# Deep Frequency Filtering for Domain Generalization

Shiqi Lin[1*]   Zhizheng Zhang[2]   Zhipeng Huang[1*]   Yan Lu[2]   Cuiling Lan[2]   Peng Chu[2]

Quanzeng You[2]   Jiang Wang[2]   Zicheng Liu[2]   Amey Parulkar[2]   Viraj Navkal[2]   Zhibo Chen[1]

[1]University of Science and Technology of China    [2]Microsoft

{linsq047,hzp1104}@mail.ustc.edu.cn  chenzhibo@ustc.edu.cn

{zhizzhang,yanlu,culan,pengchu,quyou,jiangwang,zliu,amey.parulkar,vnavkal}@microsoft.com

## Abstract

*Improving the generalization ability of Deep Neural Networks (DNNs) is critical for their practical uses, which has been a longstanding challenge. Some theoretical studies have uncovered that DNNs have preferences for some frequency components in the learning process and indicated that this may affect the robustness of learned features. In this paper, we propose Deep Frequency Filtering (DFF) for learning domain-generalizable features, which is the first endeavour to explicitly modulate the frequency components of different transfer difficulties across domains in the latent space during training. To achieve this, we perform Fast Fourier Transform (FFT) for the feature maps at different layers, then adopt a light-weight module to learn attention masks from the frequency representations after FFT to enhance transferable components while suppressing the components not conducive to generalization. Further, we empirically compare the effectiveness of adopting different types of attention designs for implementing DFF. Extensive experiments demonstrate the effectiveness of our proposed DFF and show that applying our DFF on a plain baseline outperforms the state-of-the-art methods on different domain generalization tasks, including close-set classification and open-set retrieval.*

## 1. Introduction

Domain Generalization (DG) seeks to break through the *i.i.d.* assumption that training and testing data are identically and independently distributed. This assumption does not always hold in reality since domain gaps are commonly seen between the training and testing data. However, collecting enough training data from all possible domains is costly and even impossible in some practical environments. Thus, learning generalizable feature representations is of

high practical value for both industry and academia.

Recently, a series of research works [97] analyze deep learning from the frequency perspective. These works, represented by the F-Principle [93], uncover that there are different preference degrees of DNNs for the information of different frequencies in their learning processes. Specifically, DNNs optimized with stochastic gradient-based methods tend to capture low-frequency components of the training data with a higher priority [92] while exploiting high-frequency components to trade the robustness (on unseen domains) for the accuracy (on seen domains) [83]. This observation indicates that different frequency components are of different transferability across domains.

In this work, we seek to learn generalizable features from a frequency perspective. To achieve this, we conceptualize Deep Frequency Filtering (DFF), which is a new technique capable of enhancing the transferable frequency components and suppressing the ones not conducive to generalization in the latent space. With DFF, the frequency components of different cross-domain transferability are dynamically modulated in an end-to-end manner during training. This is conceptually simple, easy to implement, yet remarkably effective. In particular, for a given intermediate feature, we apply Fast Fourier Transform (FFT) along its spatial dimensions to obtain the corresponding frequency representations where different spatial locations correspond to different frequency components. In such a frequency domain, we are allowed to learn a spatial attention map and multiply it with the frequency representations to filter out the components adverse to the generalization across domains.

The attention map above is learned in an end-to-end manner using a lightweight module, which is instance-adaptive. As indicated in [83, 92], low-frequency components are relatively easier to be generalized than high-frequency ones while high-frequency components are commonly exploited to trade robustness for accuracy. Although this phenomenon can be observed consistently over different instances, it does not mean that high-frequency com-

---

ponents have the same proportion in different samples or have the same degree of effects on the generalization ability. Thus, we experimentally compare the effectiveness of task-wise filtering with that of instance-adaptive filtering. Here, the task-wise filtering uses a shared mask over all instances while the instance-adaptive filtering uses unshared masks. We find the former one also works but is inferior to our proposed design by a clear margin. As analyzed in [11], the spectral transform theory [38] shows that updating a single value in the frequency domain globally affects all original data before FFT, rendering frequency representation as a global feature complementary to the local features learned through regular convolutions. Thus, a two-branch architecture named Fast Fourier Convolution (FFC) is introduced in [38] to exploit the complementarity of features in the frequency and original domains with an efficient ensemble. To evaluate the effectiveness of our proposed DFF, we choose this two-branch architecture as a base architecture and apply our proposed frequency filtering mechanism to its spectral transform branch. Note that FFC provides an effective implementation for frequency-space convolution while we introduce a novel frequency-space attention mechanism. We evaluate and demonstrate our effectiveness on top of it.

Our contributions can be summarized in the following:

- We discover that the cross-domain generalization ability of DNNs can be significantly enhanced by a simple learnable filtering operation in the frequency domain.
- We propose an effective Deep Frequency Filtering (DFF) module where we learn an instance-adaptive spatial mask to dynamically modulate different frequency components during training for learning generalizable features.
- We conduct an empirical study for the comparison of different design choices on implementing DFF, and find that the instance-level adaptability is required when learning frequency-space filtering for domain generalization.

## 2. Related Work

### 2.1. Domain Generalization

Domain Generalization (DG) aims to improve the generalization ability of DNNs from source domains to unseen domains, which is widely needed in different application scenarios. The challenges of DG have been addressed from *data*, *model*, and *optimization* perspectives. From the *data* perspective, augmentation [31, 78, 98] and generation [73, 81, 113] technologies are devised to increase the diversity of training samples so as to facilitate generalization. From the *model* perspective, some efforts are made to enhance the generalization ability by carefully devising the normalization operations in DNNs [53, 64, 72] or adopting an ensemble of multiple expert models [57, 114]. From the *optimization* perspective, there are many works designing different training strategies to learn generalizable fea-

tures. which is a dominant line in this field. To name a few, some works learn domain-invariant feature representations through explicit feature alignment [23, 37, 62], adversarial learning [22, 24, 49], gradient-based methods [3, 40, 50], causality-based methods [48] or meta-learning based method [86], *etc*. In this work, we showcase a conceptually simple operation, *i.e.*, learnable filtration in the frequency domain, can significantly strengthen the generalization performance on unseen domains, verified on both the close-set classification and open-set retrieval tasks.

