.48×** |
| | PVT-Small [79] | $224^2$ | 25M | 3.8G | 79.8% | 79.9% | **80.4%** | 1.55× | **1.56×** |
| | PVT-Medium [79] | $224^2$ | 44M | 6.7G | 81.2% | 81.8% | **81.8%** | 1.54× | **1.49×** |
| | PVT-Large [79] | $224^2$ | 61M | 9.8G | 81.7% | 82.3% | **82.3%** | 1.54× | **1.53×** |
| | Swin-Tiny [44] | $224^2$ | 28M | 4.5G | 81.3% | 81.3% | **81.3%** | 1.55× | **1.49×** |
| | Swin-Small [44] | $224^2$ | 50M | 8.7G | 83.0% | 83.1% | **83.2%** | 1.54× | **1.50×** |
| | Swin-Base [44] | $224^2$ | 88M | 15.4G | 83.5% | 83.4% | **83.6%** | 1.54× | **1.50×** |
| | CSWin-Tiny [14] | $224^2$ | 23M | 4.3G | 82.7% | 82.7% | **82.8%** | 1.56× | **1.55×** |
| | CSWin-Small [14] | $224^2$ | 35M | 6.9G | 83.6% | 83.4% | **83.6%** | 1.56× | **1.51×** |
| | CSWin-Base [14] | $224^2$ | 78M | 15.0G | 84.2% | 84.3% | **84.3%** | 1.55× | **1.56×** |

Table 5. **Results on ImageNet-1K (IN-1K).** We train the models w/ or w/o EfficientTrain on the IN-1K training set, and report the accuracy on the IN-1K validation set. For fair comparisons, we also report the baselines in the original papers. The training wall-time is benchmarked on NVIDIA 3090 GPUs.

| Epochs | Low-frequency Cropping | RandAug |
|---|---|---|
| $1^{st}$ – $180^{th}$ | $B = 160$ | |
| $181^{th}$ – $240^{th}$ | $B = 192$ | $m = 0 \to 9$ Increase linearly. |
| $241^{th}$ – $300^{th}$ | $B = 224$ | |

Table 4. **The EfficientTrain curriculum**. Here the standard 300-epoch training configuration in [44,45] is considered.

gue that the augmented training data provides a combination of both the information from the original samples and the information introduced by the augmentation operations. The original information may be 'easier-to-learn' as its features are drawn from the real distribution. This assumption is supported by the observation that data augmentation is mainly influential at the later stages of training [68].

Therefore, following our generalized formulation of curriculum learning in Section 3, we may design the curriculum by adopting a weaker-to-stronger data augmentation strategy during training. For example, in Table 3, we investigate varying the magnitude of the rand augmentation (RandAug) [11] in the first half training process. One can observe that this idea improves the accuracy, while the gains are compatible with the low-frequency cropping operation.

### 4.3. A Unified Training Curriculum

Finally, a unified efficient training curriculum is designed by integrating the techniques introduced in Sections 4.1 and 4.2. In specific, we first let the magnitude $m$ of the rand augmentation be a linear function of the training epoch $t$, namely $m = (t/T) \times m_0$, with other data augmentation techniques unchanged. Although being simple, this configuration leads to consistent and significant empirical improvements. Note that we adopt $m_0 = 9$ following the common practice [14,44,45,69,79]. Then we leverage a greedy-search algorithm to determine the schedule of $B$ (*i.e.*, the window size for low-frequency cropping) during training, for which the details are deferred to Appendix B.

Our proposed curriculum is shown in Table 4, which is named as *EfficientTrain*. Notably, despite the simplicity of EfficientTrain, it is well-generalizable and consistently ef-

| Models | Input Size (inference) | Top-1 Accuracy (baseline / **EfficientTrain**) | | | Wall-time Training Speedup |
|---|---|---|---|---|---|
| | | 100 epochs | 200 epochs | 300 epochs | |
| DeiT-Tiny [69] | $224^2$ | 65.8%/**68.1%** | 70.5%/**71.8%** | 72.5%/**73.3%** | 1.55× |
| DeiT-Small [69] | $224^2$ | 75.5%/**76.4%** | 79.0%/**79.1%** | 80.3%/**80.3%** | 1.51× |
| Swin-Tiny [44] | $224^2$ | 78.4%/**78.5%** | 80.6%/**80.6%** | 81.3%/**81.3%** | 1.49× |
| Swin-Small [44] | $224^2$ | 80.6%/**80.7%** | 82.7%/82.6% | 83.1%/**83.2%** | 1.50× |
| Swin-Base [44] | $224^2$ | 80.7%/**81.1%** | 83.2%/**83.2%** | 83.4%/**83.6%** | 1.50× |

(a) **Reducing training cost with the same number of epochs.**

| Models | Input Size (inference) | Top-1 Accuracy (ep: epochs) | | Wall-time Training Speedup |
|---|---|---|---|---|
| | | Baseline_{300ep} | **EfficientTrain_{450ep}** | |
| DeiT-Tiny [69] | $224^2$ | 72.5% | **74.3%** (+1.8) | 1.04× |
| DeiT-Small [69] | $224^2$ | 80.3% | **80.9%** (+0.6) | 1.01× |
| Swin-Tiny [44] | $224^2$ | 81.3% | **81.7%** (+0.4) | 1.00× |

(b) **For smaller models: higher accuracy with the same training cost.**

| Models | Input Size (inference) | Top-1 Accuracy (ep: epochs) | | **Wall-time Training Speedup** |
|---|---|---|---|---|
| | | Baseline_{300ep} | **EfficientTrain_{200ep}** | |
| Swin-Small [44] | $224^2$ | 83.1% | 82.6% (-0.5) | **2.26×** |
| Swin-Base [44] | $224^2$ | 83.4% | 83.2% (-0.2) | **2.25×** |

(c) **For larger models: more efficient training with an acceptable accuracy drop** (less than 0.6% of the fully converged accuracy).

Table 6. **Results on ImageNet-1K with varying training epochs.**

fective. It can be straightforwardly deployed on top of common visual backbones under varying training settings, and contributes to a significantly improved training efficiency.

## 5. Experiments

In this section, we will present a comprehensive empirical evaluation to validate the effectiveness of EfficientTrain.

**Datasets.** Our main results are based on the large-scale ImageNet-1K/22K [12] benchmark datasets, which consist of ∼1.28M/∼14.2M images in 1K/∼22K classes. We also conduct experiments on COCO-2017 [42], CIFAR-10/100 [36], Oxford Flowers-102 [52], and Stanford Dogs [34] to verify the transferability of our learned models.

**Models.** A wide variety of visual backbones are considered in our experiments, including ResNet [27], ConvNeXt [45], DeiT [69], PVT [79], Swin [44] and CSWin [14] Transformers. We adopt the training pipeline in [44, 45], where EfficientTrain only modifies the terms mentioned in Table 4. Unless otherwise specified, we report the results of

| Model | Input Size (inference) | #Param. | #FLOPs | Top-1 Accuracy (fine-tuned to ImageNet-1K) Baseline | EfficientTrain | Wall-time Pre-training Cost (in GPU-days) Baseline | EfficientTrain | Time Saved (for an 8-GPU node) |
|---|---|---|---|---|---|---|---|---|
| ConvNeXt-Base [45] | $224^2$ | 89M | 15.4G | 85.6% | **85.7%** | 170.6 | **114.8** (↓1.49×) | *6.98 Days* |
| | $384^2$ | 89M | 45.1G | 86.7% | **86.8%** | | | |
| ConvNeXt-Large [45] | $224^2$ | 198M | 34.4G | **86.4%** | 86.3% | 347.6 | **225.5** (↓1.54×) | *15.26 Days* |
| | $384^2$ | 198M | 101.0G | 87.3% | **87.3%** | | | |
| CSWin-Base [14] | $224^2$ | 78M | 15.0G | 85.5% | **85.6%** | 238.9 | **157.7** (↓1.52×) | *10.16 Days* |
| | $384^2$ | 78M | 47.0G | 86.7% | **87.0%** | | | |
| CSWin-Large [14] | $224^2$ | 173M | 31.5G | 86.5% | **86.6%** | 469.5 | **307.7** (↓1.53×) | *20.23 Days* |
| | $384^2$ | 173M | 96.8G | 87.6% | **87.8%** | | | |

Table 7. **Results with ImageNet-22K (IN-22K) pre-training.** The models are pre-trained on IN-22K w/ or w/o EfficientTrain, fine-tuned on the ImageNet-1K (IN-1K) training set, and evaluated on the IN-1K validation set. The wall-time pre-training cost is benchmarked on NVIDIA 3090 GPUs.