### 2.2. Frequency Domain Learning

Frequency analysis has been widely used in conventional digital image processing for decades [4, 66]. Recently, frequency-based operations, *e.g.*, Fourier transform, set forth to be incorporated into deep learning methods for different purposes in four aspects: 1) accelerating the training or facilitating the optimization of DNNs [13, 45, 59, 63, 67, 68]; 2) achieving effective data augmentation [34, 54, 91, 95]; 3) learning informative representations of non-local receptive fields [11,58,70,77,96]; 4) helping analyze and understand some behaviors of DNNs [83,92,93,97] as a tool. As introduced before, prior theoretical studies from the frequency perspective uncover that different frequency components are endowed with different priorities during training and contribute differently to the feature robustness. This inspires us to enhance the generalization ability of DNNs through modulating different frequency components.

In [11], for intermediate features, a $1 \times 1$ convolution in the frequency domain after FFT to learn global representations. However, such global representations capture global characteristics while losing local ones, thus have been demonstrated complementary with the features learned in the original latent space. To address this, a two-branch architecture is proposed in [11] to fuse these two kinds of features. This problem also exists in our work but is not our focus. Thereby, we adopt our proposed frequency filtering operation in the spectral transform branch of the two-branch architecture proposed in [11] for effectiveness evaluation. Besides, in [70], a learnable filter layer is adopted to self-attention (*i.e.*, transformer) to mix tokens representing different spatial locations, which may seem similar with ours at its first glance but is actually not. The learnable filter in [70] is implemented by network parameters while that of ours is the network output thus instance-adaptive. We theoretically analyze and experimentally compare them in the following sections. Besides, with a different purpose from token mixing, we are devoted to improve the generalization ability of DNNs.

## 2.3. Attention Mechanisms

Attention has achieved great success in many visual tasks. It can be roughly categorized into selective attention [5, 32, 36, 69, 84, 88, 102, 104, 115] and self-attention [1, 6, 7, 18, 21, 36, 65, 85] upon their working mechanisms. The former one explicitly learns a mask to enhance task-beneficial features and suppress task-unrelated features. In contrast, self-attention methods commonly take affinities of tokens as the attentions weights to refine the token representations via message passing, wherein the attention weights can be understood to model the importance of other tokens for the query token. Our proposed frequency filtering is implemented with a simple selective attention applied in the frequency domain for the intermediate features of DNNs. There have been a few primary attempts [69, 115] exploiting frequency representations to learn more effective attention. In these works [69, 115], channel attention weights are proposed to learned from multiple frequency components of 2D DCT, where they are still used to modulate channels in the original feature space. In our work, we further investigate the frequency filtering where the the learning and using of attention weights are both in the frequency domain. We make the first endeavour to showcase the effectiveness of such a conceptually simple mechanism for the DG field, and would leave more delicate designs of attention model architectures for the frequency domain in our future work.

## 3. Deep Frequency Filtering

### 3.1. Problem Definition and Core Idea

In this paper, we aim to reveal that the generalization ability of DNNs to unseen domains can be significantly enhanced through an extremely simple mechanism, *i.e.*, an explicit frequency modulation in the latent space, named Deep Frequency Filtering (DFF). To shed light on this core idea, we first introduce the problem definition of Domain Generalization (DG) as preliminaries. Given $K$ source domains $\mathcal{D}_s = \{D_s^1, D_s^2, \cdots, D_s^K\}$, where $D_s^k = (\mathbf{x}_i^k, y_i^k)_{i=1}^{N_k}$ denotes the $k$-th domain consisting of $N_k$ samples $\mathbf{x}_i^k$ with their corresponding labels $y_i^k$, the goal of DG is to enable the model trained on source domains $\mathcal{D}_s$ perform as well as possible on unseen target domains $\mathcal{D}_t$, without additional model updating using the data in $\mathcal{D}_t$. When different domains share the same label space, it corresponds to a closed-set DG problem, otherwise an open-set problem.

As introduced before, the studies for the behaviors of DNNs from the frequency perspective [83, 92, 93] have uncovered that the DNNs have different preferences for different frequency components of the learned intermediate features. The frequency characteristics affect the trade-off between robustness and accuracy [83]. This inspires us to improve the generalization ability of DNNs through modulating different frequency components of different trans-

fer difficulties across domains during training. Achieved by a simple filtering operation, transferable frequency components are enhanced while the components prejudice to cross-domain generalization are suppressed.

### 3.2. Latent Frequency Representations

Different from previous frequency-based methods [54, 91, 95] applied in the pixel space (*i.e.*, the side of inputs), we adopt our proposed filtering operation in the latent space. In this section, we briefly recall a conventional signal processing tool Fast Fourier Transform (FFT). We adopt it for obtaining the feature representations in the frequency domain, then discuss the characteristics of such representations.

Given the intermediate features $\mathbf{X} \in \mathbb{R}^{C \times H \times W}$, we perform a 2D fast Fourier transform (*i.e.*, an accelerated version [15] of 2D discrete Fourier transform) for each channel independently to get the corresponding frequency representations $\mathbf{X}_F \in \mathbb{R}^{2C \times H \times (\lfloor \frac{W}{2} \rfloor + 1)}$. We formulate this transform $\mathbf{X}_F = FFT(\mathbf{X})$ as below (where the channel dimension is omitted for brevity):

$$\mathbf{X}_F(x, y) = \sum_{h=0}^{H-1} \sum_{w=0}^{W-1} \mathbf{X}(h, w) e^{-j2\pi \left(x \frac{h}{H} + y \frac{w}{W}\right)}. \quad (1)$$

The frequency representation $\mathbf{X}_F$ can be transferred to the original feature space via an inverse FFT, succinctly expressed as $\mathbf{X} = iFFT(\mathbf{X}_F)$, which can be formulated as:

$$\mathbf{X}(h, w) = \frac{1}{H \cdot W} \sum_{h=0}^{H-1} \sum_{w=0}^{W-1} \mathbf{X}_F(x, y) e^{j2\pi \left(x \frac{h}{H} + y \frac{w}{W}\right)}. \quad (2)$$

The $\mathbf{X}_F \in \mathbb{R}^{2C \times H \times (\lfloor \frac{W}{2} \rfloor + 1)}$ above denotes the frequency representation of $\mathbf{X} \in \mathbb{R}^{C \times H \times W}$, which concatenates the real and imaginary parts after FFT (each one has $C$ channels). Besides, thanks to the conjugate symmetric property of FFT, $\mathbf{X}_F$ only needs retain the half of spatial dimensions thus has a spatial resolution of $H \times (\lfloor \frac{W}{2} \rfloor + 1)$. For the frequency representation $\mathbf{X}_F$, there are two utilizable properties: 1) Different frequency components of the original feature $\mathbf{X}$ are decomposed into elements at different spatial locations of $\mathbf{X}_F$, which could be viewed as a frequency-based disentanglement and re-arrangement for $\mathbf{X}$. This property makes the learning in frequency domain efficient in practice, and more importantly, allows us to achieve frequency filtering with a simple devised spatial attention module. 2) $\mathbf{X}_F$ is a naturally global feature representation, as discussed in [11], which can facilitate the suppression of globally distributed domain-specific information, such as illumination, imaging noises, *etc*. Next, we shed light on the specific filtering operation on $\mathbf{X}_F$.

### 3.3. Latent-space Frequency Filtering

Our goal is to adaptively modulate different frequency components over different network depths during training.
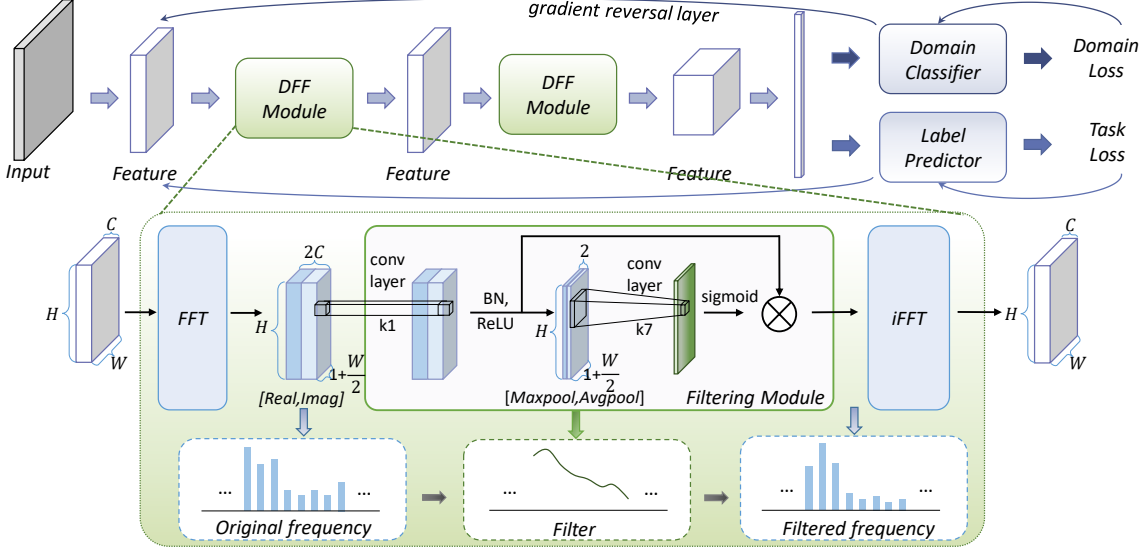
Figure 1. Illustration of our proposed Deep Frequency Filtering (DFF) module. DFF learns an instance-adaptive spatial mask to dynamically modulate different frequency components during training for learning generalizable features, which consists of three key operations: a 2D fast Fourier transform (FFT) to convert the input spatial features to the frequency domain, an filtering module to enhance the transferable components while suppressing the generalization-detrimental ones, and a 2D inverse FFT (iFFT) to map the features back to the orginal latent space.

We thus propose to apply a frequency filtering operation on $\mathbf{X}_F$ to enhance the transferable components while suppressing the generalization-detrimental ones. Thanks to the first hallmark of $\mathbf{X}_F$ discussed in Sec. 3.2, the frequency filtering operation is allowed to be implemented with a spatial attention on $\mathbf{X}_F$. Given a frequency representation $\mathbf{X}_F \in \mathbb{R}^{2C \times H \times (\lfloor \frac{W}{2} \rfloor + 1)}$, the proposed frequency filtering mechanism is formulated as follows:

$$\mathbf{X}'_F = \mathbf{X}_F \otimes \mathbf{M}_S(\mathbf{X}_F), \tag{3}$$

where $\otimes$ denotes element-wise multiplication. $\mathbf{M}_S(\cdot)$ refers to the attention module to learn a spatial mask with a resolution of $H \times (\lfloor \frac{W}{2} \rfloor + 1)$. This mask is copied along the channel dimension of $\mathbf{X}_F$ accordingly for the element-wise multiplication, filtering out the components adverse to the generalization in $\mathbf{X}_F$. The frequency feature after filtering is denoted by $\mathbf{X}'_F$. Our contributions lie in revealing such a frequency filtering operation in the latent space can bring impressive improvements for DG, although using a lightweight attention architecture designed for the features in original latent space [88] to implement $\mathbf{M}_S(\cdot)$. This provides another alternative to the field of DG, which is conceptually simple, significantly effective, but previously unexplored. Besides, we further conduct an empirical study to investigate the effectiveness of different attention types in implementing our conceptualized deep frequency filtering in the experiment section. The specific architecture design for the attention module is not our current focus, but is worth being explored in the future.