| Models | Training Approach | Training Epochs | Aug-Regs | Top-1 Accuracy | Wall-time Training Speedup |
|---|---|---|---|---|---|
| ResNet-18 [27] | Curriculum by Smoothing [62] *(NeurIPS'20)* | 90 | ✗ | 71.0% | 1.00× |
| | **EfficientTrain** | 90 | ✗ | **71.0%** | **1.48×** |
| ResNet-50 [27] | Self-paced Learning [38] *(NeurIPS'10)* | 200 | ✗ | 73.2% | 1.15× |
| | Minimax Curriculum Learning [96] *(ICLR'18)* | 200 | ✗ | 75.1% | 1.97× |
| | DIH Curriculum Learning [97] *(NeurIPS'20)* | 200 | ✗ | 76.3% | 2.45× |
| | **EfficientTrain** | 200 | ✗ | **77.5%** | 1.44× |
| | CurriculumNet [23] *(ECCV'18)* | 90 | ✗ | 76.1% | <2.22× |
| | Label-sim. Curriculum Learning [13] *(ECCV'20)* | 90 | ✗ | 76.9% | 2.22× |
| | **EfficientTrain** | 90 | ✗ | **77.0%** | **3.21×** |
| DeiT-Tiny [69] | Progressive Learning [67] *(ICML'21)* | 350 | ✓ | 78.4% | 1.21× |
| | **EfficientTrain** | 300 | ✓ | **79.4%** | **1.44×** |
| | Auto Progressive Learning [39] *(CVPR'22)* | 300 | ✓ | 72.4% | 1.51× |
| | **EfficientTrain** | 300 | ✓ | **73.3%** | **1.55×** |
| DeiT-Small [69] | Progressive Learning [67] *(ICML'21)* | 100 | ✓ | 72.6% | 1.54× |
| | Auto Progressive Learning [39] *(CVPR'22)* | 100 | ✓ | 74.4% | 1.41× |
| | **EfficientTrain** | 100 | ✓ | **76.4%** | 1.51× |
| | Budgeted Training† [40] *(ICLR'20)* | 225 | ✓ | 79.6% | 1.33× |
| | Progressive Learning† [67] *(ICML'21)* | 300 | ✓ | 79.5% | 1.49× |
| | Auto Progressive Learning [39] *(CVPR'22)* | 300 | ✓ | 79.8% | 1.42× |
| | DeiT III [70] *(ECCV'22)* | 300 | ✓ | 79.9% | 1.00× |
| | **EfficientTrain** | 300 | ✓ | **80.3%** | 1.51× |
| CSWin-Tiny [14] | Progressive Learning† [67] *(ICML'21)* | 300 | ✓ | 82.3% | 1.51× |
| | **EfficientTrain** | 300 | ✓ | **82.8%** | **1.55×** |
| | FixRes† [72] *(NeurIPS'19)* | 300 | ✓ | 82.9% | 1.00× |
| | **EfficientTrain+FixRes** [72] *(NeurIPS'19)* | 300 | ✓ | **83.1%** | **1.55×** |
| CSWin-Small [14] | Progressive Learning† [67] *(ICML'21)* | 300 | ✓ | 83.3% | 1.48× |
| | **EfficientTrain** | 300 | ✓ | **83.6%** | **1.51×** |
| | FixRes† [72] *(NeurIPS'19)* | 300 | ✓ | 83.7% | 1.00× |
| | **EfficientTrain+FixRes** [72] *(NeurIPS'19)* | 300 | ✓ | **83.8%** | **1.51×** |

Table 8. **EfficientTrain v.s. state-of-the-art efficient training algorithms on ImageNet-1K.** Here 'AugRegs' denotes the holistic combination of various model regularization and data augmentation techniques, which is widely applied to train deep networks effectively [14, 44, 45, 69, 79]. In particular, some baselines are not developed on top of this state-of-the-art training pipeline. To ensure a fair comparison with them, we also implement our method without AugRegs (notably, the 'RandAug' in Table 4 is removed in this scenario as well). †: our reproduced baselines.

our implementation for both our method and the baselines. More implementation details can be found in Appendix A.

## 5.1. Supervised Learning

**Training various visual backbones on ImageNet-1K.** Table 5 presents the results of applying our method to train representative deep networks on ImageNet-1K. EfficientTrain achieves consistent improvements across different models, *i.e.*, it reaches a competitive or better validation accuracy compared to the baselines (*e.g.*, 83.6% v.s. 83.4% on the CSWin-Small network), while saving the training cost by $1.5 - 1.6×$. Importantly, the practical speedup is largely consistent with the theoretical results. The detailed training runtime (GPU-hours) is deferred to Appendix C.2.

**Adaptation to varying training epochs.** The Efficient-Train curriculum can be conveniently adapted to a varying length of training schedule, *i.e.*, by simply scaling the indices of the training epochs in Table 4. As shown in Table 6 (a), the advantage of EfficientTrain is even more significant with fewer training epochs. For example, it outperforms the baseline by 0.9% (76.4% v.s. 75.5%) on the 100-epoch trained DeiT-Small (the speedup is the same as 300 epochs). We attribute this to the greater importance of efficient training algorithms in the scenarios of limited training resources.

Tables 6 (b) and 6 (c) further compare EfficientTrain and the baselines from the lens of different training epochs. Our method improves the accuracy significantly when leveraging the same training wall-time as the baselines, especially for smaller models (*e.g.*, by 1.8% for DeiT-Tiny). In addition, we find that larger models can be trained effectively with EfficientTrain using less epochs, resulting in $> 2×$ training speedup with an acceptable accuracy drop.

**ImageNet-22K pre-training.** One of the important advantages of modern visual backbones is their excellent scalability with the growing size of training data [14, 45]. To this end, we further verify the effectiveness of EfficientTrain on the larger ImageNet-22K benchmark dataset. The results are provided in Table 7, where the models are pre-trained on ImageNet-22K, and evaluated by being fine-tuned to ImageNet-1K. EfficientTrain is implemented at the pre-training stage, which accounts for the vast majority of the total computation/time cost. One can observe that, similar to ImageNet-1K, our method performs at least on par with the baselines on top of both ConvNets and vision Transformers, while achieving a significant training speedup. A highlight from the results is that EfficientTrain saves a considerable amount of real training time, *e.g.*, 162 GPU-days (307.7 v.s. 469.5) for CSWin-Large, which correspond to ∼20 days for a standard computational node with 8 GPUs.

**Comparisons with state-of-the-art efficient training methods** are summarized in Table 8. EfficientTrain is compared to the recently proposed sample-wise [23, 38, 96, 97] and regularization-wise [13, 62] curriculum learning approaches, as well as the progressive learning algorithms [39, 67]. It can be observed that our method outperforms all these competitive baselines in terms of either the accuracy or the training speedup. Moreover, the simplicity of Effi-

| Method | #Param. | Pre-training Epochs | Top-1 Accuracy (fine-tuning) Baseline | EfficientTrain | Pre-training Speedup Computation | Wall-time |
|---|---|---|---|---|---|---|
| MAE (ViT-B) [25] | 86M | 1600 | 83.5%† | **83.6%** | 1.54× | **1.52×** |
| MAE (ViT-L) [25] | 307M | 1600 | **85.9%†** | 85.8% | 1.53× | **1.55×** |

Table 9. **Self-supervised learning results on top of MAE [25]**. Following [25], the models are pre-trained on ImageNet-1K w/ or w/o EfficientTrain, and evaluated by end-to-end fine-tuning. †: our reproduced baseline.

| Backbone | Pre-training | Top-1 Accuracy (fine-tuned to downstream datasets) CIFAR-10 | CIFAR-100 | Flowers-102 | Stanford Dogs |
|---|---|---|---|---|---|
| DeiT-Small [69] | Baseline | 98.39% | 88.65% | 96.57% | 90.72% |
| | **EfficientTrain** | **98.47%** | **88.93%** | **96.62%** | **91.12%** |

Table 10. **Transferability to downstream image recognition tasks.** The backbone model is pre-trained on ImageNet-1K w/ or w/o EfficientTrain, and fine-tuned to the downstream datasets to report the accuracy.

| Method | Backbone | Pre-training | $AP^{box}$ | $AP^{box}_{50}$ | $AP^{box}_{75}$ | $AP^{mask}$ | $AP^{mask}_{50}$ | $AP^{mask}_{75}$ |
|---|---|---|---|---|---|---|---|---|
| RetinaNet [41] (1× schedule) | Swin-Tiny [44] | Baseline | 41.7 | 62.9 | 44.5 | – | – | – |
| | | **EfficientTrain** | **41.8** | 63.3 | 44.6 | – | – | – |
| Mask-RCNN [26] (1× schedule) | Swin-Tiny [44] | Baseline | 43.3 | 66.2 | 47.2 | **39.6** | 63.0 | 42.4 |
| | | **EfficientTrain** | **43.4** | 66.2 | 47.7 | 39.5 | 62.8 | 42.3 |
| Cascade Mask-RCNN [4] (1× schedule) | Swin-Tiny [44] | Baseline | 48.1 | 67.0 | 52.1 | 41.5 | 64.2 | 44.9 |
| | | **EfficientTrain** | **48.2** | 67.5 | 52.3 | **41.8** | 64.6 | 45.0 |
| | Swin-Small [44] | Baseline | 50.0 | 69.1 | 54.4 | 43.1 | 66.3 | 46.3 |
| | | **EfficientTrain** | **50.6** | 69.7 | 55.2 | **43.6** | 66.9 | 47.3 |
| | Swin-Base [44] | Baseline | 50.9 | 70.2 | 55.5 | 44.0 | 67.4 | 47.4 |
| | | **EfficientTrain** | **51.3** | 70.6 | 56.0 | **44.2** | 67.7 | 47.7 |
| Cascade Mask-RCNN [4] (3× schedule) | Swin-Small [44] | Baseline | 52.0 | 70.7 | 56.3 | 44.9 | 68.2 | 48.9 |
| | | **EfficientTrain** | **52.0** | 70.8 | 56.5 | **44.9** | 68.2 | 48.8 |
| | Swin-Base [44] | Baseline | 52.1 | 71.0 | 56.5 | 45.0 | 68.3 | 48.6 |
| | | **EfficientTrain** | **52.2** | 70.9 | 56.8 | **45.3** | 68.5 | 49.0 |

Table 11. **Object detection and instance segmentation on COCO.** We implement several representative detection/segmentation algorithms on top of the backbones pre-trained w/ or w/o EfficientTrain on ImageNet-1K.

| Bicubic Image Down-sampling | Low-frequency Cropping | Linear RandAug | Training Speedup | Top-1 Accuracy DeiT-Small | Swin-Tiny | Swin-Small |
|---|---|---|---|---|---|---|
| | | | 1.0× | 80.3% | 81.3% | 83.1% |
| | | ✓ | 1.0× | **80.5%** | **81.4%** | 83.2% |
| | ✓ | | ~**1.5×** | 80.0% | 81.1% | 83.1% |
| ✓ | | ✓ | ~**1.5×** | 80.3% | 81.0% | 83.0% |
| | ✓ | ✓ | ~**1.5×** | 80.3% | 81.3% | **83.2%** |

Table 12. **Ablation study of EfficientTrain**, *i.e.*, ablating or replacing the low-frequency cropping and the linearly increased magnitude of RandAug.

cientTrain by fine-tuning them on downstream classification datasets. The results are reported in Table 10. Notably, here the 32×32 images in CIFAR-10/100 [36] are resized to 224×224 for pre-processing following [69], such that the discriminative patterns are mainly distributed within the lower-frequency components. On the contrary, Flowers-102 [52] and Stanford Dogs [34] are fine-grained visual recognition datasets where the high-frequency clues contain important discriminative information (*e.g.*, the detailed facial/skin characteristics to distinguish between the species of dogs). One can observe that EfficientTrain yields competitive or better transfer learning performance than the baselines on both types of datasets. In other words, although our method learns to exploit the lower/higher-frequency information via an ordered curriculum, the finally obtained models are able to leverage both types of information effectively.

**Object detection & instance segmentation on COCO.** To further investigate transferring our pre-trained models to more complex computer vision tasks, we initialize the backbones of several representative detection and segmentation algorithms with the models pre-trained using EfficientTrain. The results are reported in Table 11. EfficientTrain performs at least on par with the baselines, despite the significantly reduced pre-training cost (see: Table 5).
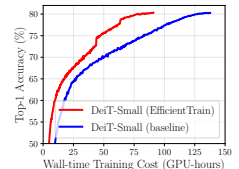
## 5.4. Discussions

**Ablation study** results are summarized in Table 12. The major gain of training efficiency comes from our low-frequency cropping operation, which effectively reduces the training cost at the price of a slight drop of accuracy. Adopting weaker RandAug at earlier epochs moderately improves the performance on top of it. Besides, replacing the low-frequency cropping with the bicubic image down-sampling degrades the final accuracy, since down-sampling cannot filter out all the higher-frequency information (see: Proposition 1), leading to a sub-optimal implementation of our idea.

**The accuracy during training** is shown in Figure 7. The horizontal axis denotes the wall-time training cost. The low-frequency cropping in EfficientTrain is performed on both the training and test inputs. Our approach is able to learn discriminative representations more efficiently at early epochs.



Figure 7. **Curves of accuracy during training**.

**From-scratch training with EfficientTrain on other datasets** is discussed in Appendix C.3.

cientTrain enables it to be conveniently adapted to different models and training settings, which may be non-trivial for other methods (*e.g.*, the network-growing method in [39]).

**Compatibility with FixRes.** FixRes [72] reveals that there exists a discrepancy in the scale of images between the training and test inputs, and thus the inference with a larger resolution will yield a better test accuracy. However, EfficientTrain does not leverage the gains of FixRes, despite the smaller input size at earlier training stages. Importantly, EfficientTrain adopts the original inputs (*e.g.*, $224^2$) at the final training stage, and hence the finally-obtained model resembles the $224^2$-trained networks, while FixRes is orthogonal to it. This fact can be confirmed by both the direct empirical evidence in Table 8 (see: FixRes v.s. EfficientTrain + FixRes on top of the state-of-the-art CSWin Transformers [14]), and the results in Table 7 (see: Input Size=$384^2$).

## 5.2. Self-supervised Learning

**Results with Masked Autoencoders (MAE).** Since our method only modifies the pre-processing of the model inputs during training, it can also be conveniently applied to self-supervised learning algorithms. As a representative example, we investigate deploying EfficientTrain on top of MAE [25] in Table 9. Our method reduces the training cost significantly while preserving a competitive accuracy.

## 5.3. Transfer Learning

**Downstream image recognition tasks.** We first evaluate the transferability of the models trained with Effi-

# 6. Conclusion

This paper investigated a novel generalized curriculum learning approach. The proposed algorithm, *EfficientTrain*, always leverages all the data at any training stage, but only exposes the 'easier-to-learn' patterns of each example at the beginning of training, and gradually introduces more difficult patterns as learning progresses. Our method significantly improves the training efficiency of state-of-the-art deep networks on the large-scale ImageNet-1K / 22K datasets, for both supervised and self-supervised learning.

# Acknowledgements

# References

[1] Guillaume Alain, Alex Lamb, Chinnadhurai Sankar, Aaron Courville, and Yoshua Bengio. Variance reduction in sgd by distributed importance sampling. *arXiv preprint arXiv:1511.06481*, 2015. 3

[2] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *ICML*, pages 41–48, 2009. 1, 2, 3

[3] Stefan Braun, Daniel Neil, and Shih-Chii Liu. A curriculum learning method for improved noise robustness in automatic speech recognition. In *EUSIPCO*, pages 548–552, 2017. 3

[4] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: high quality object detection and instance segmentation. *TPAMI*, 43(5):1483–1498, 2019. 8, 13

[5] Fergus W Campbell and John G Robson. Application of fourier analysis to the visibility of gratings. *The Journal of physiology*, 197(3):551, 1968. 3

[6] Tianqi Chen, Ian Goodfellow, and Jonathon Shlens. Net2net: Accelerating learning via knowledge transfer. In *ICLR*, 2015. 2

[7] Xinlei Chen and Abhinav Gupta. Webly supervised learning of convolutional networks. In *ICCV*, pages 1431–1439, 2015. 2

[8] Yunpeng Chen, Haoqi Fan, Bing Xu, Zhicheng Yan, Yannis Kalantidis, Marcus Rohrbach, Shuicheng Yan, and Jiashi Feng. Drop an octave: Reducing spatial redundancy in convolutional neural networks with octave convolution. In *ICCV*, pages 3435–3444, 2019. 3

[9] Corinna Cortes, Mehryar Mohri, and Afshin Rostamizadeh. Algorithms for learning kernels based on centered alignment. *JMLR*, 13:795–828, 2012. 5

[10] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation strategies from data. In *CVPR*, pages 113–123, 2019. 5

[11] Ekin Dogus Cubuk, Barret Zoph, Jon Shlens, and Quoc Le. Randaugment: Practical automated data augmentation with a reduced search space. In *NeurIPS*, pages 18613–18624, 2020. 6, 12

[12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009. 2, 6, 12, 13, 15, 16

[13] Ürün Dogan, Aniket Anand Deshmukh, Marcin Bronislaw Machura, and Christian Igel. Label-similarity curriculum learning. In *ECCV*, pages 174–190, 2020. 7

[14] Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Weiming Zhang, Nenghai Yu, Lu Yuan, Dong Chen, and Baining Guo. Cswin transformer: A general vision transformer backbone with cross-shaped windows. In *CVPR*, 2022. 2, 5, 6, 7, 8, 12, 13, 16

[15] 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. In *ICLR*, 2021. 1, 15

[16] Jeffrey L Elman. Learning and development in neural networks: The importance of starting small. *Cognition*, 48(1):71–99, 1993. 2, 3

[17] Yang Fan, Fei Tian, Tao Qin, Xiang-Yang Li, and Tie-Yan Liu. Learning to teach. In *ICLR*, 2018. 1, 2, 3

[18] Yoav Freund, Robert E Schapire, et al. Experiments with a new boosting algorithm. In *ICML*, pages 148–156, 1996. 3

[19] Chen Gong, Dacheng Tao, Stephen J Maybank, Wei Liu, Guoliang Kang, and Jie Yang. Multi-modal curriculum learning for semi-supervised image classification. *TIP*, 25(7):3249–3260, 2016. 1