Here, we introduce an extremely simple instantiation for Eq. (3). We use this for verifying our proposed concept *deep frequency filtering* in this paper. For $\mathbf{X}_F \in \mathbb{R}^{2C \times H \times (\lfloor \frac{W}{2} \rfloor + 1)}$ that consists of real and imaginary parts after FFT, inspired by the attention architecture design in [88], we first adopt a $1 \times 1$ convolutional layer followed by Batch Normalization (BN) and ReLU activation to project $\mathbf{X}_F$ to an embedding space for the subsequent filtration. After embedding, as shown in Fig. 1, we follow the spatial attention architecture design in [88] to aggregate the information of $\mathbf{X}_F$ over channels using both average-pooling and max-pooling operations along the channel axis, generating two frequency descriptors denoted by $\mathbf{X}_F^{avg}$ and $\mathbf{X}_F^{max}$, respectively. These two descriptors can be viewed as two compact representations of $\mathbf{X}_F$ in which the information of each frequency component is compressed separately by the pooling operations while the spatial discriminability is still preserved. We then concatenate $\mathbf{X}_F^{avg}$ with $\mathbf{X}_F^{max}$ and use a large-kernel $7 \times 7$ convolution layer followed by a sigmoid function to learn the spatial mask. Mathematically, this instantiation can be formulated as:

$$\mathbf{X}'_F = \mathbf{X}_F \otimes \sigma(\mathrm{Conv}_{7\times7}([\mathrm{AvgPool}(\mathbf{X}_F), \mathrm{MaxPool}(\mathbf{X}_F)])), \tag{4}$$

where $\sigma$ denotes the sigmoid function. The $[\cdot, \cdot]$ is a concatenation operation. $\mathrm{AvgPool}(\cdot)$ and $\mathrm{MaxPool}(\cdot)$ denote the average and max pooling operations, respectively. $\mathrm{Conv}_{7\times7}(\cdot)$ is a convolution layer with the kernel size of 7. Albeit using a large-size kernel, the fea-

ture $[\mathrm{AvgPool}(\mathbf{X}_F), \mathrm{MaxPool}(\mathbf{X}_F)]$ has only two channels through the information squeeze by pooling operations such that this step is still very computationally efficient in practice. We omit the embedding of $\mathbf{X}_F$ in this formulation for brevity. We believe using more complex attention architectures, such as [17, 61, 104], is of the potentials to achieve higher improvements, and we expect more effective instantiations of our conceptualized Deep Frequency Filtering.

**Discussion.** The proposed Deep Frequency Filtering is conceptually new design to achieve instance-adaptive frequency modulation in the latent space of DNNs. It also corresponds to a novel neural operation albeit using an off-the-shelf architecture design as an exampled instantiation. Compared to prior frequency-domain works [11, 70], we make the first endeavour to introduce an explicit instance-adaptive frequency selection mechanism into the optimization of DNNs. From the perspective of attention, conventional attention designs [32, 85, 88, 104] learn masks from deep features in the original latent space, and adopt the learned masks to these features themselves to achieve feature modulation. FcaNet [69] strives to a further step by learning channel attention weights from the results of frequency transform. But the learned attention weights are still used for the original features. In this aspect, we are the first to learn attention weights from frequency representations and also use the learned masks in the frequency domain to achieve our conceptualized frequency filtering.

### 3.4. Post-filtering Feature Restitution

The features captured in the frequency domain have been demonstrated to be global and complementary to the local ones captured in the original latent space in [11]. Thus, a simple two-branch is designed to exploit this complementarity to achieve an ensemble of both local and global features in [11]. This architecture is naturally applicable to the restitution of complementary local features as a post-filtering refinement in the context of our proposed concept. We thus evaluate the effectiveness of our proposed method on top of the two-branch architecture in [11]. Specifically, similar to [11], we split the given intermediate feature $\mathbf{X} \in \mathbb{R}^{C \times H \times W}$ along its channel dimension into $\mathbf{X}^g \in \mathbb{R}^{rC \times H \times W}$ and $\mathbf{X}^l \in \mathbb{R}^{(1-r)C \times H \times W}$. The two-branch architecture can be formulated as:

$$Y^l = f_l(\mathbf{X}^l) + f_{g \to l}(\mathbf{X}^g), \quad Y^g = f_g(\mathbf{X}^g) + f_{l \to g}(\mathbf{X}^l),$$
(5)

where $f_l(\cdot)$, $f_g(\cdot)$, $f_{l \to g}(\cdot)$ and $f_{g \to l}(\cdot)$ denote four different transformation functions. Among them, $f_l(\cdot)$, $f_{l \to g}(\cdot)$ and $f_{g \to l}(\cdot)$ are three regular convolution layers. In [11], the $f_g(\cdot)$ corresponds to the spectral transform implemented by their proposed convolution operation in the frequency domain. On top of it, we evaluate the effectiveness of our pro-

posed DFF by adding this operation into the spectral transform branch of this architecture to achieve an explicit filtering operation in the frequency domain. The contribution on this two-branch architecture design belongs to [11].

### 3.5. Model Training

In addition to commonly used task-related loss functions (for classification or retrieval), we train a domain classifier with a domain classification loss and adopt a gradient reversal layer [22]. These are commonly used in DG research works for explicitly encouraging the learning of domain-invariant features and the suppression for features conducive to domain generalization. The feature extractor is optimized for minimizing the task losses while maximizing the domain classification loss simultaneously.

## 4. Experiments

### 4.1. Datasets and Settings

We evaluate the effectiveness of our proposed Deep Frequency Filtering (DFF) for Domain Generalization (DG) on **Task-1**: the close-set classification task and **Task-2**: the open-set retrieval task, *i.e.*, person re-identification (ReID).

For Task-1, Office-Home dataset [41] is a commonly used domain generalization (DG) benchmark on the task of classification. It consists of 4 domains (*i.e.*, Art (Ar), Clip Art (Cl), Product (Pr), Real-World (Rw)). Among them, three domains are used for training and the remaining one is considered as the unknown target domain for testing.

For Task-2, person ReID is a representative open-set retrieval task, where different domains do not share their label space. i) following [53, 106], we take four large-scale datasets (CUHK-SYSU (CS) [89], MSMT17 (MS) [87], CUHK03 (C3) [47] and Market-1501 (MA) [108]). For evaluation, a model is trained on three domains and tested on the remaining one. ii) several large-scale ReID datasets *e.g.*, CUHK02 [46], CUHK03, Market-1501 and CUHK-SYSU, are viewed as multiple source domains. Each small-scale ReID dataset including VIPeR [25], PRID [30], GRID [55] and iLIDS [109] is used as an unseen target domain. To comply with the General Ethical Conduct, we exclude DukeMTMC from the source domains.