[20] Linyuan Gong, Di He, Zhuohan Li, Tao Qin, Liwei Wang, and Tieyan Liu. Efficient training of bert by progressively stacking. In *ICML*, pages 2337–2346, 2019. 2

[21] Siddharth Gopal. Adaptive sampling for sgd by exploiting side information. In *ICML*, pages 364–372, 2016. 3

[22] Alex Graves, Marc G Bellemare, Jacob Menick, Remi Munos, and Koray Kavukcuoglu. Automated curriculum learning for neural networks. In *ICML*, pages 1311–1320, 2017. 1, 2, 3

[23] Sheng Guo, Weilin Huang, Haozhi Zhang, Chenfan Zhuang, Dengke Dong, Matthew R Scott, and Dinglong Huang. Curriculumnet: Weakly supervised learning from large-scale web images. In *ECCV*, pages 135–150, 2018. 7

[24] Guy Hacohen and Daphna Weinshall. On the power of curriculum learning in training deep networks. In *ICML*, pages 2535–2544, 2019. 1, 2, 3

[25] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, pages 16000–16009, 2022. 2, 8

[26] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, pages 2961–2969, 2017. 8

[27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 1, 2, 6, 7, 12, 16

[28] Gao Huang, Zhuang Liu, Geoff Pleiss, Laurens Van Der Maaten, and Kilian Weinberger. Convolutional networks with dense connectivity. *TPAMI*, pages 1–1, 2019. 1

[29] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *ECCV*, pages 646–661, 2016. 12, 13

[30] Lu Jiang, Deyu Meng, Teruko Mitamura, and Alexander G Hauptmann. Easy samples first: Self-paced reranking for zero-example multimedia search. In *ACM MM*, pages 547–556, 2014. 2

[31] Lu Jiang, Deyu Meng, Qian Zhao, Shiguang Shan, and Alexander G Hauptmann. Self-paced curriculum learning. In *AAAI*, 2015. 1

[32] Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *ICML*, pages 2304–2313, 2018. 1, 2, 3

[33] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *ICLR*, 2018. 2

[34] Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Fei-Fei Li. Novel dataset for fine-grained image categorization: Stanford dogs. In *CVPRW*, 2011. 6, 8, 15, 16

[35] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *ICML*, pages 3519–3529, 2019. 5

[36] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 6, 8, 15, 16

[37] Kai A Krueger and Peter Dayan. Flexible shaping: How learning in small steps helps. *Cognition*, 110(3):380–394, 2009. 2, 3

[38] M Kumar, Benjamin Packer, and Daphne Koller. Self-paced learning for latent variable models. In *NeurIPS*, 2010. 2, 3, 7

[39] Changlin Li, Bohan Zhuang, Guangrun Wang, Xiaodan Liang, Xiaojun Chang, and Yi Yang. Automated progressive learning for efficient training of vision transformers. In *CVPR*, 2022. 1, 2, 7, 8

[40] Mengtian Li, Ersin Yumer, and Deva Ramanan. Budgeted training: Rethinking deep neural network training under resource constraints. In *ICLR*, 2020. 7

[41] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, pages 2980–2988, 2017. 8, 13

[42] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, pages 740–755, 2014. 6

[43] Fengbei Liu, Yu Tian, Yuanhong Chen, Yuyuan Liu, Vasileios Belagiannis, and Gustavo Carneiro. Acpl: Anti-curriculum pseudo-labelling for semi-supervised medical image classification. In *CVPR*, pages 20697–20706, 2022. 3

[44] 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 *ICCV*, pages 10012–10022, 2021. 1, 2, 5, 6, 7, 8, 12, 14, 16

[45] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *CVPR*, 2022. 2, 5, 6, 7, 12, 13, 16

[46] Raphael Gontijo Lopes, Dong Yin, Ben Poole, Justin Gilmer, and Ekin D Cubuk. Improving robustness without sacrificing accuracy with patch gaussian augmentation. *arXiv preprint arXiv:1906.02611*, 2019. 2

[47] Ilya Loshchilov and Frank Hutter. Online batch selection for faster training of neural networks. *arXiv preprint arXiv:1511.06343*, 2015. 3

[48] Stéphane Mallat. *A wavelet tour of signal processing*. Elsevier, 1999. 3

[49] Tambet Matiisen, Avital Oliver, Taco Cohen, and John Schulman. Teacher–student curriculum learning. *IEEE transactions on neural networks and learning systems*, 31(9):3732–3740, 2019. 2

[50] Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed precision training. In *ICLR*, 2018. 12

[51] Pietro Morerio, Jacopo Cavazza, Riccardo Volpi, René Vidal, and Vittorio Murino. Curriculum dropout. In *ICCV*, pages 3544–3552, 2017. 2

[52] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, pages 722–729, 2008. 6, 8, 15, 16

[53] Sayak Paul and Pin-Yu Chen. Vision transformers are robust learners. *arXiv preprint arXiv:2105.07581*, 2(3), 2021. 2

[54] Te Pi, Xi Li, Zhongfei Zhang, Deyu Meng, Fei Wu, Jun Xiao, and Yueting Zhuang. Self-paced boost learning for classification. In *IJCAI*, pages 1932–1938, 2016. 3

[55] Emmanouil Antonios Platanios, Otilia Stretcu, Graham Neubig, Barnabás Póczos, and Tom M Mitchell. Competence-based curriculum learning for neural machine translation. In *NAACL-HLT (1)*, 2019. 2

[56] Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging. *SIAM journal on control and optimization*, 30(4):838–855, 1992. 12

[57] Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? In *NeurIPS*, 2021. 13

[58] Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for robust deep learning. In *ICML*, pages 4334–4343, 2018. 2

[59] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the masses. *arXiv preprint arXiv:2104.10972*, 2021. 13

[60] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. In *CVPR*, pages 761–769, 2016. 3

[61] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2014. 2

[62] Samarth Sinha, Animesh Garg, and Hugo Larochelle. Curriculum by smoothing. In *NeurIPS*, pages 21653–21664, 2020. 2, 7

[63] Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu Sebe. Curriculum learning: A survey. *International Journal of Computer Vision*, pages 1–40, 2022. 1, 2, 3

[64] Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for deep learning in nlp. In *ACL*, pages 3645–3650, 2019. 1

[65] Wim Sweldens. The lifting scheme: A construction of second generation wavelets. *SIAM journal on mathematical analysis*, 29(2):511–546, 1998. 3

[66] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, pages 2818–2826, 2016. 12

[67] Mingxing Tan and Quoc Le. Efficientnetv2: Smaller models and faster training. In *ICML*, pages 10096–10106, 2021. 2, 7

[68] Keyu Tian, Chen Lin, Ming Sun, Luping Zhou, Junjie Yan, and Wanli Ouyang. Improving auto-augment via augmentation-wise weight sharing. In *NeurIPS*, pages 19088–19098, 2020. 6

[69] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, pages 10347–10357, 2021. 2, 4, 5, 6, 7, 8, 12, 15, 16

[70] Hugo Touvron, Matthieu Cord, and Herve Jegou. Deit iii: Revenge of the vit. In *ECCV*, 2022. 7

[71] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *ICCV*, pages 32–42, 2021. 12

[72] Hugo Touvron, Andrea Vedaldi, Matthijs Douze, and Hervé Jégou. Fixing the train-test resolution discrepancy. In *NeurIPS*, 2019. 2, 7, 8

[73] Radu Tudor Ionescu, Bogdan Alexe, Marius Leordeanu, Marius Popescu, Dim P Papadopoulos, and Vittorio Ferrari. How hard can it be? estimating the difficulty of visual search in an image. In *CVPR*, pages 2157–2166, 2016. 2

[74] Jonathan G Tullis and Aaron S Benjamin. On the effectiveness of self-paced learning. *Journal of memory and language*, 64(2):109–118, 2011. 2

[75] Guangcong Wang, Xiaohua Xie, Jianhuang Lai, and Jiaxuan Zhuo. Deep growing learning. In *ICCV*, pages 2812–2820, 2017. 2

[76] Haohan Wang, Xindi Wu, Zeyi Huang, and Eric P Xing. High-frequency component helps explain the generalization of convolutional neural networks. In *CVPR*, pages 8684–8694, 2020. 2, 4

[77] Jiachun Wang, Fajie Yuan, Jian Chen, Qingyao Wu, Min Yang, Yang Sun, and Guoxiao Zhang. Stackrec: Efficient training of very deep sequential recommender models by iterative stacking. In *ACM SIGIR*, pages 357–366, 2021. 2

[78] Wei Wang, Isaac Caswell, and Ciprian Chelba. Dynamically composing domain-data selection with clean-data selection by "co-curricular learning" for neural machine translation. In *ACL*, pages 1282–1292, 2019. 3

[79] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *ICCV*, pages 568–578, 2021. 2, 6, 7, 12, 16

[80] Xin Wang, Yudong Chen, and Wenwu Zhu. A survey on curriculum learning. *TPAMI*, 2021. 1, 2, 3

[81] Tao Wei, Changhu Wang, and Chang Wen Chen. Modularized morphing of deep convolutional neural networks: A graph approach. *IEEE Transactions on Computers*, 70(2):305–315, 2020. 2

[82] Tao Wei, Changhu Wang, Yong Rui, and Chang Wen Chen. Network morphism. In *ICML*, pages 564–572, 2016. 2

[83] Yunchao Wei, Xiaodan Liang, Yunpeng Chen, Xiaohui Shen, Ming-Ming Cheng, Jiashi Feng, Yao Zhao, and Shuicheng Yan. Stc: A simple to complex framework for weakly-supervised semantic segmentation. *TPAMI*, 39(11):2314–2320, 2016. 2

[84] Daphna Weinshall, Gad Cohen, and Dan Amir. Curriculum learning by transfer learning: Theory and experiments with deep networks. In *ICML*, pages 5238–5246, 2018. 2

[85] Zhuofan Xia, Xuran Pan, Shiji Song, Li Erran Li, and Gao Huang. Vision transformer with deformable attention. In *CVPR*, 2022. 13

[86] Dong Yin, Raphael Gontijo Lopes, Jon Shlens, Ekin Dogus Cubuk, and Justin Gilmer. A fourier perspective on model robustness in computer vision. In *NeurIPS*, 2019. 2

[87] Yuning You, Tianlong Chen, Zhangyang Wang, and Yang Shen. L2-gcn: Layer-wise and learned efficient training of graph convolutional networks. In *CVPR*, pages 2127–2135, 2020. 2

[88] 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 *ICCV*, pages 6023–6032, 2019. 12

[89] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In *CVPR*, 2022. 1

[90] Dingwen Zhang, Junwei Han, Long Zhao, and Deyu Meng. Leveraging prior-knowledge for weakly supervised object detection under a collaborative self-paced curriculum learning framework. *IJCV*, 127(4):363–380, 2019. 2

[91] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *ICLR*, 2018. 12

[92] Minjia Zhang and Yuxiong He. Accelerating training of transformer-based language models with progressive layer dropping. In *NeurIPS*, pages 14011–14023, 2020. 2

[93] Xuan Zhang, Gaurav Kumar, Huda Khayrallah, Kenton Murray, Jeremy Gwinnup, Marianna J Martindale, Paul McNamee, Kevin Duh, and Marine Carpuat. An empirical exploration of curriculum learning for neural machine translation. *arXiv preprint arXiv:1811.00739*, 2018. 3

[94] Xinyu Zhang, Qiang Wang, Jian Zhang, and Zhao Zhong. Adversarial autoaugment. In *ICLR*, 2019. 5

[95] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *AAAI*, pages 13001–13008, 2020. 12

[96] Tianyi Zhou and Jeff Bilmes. Minimax curriculum learning: Machine teaching with desirable difficulties and scheduled diversity. In *ICLR*, 2018. 3, 7

[97] Tianyi Zhou, Shengjie Wang, and Jeffrey Bilmes. Curriculum learning by dynamic instance hardness. In *NeurIPS*, pages 8602–8613, 2020. 3, 7

# Appendix for
# "EfficientTrain: Exploring Generalized Curriculum Learning for Training Visual Backbones"

## A. Implementation Details

### A.1. Training Models on ImageNet-1K

**Dataset.** We use the data provided by ILSVRC2012[1] [12]. The dataset includes 1.2 million images for training and 50,000 images for validation, both of which are categorized in 1,000 classes.

**Training.** Our approach is developed on top of a state-of-the-art training pipeline of deep networks, which incorporates a holistic combination of various model regularization & data augmentation techniques, and is widely applied to train recently proposed models [14, 44, 45, 69, 79]. Our training settings generally follow from [45], while we modify the configurations of weight decay, stochastic depth and exponential moving average (EMA) according to the recommendation in the original papers of different models (*i.e.*, ConvNeXt [45], DeiT [69], PVT [79], Swin Transformer [44] and CSWin Transformer [14])[2]. The detailed hyper-parameters are summarized in Table 13.

The baselines presented in Table 5 directly use the training configurations in Table 13. Based on Table 13, our proposed EfficientTrain curriculum performs low-frequency cropping and modifies the value of $m$ in RandAug during training, as introduced in Table 4. The results in Tables 6 and 8 adopt a varying number of training epochs on top of Table 5.

In addition, the low-frequency cropping operation in EfficientTrain leads to a varying input size during training. Notably, visual backbones can naturally process different sizes of inputs with no or minimal modifications. Specifically, once the input size varies, ResNets and ConvNeXts do not need any change, while vision Transformers (*i.e.*, DeiT, PVT, Swin and CSWin) only need to resize their position bias correspondingly, as suggested in their papers. Our method starts the training with small-size inputs and the reduced computational cost. The input size is switched midway in the training process, where we resize the position bias for ViTs (do nothing for ConvNets). Finally, the learning ends up with full-size inputs, as used at test time. As a consequence, the overall computational/time cost to obtain the final trained models is effectively saved.

**Inference.** Following [14, 44, 45, 69, 79], we use a crop ratio of 0.875 and 1.0 for the inference input size of $224^2$ and $384^2$, respectively.

| Training Config | Values / Setups |
|---|---|
| Input size | $224^2$ |
| Weight init. | Truncated normal (0.2) |
| Optimizer | AdamW |
| Optimizer hyper-parameters | $\beta_1, \beta_2$=0.9, 0.999 |
| Initial learning rate | 4e-3 |
| Learning rate schedule | Cosine annealing |
| Weight decay | 0.05 |
| Batch size | 4,096 |
| Training epochs | 300 |
| Warmup epochs | 20 |
| Warmup schedule | linear |
| RandAug [11] | (9, 0.5) |
| Mixup [91] | 0.8 |
| Cutmix [88] | 1.0 |
| Random erasing [95] | 0.25 |
| Label smoothing [66] | 0.1 |
| Stochastic depth [29] | Following the values in original papers [14, 44, 45, 69, 79]. |
| Layer scale [71] | 1e-6 (ConvNeXt [45]) / None (others) |
| Gradient clip | 5.0 (DeiT [69], PVT [79] and Swin [44]) / None (others) |
| Exp. mov. avg. (EMA) [56] | 0.9999 (ConvNeXt [45] and CSWin [14]) / None (others) |
| Auto. mix. prec. (AMP) [50] | Inactivated (ConvNeXt [45]) / Activated (others) |

Table 13. **Basic training hyper-parameters for the models in Table 5.**

---

[1] https://image-net.org/index.php
[2] The training of ResNet [27] follows the recipe provided in [45].

## A.2. ImageNet-22K Pre-training

**Dataset and pre-processing.** In our experiments, the officially released processed version of ImageNet-22K[3] [12, 59] is used. The original ImageNet-22K dataset is pre-processed by resizing the images (to reduce the dataset's memory footprint from 1.3TB to ∼250GB) and removing a small number of samples. The processed dataset consists of ∼13M images in ∼19K classes. Note that this pre-processing procedure is officially recommended and accomplished by the official website.

**Pre-training.** We pre-train CSWin-Base/Large and ConvNeXt-Base/Large on ImageNet-22K. The pre-training process basically follows the training configurations of ImageNet-1K (*i.e.*, Table 13), except for the differences presented in the following. The number of training epochs is set to 120 with a 5-epoch linear warm-up. For all the four models, the maximum value of the increasing stochastic depth regularization [29] is set to 0.1 [14, 45]. Following [14], the initial learning rate for CSWin-Base/Large is set to 2e-3, while the weight-decay coefficient for CSWin-Base/Large is set to 0.05/0.1. Following [45], we do not leverage the exponential moving average (EMA) mechanism. To ensure a fair comparison, we report the results of our implementation for both baselines and EfficientTrain, where they adopt exactly the same training settings (apart from the configurations modified by EfficientTrain itself).

**Fine-tuning.** We evaluate the ImageNet-22K pre-trained models by fine-tuning them and reporting the corresponding accuracy on ImageNet-1K. The fine-tuning process of ConvNeXt-Base/Large follows their original paper [45]. The fine-tuning of CSWin-Base/Large adopts the same setups as ConvNeXt-Base/Large. We empirically observe that this setting achieves a better performance than the original fine-tuning pipeline of CSWin-Base/Large in [14].

## A.3. Object Detection and Segmentation on COCO

Our implementation of RetinaNet [41] follows from [85]. Our implementation of Cascade Mask-RCNN [4] is the same as [45].

## A.4. Experiments in Section 4

In particular, the experimental results provided in Section 4 are based on the training settings listed in Table 13 as well, expect for the specified modifications (*e.g.*, with the low-passed filtered inputs). The computing of CKA feature similarity follows [57].

---

[3]https://image-net.org/data/imagenet21k_resized.tar.gz

# B. Algorithm to Search for the Curriculum

To determine a proper schedule for $B$ (*i.e.*, the window size for low-frequency cropping), we propose a greedy search algorithm, which is shown in Algorithm 1. We divide the training process into several stages and solve a value of $B$ for each stage. The algorithm starts from the last stage, minimizing $B$ under the constraint of not degrading the performance compared to the baseline (trained with $B = 224$). In implementation, we consider the standard 300-epoch training setting proposed in [44]. To ensure the generalization performance, we restrict the accuracy in terms of both Swin-Tiny and DeiT-Small.