We adopt ResNet-18 and ResNet-50 [26] as our backbone for Task-1 and Task-2, respectively. All reported results are obtained by the averages of five runs. We provide more implementation details in the supplementary material.

### 4.2. Ablation Study

#### 4.2.1 The effectiveness of DFF

To investigate the effectiveness of our proposed Deep Frequency Filtering (DFF), we compare it with the ResNet baselines (*Base*) and the ResNet-based FFC [11] models

Table 1. Performance comparisons of our proposed Deep Frequency Filtering (DFF) with the baselines and the models with deep filtering in the original feature space. "*Base*" refers to the vanilla ResNet baseline. In "*SBase*", we use ResNet-based FFC in [11], serving as a strong baseline. "*Ori-F*" refers to a filtering operation in the original feature space, adopted in the local branch $f_l$ and global branch $f_g$, respectively. When adopted in $f_g$, the FFT/iFFT operations are discarded from $f_g$ so that the filtering is in the original feature space instead of the frequency space. "*Fre-F*" represents our proposed Deep Frequency Filtering.

| Method | Source→Target | | | | | |
| | MS+CS+C3→MA | | MS+MA+CS→C3 | | MA+CS+C3→MS | |
| | mAP | R1 | mAP | R1 | mAP | R1 |
|---|---|---|---|---|---|---|
| Base | 59.4 | 83.1 | 30.3 | 29.1 | 18.0 | 41.9 |
| SBase (FFC) | 66.2 | 84.7 | 35.8 | 35.4 | 19.4 | 44.8 |
| Ori-F in $f_l$ | 66.9 | 85.0 | 36.2 | 35.9 | 19.8 | 45.1 |
| Ori-F in $f_g$ | 61.9 | 83.5 | 32.7 | 31.9 | 18.4 | 42.8 |
| Fre-F (Ours) | **71.1** | **87.1** | **41.3** | **41.1** | **25.1** | **50.5** |

Table 2. Performance comparisons of different implementations of the frequency filtering operation. "*Task.*" refers to the filtering operation using a task-level attention mask where the mask is implemented with network parameters and is shared over different instances. "*Ins.*" denotes the filtering operation using learned instance-adaptive masks. "*C*" and "*S*" represents the filtering performed along the channel and spatial dimensions, respectively.

| Method | Source→Target | | | | | |
| | MS+CS+C3→MA | | MS+MA+CS→C3 | | MA+CS+C3→MS | |
| | mAP | R1 | mAP | R1 | mAP | R1 |
|---|---|---|---|---|---|---|
| Base | 59.4 | 83.1 | 30.3 | 29.1 | 18.0 | 41.9 |
| Task.(C) | 62.7 | 80.0 | 32.1 | 31.4 | 19.5 | 44.9 |
| Task.(S) | 68.6 | 85.8 | 37.0 | 36.3 | 20.8 | 45.4 |
| Ins.(C) | 69.8 | 86.2 | 36.4 | 35.9 | 21.0 | 45.7 |
| Ins.(S) (Ours) | **71.1** | **87.1** | **41.3** | **41.1** | **25.1** | **50.5** |

(*SBase*) that serve as the strong baselines, respectively. The experiment results are in Table 1. There are two crucial observations: 1) Our *Fre-F (DFF)* consistently outperforms *Base* by a clear margin. This demonstrates the effectiveness of our proposed method. With a simple instantiation, our proposed Deep Frequency Filtering operation can significantly improve the generalization ability of DNNs. 2) Our *Fre-F (DFF)* brings more improvements than *SBase (FFC)*. This indicates that frequency-domain convolutions are inadequate for model generalization. Instead, DFF is an effective solution for classification generalization task.

### 4.2.2 Frequency-domain v.s. Original-domain filtering

Feature filtering operations can be implemented in either the original feature domain or the frequency domain. We conduct experiments to study the impact of different implementation in both domains. Table 1 indicates that the proposed frequency-domain filtering (*Fre-F*) outperforms the original-domain feature filtering (*Ori-F*). This demonstrates the importance of modulating different frequency components in the latent space for domain generalization.

Table 3. The ablation results on the influence of the two-branch architecture. "*FFC*" refers to the frequency-domain convolution work [11] without our proposed filtering operation. "*only $f_g$*" denotes the setting in which we adopt the *FFC* or our proposed *DFF* over the complete feature without the splitting along the channel dimension, corresponding to using a single branch architecture. "*$f_l + f_g$*" represents the setting in which we split the feature along the channel dimension and adopt frequency-domain operations (*FFC* or our *DFF*) on one half of them.

| Method | Source→Target | | | | | |
| | MS+CS+C3→MA | | MS+MA+CS→C3 | | MA+CS+C3→MS | |
| | mAP | R1 | mAP | R1 | mAP | R1 |
|---|---|---|---|---|---|---|
| Base | 59.4 | 83.1 | 30.3 | 29.1 | 18.0 | 41.9 |
| FFC(only $f_g$) | 60.1 | 80.4 | 26.1 | 24.8 | 17.3 | 40.2 |
| FFC($f_l + f_g$) | 66.2 | 84.7 | 35.8 | 35.4 | 19.4 | 44.8 |
| Ours(only $f_g$) | 64.2 | 83.4 | 29.3 | 28.1 | 17.6 | 40.4 |
| Ours($f_l + f_g$) | **71.1** | **87.1** | **41.3** | **41.1** | **25.1** | **50.5** |

### 4.2.3 The importance of instance-adaptive attention

Our proposed DFF is allowed to be implemented using a spatial attention on the frequency representations of features. In Table 2, we compare the performance of adopting task-level and instance-level attention mask as well as channel and spatial filtering, respectively. We can observe that *Ins. (S)* clearly and consistently outperforms *Task. (S)* over all settings on both tasks. The results suggest that using instance-adaptive attention mask is essential for implementing our proposed DFF. We observe that different instances correspond to diversified frequency components in the feature space. Thus, a task-level attention mask is weak for all instances, which may explain the performance gaps in Table 2. We need to perform the modulation of feature frequency components with instance-adaptive weights.

### 4.2.4 Spatial v.s. Channel

DFF can be implemented with a spatial attention module, since different spatial positions of the frequency representations correspond to different frequencies. Admittedly, we can also adopt the channel attention to the frequency representations, which can be viewed as a representation refinement within each frequency component rather than perform Deep Frequency Filtering. The results in Table 2 show that spatial attention consistently outperforms channel attention in both task-level and instance-level settings. This indicates that frequency filtering (or selection) is more important than the refinement of frequency-domain representations for domain generalization.

### 4.2.5 The Influence of the Two-branch Architecture

As mentioned above, we adopt the two-branch architecture proposed by FFC [12] as a base architecture and apply our frequency filtering mechanism to the spectral transform branch. As shown in Table 3, the performances of

Table 4. Performance (classification accuracy %) comparison with the state-of-the-art methods on close-set classification task. We use ResNet-18 as backbone. Best in bold.

| Method | Source→Target | | | | Avg |
|---|---|---|---|---|---|
| | Cl,Pr,Rw→Ar | Ar,Pr,Rw→Cl | Ar,Cl,Rw→P | Ar,Cl,Pr→Rw | |
| Baseline | 52.2 | 45.9 | 70.9 | 73.2 | 60.5 |
| CCSA [62] | 59.9 | 49.9 | 74.1 | 75.7 | 64.9 |
| D-SAM [20] | 58.0 | 44.4 | 69.2 | 71.5 | 60.8 |
| MMD-AAE [43] | 56.5 | 47.3 | 72.1 | 74.8 | 62.7 |
| CrossGrad [73] | 58.4 | 49.4 | 73.9 | 75.8 | 64.4 |
| JiGen [8] | 53.0 | 47.5 | 71.5 | 72.8 | 61.2 |
| RSC [35] | 58.4 | 47.9 | 71.6 | 74.5 | 63.1 |
| MixStyle [114] | 58.7 | 53.4 | 74.2 | 75.9 | 65.5 |
| Ours | **65.4** | **53.7** | **74.9** | **76.5** | **67.6** |

Table 5. Performance (%) comparison with the state-of-the-art methods on the open-set person ReID task. "Source" refers to the multiple training datasets, i.e., Market-1501 (MA), DukeMTMC-reID (D), CUHK-SYSU (CS), CUHK03 (C3) and CUHK02 (C2). "All" represents using MA+D+CS+C3+C2 as source domains. We do not include DukeMTMC-reID (D) in the training domains since this dataset has been discredited by the creators, denoted as "All w/o D".

| Method | Source | Target:VIPeR(V) | | Target:PRID(P) | | Target:GRID(G) | | Target:iLIDS(I) | | Mean of V,P,G,I | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | R1 | mAP | R1 | mAP | R1 | mAP | R1 | mAP | R1 | mAP |
| CDEL [52] | All | 38.5 | - | 57.6 | - | 33.0 | - | 62.3 | - | 47.9 | - |
| DIMN [74] | All | 51.2 | 60.1 | 39.2 | 52.0 | 29.3 | 41.1 | 70.2 | 78.4 | 47.5 | 57.9 |
| DDAN [9] | All | 56.5 | 60.8 | 62.9 | 67.5 | 46.2 | 50.9 | 78.0 | 81.2 | 60.9 | 65.1 |
| RaMoE [16] | All | 56.6 | 64.6 | 57.7 | 67.3 | 46.8 | 54.2 | 85.0 | 90.2 | 61.5 | 69.1 |
| SNR [37] | All | 49.2 | 58.0 | 47.3 | 60.4 | 39.4 | 49.0 | 77.3 | 84.0 | 53.3 | 62.9 |
| CBN [116] | All | 49.0 | 59.2 | 61.3 | 65.7 | 43.3 | 47.8 | 75.3 | 79.4 | 57.2 | 63.0 |
| Person30K [2] | All | 53.9 | 60.4 | 60.6 | 68.4 | 50.9 | 56.6 | 79.3 | 83.9 | 61.1 | 67.3 |
| DIR-ReID [100] | All | 58.3 | 62.9 | 71.1 | 75.6 | 47.8 | 52.1 | 74.4 | 78.6 | 62.9 | 67.3 |
| MetaBIN [14] | All | 56.2 | 66.0 | 72.5 | 79.8 | 49.7 | 58.1 | 79.7 | 85.5 | 64.5 | 72.4 |
| QAConv$_{50}$ [51] | All w/o D | 57.0 | 66.3 | 52.3 | 62.2 | 48.6 | 57.4 | 75.0 | 81.9 | 58.2 | 67.0 |
| M$^3$L [107] | All w/o D | 60.8 | 68.2 | 55.0 | 65.3 | 40.0 | 50.5 | 65.0 | 74.3 | 55.2 | 64.6 |
| MetaBIN [14] | All w/o D | 55.9 | 64.3 | 61.2 | 70.8 | 50.2 | 57.9 | 74.7 | 82.7 | 60.5 | 68.9 |
| Ours | All w/o D | **65.7** | **74.2** | **71.8** | **78.6** | **56.4** | **65.5** | **83.6** | **88.3** | **69.4** | **76.7** |

*FFC(only $f_g$)* and *FFC($f_l + f_g$)* are inferior to *Ours(only $f_g$)* and *Ours($f_l + f_g$)* respectively, which indicates that simply filtering frequency components of features can bring striking improvement of the generalization ability. And the performance gap between *Ours($f_l + f_g$)* and Ours(only $f_g$) demonstrates that the two-branch structure for Post-filtering Feature Restitution can restore the globally filtered features. Furthermore, the complementary local filtering helps the learning of generalizable feature representation.

### 4.3. Comparison with the State-of-the-arts

#### 4.3.1 Performances on close-set classification (Task-1)

In Table 4, we show the comparisons with the state-of-the-art approaches for Task-1 on Office-Home dataset. All our reported results are averaged over five runs. We observe that our proposed DFF outperforms most existing DG methods and achieves **67.6%** classification accuracy on average using ResNet-18, and outperforms the second-best method MixStyle [114] by **2.2%** average classification accuracy.