---

**Algorithm 1** Algorithm to Search for Curricula.

---

1: **Input:** Number of training epochs $T$ and training stages $N$ (*i.e.*, $T/N$ epochs for each stage).
2: **Input:** Baseline accuracy $a_0$ (with $224^2$ images).
3: **To solve:** The value of $B$ for $i^{\text{th}}$ training stage: $\hat{B}_i$.
4: **Initialize:** $\hat{B}_1 = \ldots = \hat{B}_N = 224$
5: **for** $i = N - 1$ **to** $1$ **do**
6:     $\hat{B}_i = \underset{B_1=\ldots=B_i=B,\ B_j=\hat{B}_j,\ i<j\leq N}{\text{minimize}} B,$
    s.t. ValidationAccuracy$(B_1, \ldots, B_N) \geq a_0$
7: **end for**
8: **Output:** $\hat{B}_1, \ldots, \hat{B}_N$

---

# C. Additional Results

## C.1. Additional Evidence for the Frequency-inspired Curricula

**Results with high-pass filtering.** Here we consider an opposite case to Figure 3, where the high-pass filtering is performed. The consequent training curves are shown in Figure 8. One can observe that the baseline outperforms all the variants significantly from the very beginning, while the gap is not effectively closed or shrunk during the whole training process. In addition, such degradation of the performance becomes even more serious when more lower-frequency components are eliminated (*i.e.*, with larger $r$). This phenomenon is favorable to our assumption: due to the
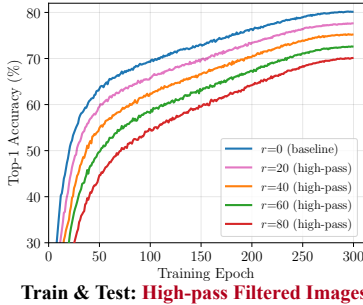


Figure 8. **Validation accuracy v.s. training epochs with high-pass filtered inputs** (DeiT-Small [69] on ImageNet-1K). We preserve the high-frequency components which are eliminated in Figure 3, and remove the low-frequency components that would have been retained.

lack of the necessary lower-frequency information, the beginning stages of training yield less discriminative features, and the remaining training process cannot make up for this deficiency.

**Comparisons of the learning process on low/high-pass filtered data.** To further validate our assumption, we train a model with a regular learning process (*i.e.*, using the original input data without any filtering), and evaluate the intermediate checkpoints using both low-pass and high-pass filtered validation sets. In particular, we adjust the bandwidth of the low/high-pass filter to make the finally trained model have the same accuracy on the two validation set (*i.e.*, 71.3%). In other words, we let the low/high-pass filtered validation set have the same amount of discriminative information, while in the two validation sets, this discriminative information mainly consists of the lower/higher-frequency components, respectively. The intermediate accuracy on the two validation sets is given in Table 14. Although the two final accuracies are equal, the accuracy on the low-pass filtered validation set grows much faster at the beginning stages of training. This favorable observation validates our assumption effectively.

| Training Epoch | $20^{th}$ | $50^{th}$ | $100^{th}$ | $200^{th}$ | $300^{th}$ (end) |
|---|---|---|---|---|---|
| Low-pass Filtered Validation Set | **46.9%** | **58.8%** | **63.1%** | **68.6%** | 71.3% |
| High-pass Filtered Validation Set | 23.9% | 43.5% | 53.5% | 64.3% | 71.3% |

Table 14. **Accuracy of DeiT-Small during the training process on the low/high-pass filtered validation set.**

## C.2. Wall-time Training Cost

The detailed wall-time training cost for the models presented in Table 5 of the paper is reported in Table 15. The numbers of GPU-hours are benchmarked on NVIDIA 3090 GPUs. The batch size for each GPU and the total number of GPUs are configured conditioned on different models, under the principle of saturating all the computational cores of GPUs.

## C.3. From-scratch training with EfficientTrain on other datasets

**Datasets.** We first summarize several representative benchmark datasets with various characteristics, as shown below.

- ImageNet-1K [12]: The large-scale natural image datasets. The discriminative information is distributed in both low-frequency and high-frequency components.

- CIFAR-10/100 [36]: Containing 32x32 colored natural images. Following the common practice in [15, 69], we resize the data to 224x224. Due to this operation, the discriminative information is distributed within the low-frequency components (mathematically, the high-frequency parts in the Fourier spectrum are empty).

- Flowers-102 [52] and Stanford Dogs [34]: Fine-grained recognition datasets. The high-frequency clues contain important discriminative information (*e.g.*, the detailed facial/skin characteristics to distinguish between the species of dogs).

**Pre-training + fine-tuning.** In the context of computer vision, the ubiquitous common practice is to first pre-train the models on large-scale datasets and then fine-tune them on the downstream datasets of interest. EfficientTrain focuses on the pre-training stage which contributes to the vast majority of the total training cost. The large-scale datasets for pre-training are scarce, and ImageNet-1K/22K is widely used in most cases. In this sense, we believe that our results on ImageNet-1K/22K may be sufficient enough to represent most practical scenarios.

| | Model | Input Size (inference) | #Param. | #FLOPs | Top-1 Accuracy (300 epochs) | | Wall-time Training Cost (in GPU-hours) | | Speedup |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | Baseline | EfficientTrain | Baseline | EfficientTrain | |
| *ConvNets* | ResNet-50 [27] | $224^2$ | 26M | 4.1G | 78.8% | **79.4%** | 205.2 | **142.2** | **1.44×** |
| | ConvNeXt-Tiny [45] | $224^2$ | 29M | 4.5G | 82.1% | **82.2%** | 379.5 | **254.1** | **1.49×** |
| | ConvNeXt-Small [45] | $224^2$ | 50M | 8.7G | 83.1% | **83.2%** | 673.5 | **449.9** | **1.50×** |
| | ConvNeXt-Base [45] | $224^2$ | 89M | 15.4G | **83.8%** | 83.7% | 997.0 | **671.6** | **1.48×** |
| *Isotropic ViTs* | DeiT-Tiny [69] | $224^2$ | 5M | 1.3G | 72.5% | **73.3%** | 60.5 | **38.9** | **1.55×** |
| | DeiT-Small [69] | $224^2$ | 22M | 4.6G | 80.3% | **80.3%** | 137.7 | **90.9** | **1.51×** |
| *Multi-stage ViTs* | PVT-Tiny [79] | $224^2$ | 13M | 1.9G | 75.5% | **75.5%** | 99.0 | **66.8** | **1.48×** |
| | PVT-Small [79] | $224^2$ | 25M | 3.8G | 79.9% | **80.4%** | 201.1 | **129.2** | **1.56×** |
| | PVT-Medium [79] | $224^2$ | 44M | 6.7G | 81.8% | **81.8%** | 310.9 | **208.3** | **1.49×** |
| | PVT-Large [79] | $224^2$ | 61M | 9.8G | 82.3% | **82.3%** | 515.9 | **337.3** | **1.53×** |
| | Swin-Tiny [44] | $224^2$ | 28M | 4.5G | 81.3% | **81.3%** | 232.3 | **155.5** | **1.49×** |
| | Swin-Small [44] | $224^2$ | 50M | 8.7G | 83.1% | **83.2%** | 360.3 | **239.4** | **1.50×** |
| | Swin-Base [44] | $224^2$ | 88M | 15.4G | 83.4% | **83.6%** | 494.5 | **329.9** | **1.50×** |
| | CSWin-Tiny [14] | $224^2$ | 23M | 4.3G | 82.7% | **82.8%** | 290.1 | **187.5** | **1.55×** |
| | CSWin-Small [14] | $224^2$ | 35M | 6.9G | 83.4% | **83.6%** | 438.5 | **291.0** | **1.51×** |
| | CSWin-Base [14] | $224^2$ | 78M | 15.0G | 84.3% | **84.3%** | 823.7 | **528.2** | **1.56×** |

Table 15. **Accuracy v.s. wall-time training cost for the deep networks trained on ImageNet-1K** (*i.e.*, corresponding to the models presented in Table 5 of the paper).

**Training from scratch.** From the perspective of pure experimental exploration, we also test deploying EfficientTrain to train the models from scratch on other datasets, as presented in Table 16. We find that EfficientTrain may slightly degrade the accuracy on the fine-grained recognition datasets, where the high-frequency details may be important to distinguish between the fine-grained categories with subtle visual differences (*e.g.*, the species of flowers/dogs). The learning process on these datasets may contradict the assumption behind EfficientTrain since here the lower-frequency components include little discriminative information. This observation may motivate future research to explore the efficient learning algorithms on the datasets dominated by high-frequency information.

| | ImageNet-1K [12] | CIFAR-10 [36] | CIFAR-100 [36] | Oxford Flowers-102 [52] | Stanford Dogs [34] |
|---|---|---|---|---|---|
| Baseline | 80.3% | 87.55% | 65.47% | **52.41%** | **37.31%** |
| EfficientTrain | **80.3%** | **88.92%** | **67.61%** | 51.01% | 36.06% |

Table 16. **Accuracy of DeiT-Small when it is trained on the corresponding datasets from scratch.**

# D. Theoretical Deduction

In this section, we theoretically demonstrate the difference between two transformations, namely the low-frequency cropping and the image down-sampling. In specific, we will show that from the perspective of signal processing, the former perfectly preserves the lower-frequency signals within a square region in the frequency domain and discards the rest, while the image obtained from the pixel-space down-sampling contains the signals mixed from both lower- and higher- frequencies.