#### 4.3.2 Performances on open-set retrieval (Task-2)

As shown in Table 5 and Table 6, our DFF model achieves significant performance improvements on all settings. Specifically, (see the Table 5), the mean performance of our method exceeds the second-best by **8.9%** in R1 ac-

curacy and **7.8%** in mAP. As shown in Table 6, our DFF outperforms the second-best by **3.4%, 6.7%, 5.2%** in R1 accuracy and **8.0%, 7.1%, 7.3%** in mAP on Market-1501, CUHK03, MSMT17, respectively. When under the setting (*i.e.*, the testing set of the seen domains are also included for training model), our DFF performs better than the second-best by **7.3%, 8.1%, 5.2%** in R1 accuracy and **13.8%, 8.5%, 6.5%** in mAP. The results demonstrate that our DFF can significantly improve the generalization ability of learned features even with a simple learned filtering operation in the frequency domain.

### 4.4. Complexity Analysis

In Table 13, we compare the complexities of our DFF and vanilla ResNet models. The GFLOPs are calculated with input size of $224 \times 224$. FFT and inverse FFT are parameter-free and our filtering design uses average-pooling and max-pooling operations to reduce computational cost. Thus, our DFF module only brings limited extra GFLOPs and parameters. Our experiment results demonstrate significant performance gain over vanilla ResNet variants.

### 4.5. Visualization of Learned Masks

We visualize the learned masks at different depths used for Deep Frequency Filtering in our proposed scheme. The visualization results are in Fig. 2. We draw two observa-

Table 6. Performance (%) comparison with the state-of-the-art methods on the open-set person ReID task. Evalustion on four large-scale person ReID benchmarks including Market-1501 (MA), Cuhk-SYSC (CS), CUHK03 (C3) and MSMT17 (MS). 'Com-' refers to combining the train and test sets of source domains for training. The $M^3L$ with a superscript * denote the model adopting IBN-Net50 as backbone. Without this superscript, ResNet-50 is taken as the backbone.

| Method | Source | Market-1501 | | Source | CUHK03 | | Source | MSMT17 | |
|--------|--------|-------------|------|--------|--------|------|--------|--------|------|
| | | mAP | R1 | | mAP | R1 | | mAP | R1 |
| SNR [37] | | 34.6 | 62.7 | | 8.9 | 8.9 | | 6.8 | 19.9 |
| $M^3L$ [107] | | 58.4 | 79.9 | | 20.9 | 31.9 | | 15.9 | 36.9 |
| $M^3L$* [107] | MS+CS+C3 | 61.5 | 82.3 | MS+CS+MA | 34.2 | 34.4 | CS+MA+C3 | 16.7 | 37.5 |
| $QAConv_{50}$ [51] | | 63.1 | 83.7 | | 25.4 | 24.8 | | 16.4 | 45.3 |
| MetaBIN [14] | | 57.9 | 80.0 | | 28.8 | 28.1 | | 17.8 | 40.2 |
| Ours | | **71.1** | **87.1** | | **41.3** | **41.1** | | **25.1** | **50.5** |
| SNR [37] | | 52.4 | 77.8 | | 17.5 | 17.1 | | 7.7 | 22.0 |
| $M^3L$ [107] | | 61.2 | 81.2 | | 32.3 | 33.8 | | 16.2 | 36.9 |
| $M^3L$* [107] | Com-(MS+CS+C3) | 62.4 | 82.7 | Com-(MS+CS+MA) | 35.7 | 36.5 | Com-(CS+MA+C3) | 17.4 | 38.6 |
| $QAConv_{50}$ [51] | | 66.5 | 85.0 | | 32.9 | 33.3 | | 17.6 | 46.6 |
| MetaBIN [14] | | 67.2 | 84.5 | | 43.0 | 43.1 | | 18.8 | 41.2 |
| Ours | | **81.0** | **92.3** | | **51.5** | **51.2** | | **25.3** | **51.8** |

Table 7. Comparison of Parameters (#Para), GFLOPs and the Top-1 classification accuracy (Acc.) on ImageNet-1K between the models equipped with DFF and their corresponding base models. "*DFF-R18/R50*" denote the ResNet-18/-50 models equipped with our DFF.

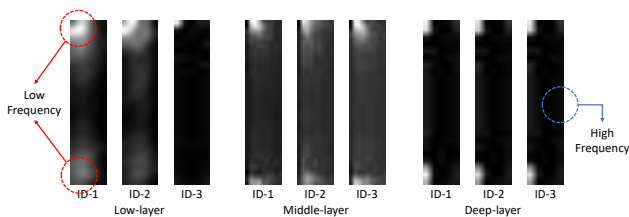| Model | #Para | GFLOPs | Acc. | Model | #Para | GFLOPs | Acc. |
|-------|-------|--------|------|-------|-------|--------|------|
| ResNet18 | 11.7M | 1.8 | 69.8 | ResNet50 | 25.6M | 4.1 | 76.3 |
| DFF-R18 | 12.2M | 2.0 | 72.3 | DFF-R50 | 27.7M | 4.5 | 77.9 |



Figure 2. Visualization of the learned spatial masks of DFF for the filtering in the frequency domain. The grayscale denotes the mask value, while "white" and "black" correspond to "1" and "0", respectively. For each mask, the left top and bottom corners correspond to the low-frequency components while the right middle corresponds to the high-frequency components. The masks at different depths are resized to the same resolution.

tions: 1) The models equipped with DFF tend to enhance relatively low frequency components while suppressing relatively high ones in the latent frequency domain, which is consistently observed over different instances (IDs). This is in line with the results of theoretical study on the relationship between frequency and robustness [83, 92, 93]. 2) The learned masks are instance-adaptive, demonstrating our proposed scheme can achieve instance-adaptive frequency modulation/filtering in the latent space.

## 4.6. Visualization of Learned Feature Maps

We compare the feature maps extracted by the model equipped with our proposed DFF and the one without DFF (see Fig. 3). We observe that the features learned by

the model equipped with DFF have higher responses for human-body regions than those learned by the baselines model without DFF. This indicates that DFF enables neural networks to focus more precisely on target regions and suppress unrelated feature components (*e.g.*, backgrounds).



Figure 3. We compare the feature maps extracted by the model without (left) and the model with DFF (right). The *lighter* the color is, the *larger* the feature value is.

Moreover, we also perform t-SNE visualization for the ReID feature vectors learned by the baseline and the model with DFF. The results and analysis are in the supplementary.