## D.1. Preliminaries

An image can be seen as a high-dimensional data point $\boldsymbol{X} \in \mathbb{R}^{C_0 \times H_0 \times W_0}$, where $C_0, H_0, W_0$ represent the number of channels, height and width. Since each channel's signals are regarded as independent, for the sake of simplicity, we can focus on a single-channel image with even edge length $\boldsymbol{X} \in \mathbb{R}^{2H \times 2W}$. Denote the 2D discrete Fourier transform as $\mathcal{F}(\cdot)$. Without loss of generality, we assume that the coordinate ranges are $\{-H, -H+1, \ldots, H-1\}$ and $\{-W, -W+1, \ldots, W-1\}$. The value of the pixel at the position $[u, v]$ in the frequency map $\boldsymbol{F} = \mathcal{F}(\boldsymbol{X})$ is computed by

$$\boldsymbol{F}[u, v] = \sum_{x=-H}^{H-1} \sum_{y=-W}^{W-1} \boldsymbol{X}[x, y] \cdot \exp\left(-j2\pi \left(\frac{ux}{2H} + \frac{vy}{2W}\right)\right).$$

Similarly, the inverse 2D discrete Fourier transform $\boldsymbol{X} = \mathcal{F}^{-1}(\boldsymbol{F})$ is defined by

$$\boldsymbol{X}[x, y] = \frac{1}{HW} \sum_{u=-H}^{H-1} \sum_{v=-W}^{W-1} \boldsymbol{F}[u, v] \cdot \exp\left(j2\pi \left(\frac{ux}{2H} + \frac{vy}{2W}\right)\right).$$

Denote the low-frequency cropping operation parametrized by output size $(2H', 2W')$ as $\mathcal{C}_{H',W'}(\cdot)$, which gives output by simple cropping:

$$\mathcal{C}_{H',W'}(\boldsymbol{F})[u, v] = \frac{H'W'}{HW} \cdot \boldsymbol{F}[u, v].$$

Note that here $u \in \{-H, -H+1, \ldots, H-1\}, v \in \{-W, -W+1, \ldots, W-1\}$, the operation just copy the central area of $\boldsymbol{F}$ with a scaling ratio. The scale ratio $\frac{H'W'}{HW}$ is a natural term from the change of total energy in the pixels, since the number of pixels shrinks by a ratio of $\frac{H'W'}{HW}$.

Also, denote the down-sampling operation parametrized by ratio $r \in (0, 1]$ as $\mathcal{D}_r(\cdot)$. For simplicity, we first consider the case where $r = \frac{1}{k}$ for an integer $k \in \mathbb{N}^+$, and then extend our conclusions to the general cases where $r \in (0, 1]$. In real applications, there are many different down-sampling strategy using different interpolation method, *i.e.*, nearest, bilinear, bicubic, etc. When $k$ is an integer, these operations can be modeled as *using a constant convolution kernel to aggregate the neighborhood pixels*. Denote this kernel's parameter as $\boldsymbol{w}_{s,t}$ where $s, t \in \{0, 1, \ldots, k-1\}$ and $\sum_{s=0}^{k-1} \sum_{t=0}^{k-1} \boldsymbol{w}_{s,t} = \frac{1}{k^2}$. Then the down-sampling operation can be represented as

$$\mathcal{D}_{1/k}(\boldsymbol{X})[x', y'] = \sum_{s=0}^{k-1} \sum_{t=0}^{k-1} \boldsymbol{w}_{s,t} \cdot \boldsymbol{X}[kx' + s, ky' + t].$$

## D.2. Propositions

Now we are ready to demonstrate the difference between the two operations and prove our claims. We start by considering shrinking the image size by $k$ and $k$ is an integer. Here the low-frequency region of an image $\boldsymbol{X} \in \mathbb{R}^{H \times W}$ refers to the signals within range $[-H/k, H/k - 1] \times [-W/k, W/k - 1]$ in $\mathcal{F}(\boldsymbol{X})$, while the rest is named as the high-frequency region.

**Proposition 1.1.** *Suppose that the original image is $\boldsymbol{X}$, and that the image generated from the low-frequency cropping operation is $\boldsymbol{X}_c = \mathcal{F}^{-1} \circ \mathcal{C}_{H/k, W/k} \circ \mathcal{F}(\boldsymbol{X}), k \in \mathbb{N}^+$. We have that all the signals in the spectral map of $\boldsymbol{X}_c$ is only from the low frequency region of $\boldsymbol{X}$, while we can always recover $\mathcal{C}_{H/k, W/k} \circ \mathcal{F}(\boldsymbol{X})$ from $\boldsymbol{X}_c$.*

**Proof.** The proof of this proposition is simple and straightforward. Take the Fourier transform on both sides of above transformation equation, we get

$$\mathcal{F}(\boldsymbol{X}_c) = \mathcal{C}_{H/k, W/k} \circ \mathcal{F}(\boldsymbol{X}).$$

Denote the spectral of $\boldsymbol{X}_c$ as $\boldsymbol{F}_c = \mathcal{F}(\boldsymbol{X}_c)$ and similarly $\boldsymbol{F} = \mathcal{F}(\boldsymbol{X})$. According to our definition of the cropping operation, we know

$$\boldsymbol{F}_c[u,v] = \frac{H/k \cdot W/k}{HW} \cdot \boldsymbol{F}[u,v] = \frac{1}{k^2} \cdot \boldsymbol{F}[u,v].$$

Hence the spectral information of $\boldsymbol{X}_c$ is simply copying $\boldsymbol{X}$'s low frequency part and conduct a uniform scaling by dividing $k^2$.

**Proposition 1.2.** *Suppose that the original image is $\boldsymbol{X}$, and that the image generated from down-sampling operation is $\boldsymbol{X}_d = \mathcal{D}_{1/k}(\boldsymbol{X}), k \in \mathbb{N}^+$. We have that signals in the spectral map of $\boldsymbol{X}_d$ have non-zero dependency on the high frequency region of $\boldsymbol{X}$.*

**Proof.** Taking the Fourier transform on both sides, we have

$$\mathcal{F}(\boldsymbol{X}_d) = \mathcal{F}(\mathcal{D}_{1/k}(\boldsymbol{X})).$$

For any $u \in [-H/k, H/k-1], v \in [-W/k, W/k-1]$, according to the definition we have

$$\mathcal{F}(\boldsymbol{X}_d)[u,v] = \sum_{x=-H/k}^{H/k-1} \sum_{y=-W/k}^{W/k-1} \mathcal{D}_{1/k}(\boldsymbol{X})[x,y] \cdot \exp\left(-j2\pi\left(\frac{kux}{2H} + \frac{kvy}{2W}\right)\right)$$

$$= \sum_{x=-H/k}^{H/k-1} \sum_{y=-W/k}^{W/k-1} \sum_{s=0}^{k-1}\sum_{t=0}^{k-1} \boldsymbol{w}_{s,t} \cdot \boldsymbol{X}[kx+s, ky+t] \cdot \exp\left(-j2\pi\left(\frac{kux}{2H} + \frac{kvy}{2W}\right)\right), \qquad (*1)$$

while at the same time we have the inverse FFT for $\boldsymbol{X}$

$$\boldsymbol{X}[x,y] = \frac{1}{2H \cdot 2W} \sum_{u'=-H}^{H-1} \sum_{v'=-W}^{W-1} \boldsymbol{F}[u',v'] \cdot \exp\left(j2\pi\left(\frac{u'x}{2H} + \frac{v'y}{2W}\right)\right). \qquad (*2)$$

Plugging (*2) into (*1), it is easy to see that essentially each $\boldsymbol{F}_d(u,v) = \mathcal{F}(\boldsymbol{X}_d)[u,v]$ is a linear combination of the signals in original image $\boldsymbol{F}[u',v']$. Namely, it can be represented as form

$$\boldsymbol{F}_d(u,v) = \sum_{u'=-H}^{H-1} \sum_{v'=-W}^{W-1} \alpha(u,v,u',v') \cdot \boldsymbol{F}(u',v').$$

Therefore, we can compute the dependency weight for any given tuple $(u,v,u',v')$ as

$$\alpha(u,v,u',v') = \frac{1}{4HW} \sum_{x=-H}^{H-1} \sum_{y=-W}^{W-1} \boldsymbol{w}_{x_r,y_r} \cdot \exp\left(-j2\pi\left(\frac{ux_p}{2H} + \frac{vy_p}{2W}\right)\right) \cdot \exp\left(j2\pi\left(\frac{u'x}{2H} + \frac{v'y}{2W}\right)\right)$$

$$= \frac{1}{4HW} \sum_{x=-H}^{H-1} \sum_{y=-W}^{W-1} \boldsymbol{w}_{x_r,y_r} \cdot \exp\left(j2\pi\left(\frac{u'x - ux_p}{2H} + \frac{v'y - vy_p}{2W}\right)\right),$$

where $x_r = x \bmod k, x_p = x - x_r$, same for $y_r, y_p$. Further deduction shows

$$\alpha(u,v,u',v') = \frac{1}{4HW} \sum_{x=-H}^{H-1} \sum_{y=-W}^{W-1} \boldsymbol{w}_{x_r,y_r} \cdot \exp\left(j2\pi\left(\frac{(u'-u)x_p + u'x_r}{2H} + \frac{(v'-v)y_p + v'y_r}{2W}\right)\right)$$

$$= \frac{1}{4HW} \sum_{x'=-H/k}^{H/k-1} \sum_{y'=-W/k}^{W/k-1} \sum_{s=0}^{k-1}\sum_{t=0}^{k-1} \boldsymbol{w}_{s,t} \cdot \exp\left(j2\pi k\left(\frac{(u'-u)x'}{2H} + \frac{(v'-v)y'}{2W}\right)\right) \cdot \exp\left(j2\pi\left(\frac{u's}{2H} + \frac{v't}{2W}\right)\right)$$

$$= \frac{1}{4HW} \sum_{x'=-H/k}^{H/k-1} \sum_{y'=-W/k}^{W/k-1} \exp\left(j2\pi k\left(\frac{(u'-u)x'}{2H} + \frac{(v'-v)y'}{2W}\right)\right) \sum_{s=0}^{k-1}\sum_{t=0}^{k-1} \boldsymbol{w}_{s,t} \cdot \exp\left(j2\pi\left(\frac{u's}{2H} + \frac{v't}{2W}\right)\right).$$

Denote $\beta(u', v') = \sum_{s=0}^{k-1} \sum_{t=0}^{k-1} \boldsymbol{w}_{s,t} \cdot \exp\left(j2\pi\left(\frac{u's}{2H} + \frac{v't}{2W}\right)\right)$, which is a constant conditioned on $(u', v')$. Then we know

$$
\begin{aligned}
\alpha(u, v, u', v') &= \frac{\beta(u', v')}{4HW} \sum_{x'=-H/k}^{H/k-1} \sum_{y=-W/k}^{W/k-1} \exp\left(j2\pi k\left(\frac{(u'-u)x'}{2H} + \frac{(v'-v)y'}{2W}\right)\right) \\
&= \frac{\beta(u', v')}{4HW} \sum_{x'=-H/k}^{H/k-1} \exp\left(j2\pi k\left(\frac{(u'-u)x'}{2H}\right)\right) \sum_{y'=-W/k}^{W/k-1} \exp\left(j2\pi k\left(\frac{(v'-v)y'}{2W}\right)\right) \\
&= \begin{cases} \frac{\beta(u', v')}{k^2}, & u'-u = a \cdot \frac{2H}{k}, v'-v = b \cdot \frac{2W}{k}, \quad a, b \in \mathbb{Z} \\ 0, & \text{otherwise} \end{cases}
\end{aligned}
\tag{*3}
$$

In general, $\beta(u', v') \neq 0$ when $u' \neq c \cdot \frac{2H}{k}, v' \neq d \cdot \frac{2W}{k}, c, d \in \mathbb{Z}, cd \neq 0$. Hence, when $\frac{2H}{k}|(u' - u), \frac{2W}{k}|(v' - v)$, we have $\alpha(u, v, u', v') \neq 0$ given $uv \neq 0$, while $\alpha(u, 0, u', 0) \neq 0, \alpha(0, v, 0, v') \neq 0$. Therefore, the image generated by down-sampling operation contains mixed information from both low frequency and high frequency, since most signals have non-zero dependency on global signals of original image.

**Proposition 1.3.** *The conclusions of Proposition 1.1 and Proposition 1.2 still hold when $k \in \mathbb{Q}^+$ is not an integer.*

**Proof.** It is obvious that Proposition 1.1 can be naturally extended to $k \in \mathbb{Q}^+$. Therefore, here we focus on Proposition 1.2. First, consider up-sample an image $\boldsymbol{X}$ by $m \in \mathbb{N}^+$ with nearest interpolation, namely

$$
\boldsymbol{X}_{\text{up}}[mx + s, my + t] = \boldsymbol{X}[x, y], \quad s, t \in \{0, 1, \ldots, m - 1\}.
$$

Taking the Fourier transform, we have

$$
\begin{aligned}
\mathcal{F}(\boldsymbol{X}_{\text{up}})[u, v] &= \sum_{x=-mH}^{mH-1} \sum_{y=-mW}^{mW-1} \boldsymbol{X}_{\text{up}}[x, y] \cdot \exp\left(-j2\pi\left(\frac{ux}{2mH} + \frac{kvy}{2mW}\right)\right) \\
&= \sum_{x=-H}^{H-1} \sum_{y=-W}^{W-1} \sum_{s=0}^{m-1} \sum_{t=0}^{m-1} \boldsymbol{X}[x, y] \cdot \exp\left(-j2\pi\left(\frac{u(mx + s)}{2mH} + \frac{v(my + t)}{2mW}\right)\right).
\end{aligned}
\tag{*4}
$$

Similar to the proof of Proposition 1.2, by plugging (*2) into (*4), it is easy to see that $\boldsymbol{F}_{\text{up}}(u, v) = \mathcal{F}(\boldsymbol{X}_{\text{up}})[u, v]$ is a linear combination of the signals in original image. Namely, we have

$$
\boldsymbol{F}_{\text{up}}(u, v) = \sum_{u'=-H}^{H-1} \sum_{v'=-W}^{W-1} \alpha_{\text{up}}(u, v, u', v') \cdot \boldsymbol{F}(u', v').
$$

Given any $(u, v, u', v')$, $\alpha_{\text{up}}(u, v, u', v')$ can be computed as

$$
\begin{aligned}
\alpha_{\text{up}}(u, v, u', v') &= \frac{1}{4HW} \sum_{x=-H}^{H-1} \sum_{y=-W}^{W-1} \sum_{s=0}^{m-1} \sum_{t=0}^{m-1} \exp\left(-j2\pi\left(\frac{u(mx + s)}{2mH} + \frac{v(my + t)}{2mW}\right)\right) \cdot \exp\left(j2\pi\left(\frac{u'x}{2H} + \frac{v'y}{2W}\right)\right) \\
&= \frac{1}{4HW} \sum_{x=-H}^{H-1} \sum_{y=-W}^{W-1} \exp\left(j2\pi\left(\frac{(u'-u)x}{2H} + \frac{(v'-v)y}{2W}\right)\right) \sum_{s=0}^{m-1} \sum_{t=0}^{m-1} \exp\left(-j2\pi\left(\frac{us}{2mH} + \frac{vt}{2mW}\right)\right).
\end{aligned}
$$

19

Denote $\beta_{\text{up}}(u,v) = \sum_{s=0}^{m-1} \sum_{t=0}^{m-1} \exp\left(-j2\pi\left(\frac{us}{2mH} + \frac{vt}{2mW}\right)\right)$, which is a constant conditioned on $(u,v)$. Then we know

$$
\begin{aligned}
\alpha_{\text{up}}(u,v,u',v') &= \frac{\beta_{\text{up}}(u,v)}{4HW} \sum_{x=-H}^{H-1} \sum_{y=-W}^{W-1} \exp\left(j2\pi\left(\frac{(u'-u)x}{2H} + \frac{(v'-v)y}{2W}\right)\right) \\
&= \frac{\beta_{\text{up}}(u,v)}{4HW} \sum_{x=-H}^{H-1} \exp\left(j2\pi\left(\frac{(u'-u)x}{2H}\right)\right) \sum_{y=-W}^{W-1} \exp\left(j2\pi\left(\frac{(v'-v)y}{2W}\right)\right) \\
&= \begin{cases} \beta_{\text{up}}(u,v), & u'-u = a \cdot 2H, v'-v = b \cdot 2W, \quad a,b \in \mathbb{Z} \\ 0, & \text{otherwise} \end{cases} .
\end{aligned}
\tag{*5}
$$

Since $-H \leq u \leq H-1, -W \leq v \leq W-1$, we have $\beta_{\text{up}}(u,v) \neq 0$. Thus, we have $\alpha(u,v,u',v') \neq 0$ when $\frac{2H}{k}|(u'-u), \frac{2W}{k}|(v'-v)$.

Now we return to Proposition 1.3. Suppose that the original image is $\boldsymbol{X}$, and that the image generated from down-sampling operation is $\boldsymbol{X}_{\text{d}} = \mathcal{D}_{1/k}(\boldsymbol{X})$, where $k \in \mathbb{Q}^+$ may not be an integer. We can always find two integers $m_0$ and $k_0$ such that $\frac{k_0}{m_0} = k$. Consider first up-sampling $\boldsymbol{X}$ by $m_0$ with nearest interpolation and then performing the down-sampling operation by $k_0$. By combining (*3) and (*5), it is easy to verify that Proposition 1.3 is true. $\qquad\square$