## 5. Conclusion

In this paper, we first conceptualize Deep Frequency Filtering (DFF) and point out that such a simple mechanism can significantly enhance the generalization ability of deep neural networks across domains. This provides a novel alternative for this field. Furthermore, we discuss the implementations of DFF and showcase the implementation with a simple spatial attention in the frequency domain can bring stunning performance improvements for DG. Extensive experiments and ablation studies demonstrate the effectiveness of our proposed method. We leave the exploration on more effective instantiations of our conceptualized DFF in the future work, and encourage more combinations and interplay between conventional signal processing and deep learning technologies.

# References

[1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *ICCV*, pages 6836–6846, 2021. 3

[2] Yan Bai, Jile Jiao, Wang Ce, Jun Liu, Yihang Lou, Xuetao Feng, and Ling-Yu Duan. Person30k: A dual-meta generalization network for person re-identification. In *CVPR*, pages 2123–2132, 2021. 7

[3] Yogesh Balaji, Swami Sankaranarayanan, and Rama Chellappa. Metareg: Towards domain generalization using meta-regularization. *NIPS*, 31, 2018. 2, 14

[4] Gregory A Baxes. *Digital image processing: principles and applications*. John Wiley & Sons, Inc., 1994. 2

[5] Irwan Bello, Barret Zoph, Ashish Vaswani, Jonathon Shlens, and Quoc V Le. Attention augmented convolutional networks. In *ICCV*, pages 3286–3295, 2019. 3

[6] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *ICML*, pages 813–824, 2021. 3

[7] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, pages 213–229, 2020. 3

[8] Fabio M Carlucci, Antonio D'Innocente, Silvia Bucci, Barbara Caputo, and Tatiana Tommasi. Domain generalization by solving jigsaw puzzles. In *CVPR*, pages 2229–2238, 2019. 7, 13, 14

[9] Peixian Chen, Pingyang Dai, Jianzhuang Liu, Feng Zheng, Qi Tian, and Rongrong Ji. Dual distribution alignment network for generalizable person re-identification. In *AAAI*, volume 6, 2021. 7, 13

[10] Xiaodong Chen, Xinchen Liu, Wu Liu, Xiao-Ping Zhang, Yongdong Zhang, and Tao Mei. Explainable person re-identification with attribute-guided metric distillation. In *ICCV*, pages 11813–11822, 2021. 15

[11] Lu Chi, Borui Jiang, and Yadong Mu. Fast fourier convolution. In *NeurIPS*, pages 4479–4488, 2020. 2, 3, 5, 6, 13, 15

[12] Lu Chi, Guiyu Tian, Yadong Mu, Lingxi Xie, and Qi Tian. Fast non-local neural networks with spectral residual learning. In *Proceedings of the 27th ACM International Conference on Multimedia*, pages 2142–2151, 2019. 6

[13] Kamran Chitsaz, Mohsen Hajabdollahi, Nader Karimi, Shadrokh Samavi, and Shahram Shirani. Acceleration of convolutional neural network using fft-based split convolutions. *arXiv preprint arXiv:2003.12621*, 2020. 2

[14] Seokeon Choi, Taekyung Kim, Minki Jeong, Hyoungseob Park, and Changick Kim. Meta batch-instance normalization for generalizable person re-identification. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 3425–3435, 2021. 7, 8, 13, 14

[15] James W Cooley and John W Tukey. An algorithm for the machine calculation of complex fourier series. *Mathematics of computation*, 19(90):297–301, 1965. 3

[16] Yongxing Dai, Xiaotong Li, Jun Liu, Zekun Tong, and Ling-Yu Duan. Generalizable person re-identification with relevance-aware mixture of experts. In *CVPR*, pages 16145–16154, 2021. 7, 14

[17] Zihang Dai, Hanxiao Liu, Quoc Le, and Mingxing Tan. Coatnet: Marrying convolution and attention for all data sizes. In *NeurIPS*, volume 34, 2021. 5

[18] 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. 3

[19] Qi Dou, Daniel C Castro, Konstantinos Kamnitsas, and Ben Glocker. Domain generalization via model-agnostic learning of semantic features. In *NeurIPS*, pages 6447–6458, 2019. 14

[20] Antonio D'Innocente and Barbara Caputo. Domain generalization with domain-specific aggregation modules. In *GCPR*, pages 187–198. Springer, 2018. 7

[21] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *ICCV*, pages 6824–6835, 2021. 3

[22] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*, pages 1180–1189. PMLR, 2015. 2, 5, 14

[23] Muhammad Ghifary, W Bastiaan Kleijn, Mengjie Zhang, and David Balduzzi. Domain generalization for object recognition with multi-task autoencoders. In *ICCV*, pages 2551–2559, 2015. 2

[24] Rui Gong, Wen Li, Yuhua Chen, and Luc Van Gool. Dlow: Domain flow for adaptation and generalization. In *CVPR*, pages 2477–2486, 2019. 2

[25] Douglas Gray and Hai Tao. Viewpoint invariant pedestrian recognition with an ensemble of localized features. In *ECCV*, pages 262–275, 2008. 5, 13

[26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5, 13, 14

[27] Lingxiao He, Xingyu Liao, Wu Liu, Xinchen Liu, Peng Cheng, and Tao Mei. Fastreid: A pytorch toolbox for general instance re-identification. *arXiv preprint arXiv:2006.02631*, 2020. 13, 14

[28] Lingxiao He and Wu Liu. Guided saliency feature learning for person re-identification in crowded scenes. In *ECCV*, pages 357–373, 2020. 15, 16

[29] Lingxiao He, Yinggang Wang, Wu Liu, He Zhao, Zhenan Sun, and Jiashi Feng. Foreground-aware pyramid reconstruction for alignment-free occluded person re-identification. In *ICCV*, pages 8450–8459, 2019. 15

[30] Martin Hirzer, Csaba Beleznai, Peter M Roth, and Horst Bischof. Person re-identification by descriptive and discriminative classification. In *SCIA*, pages 91–102, 2011. 5, 13

[31] Narges Honarvar Nazari and Adriana Kovashka. Domain generalization using shape representation. In *ECCV*, pages 666–670, 2020. 2

[32] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, pages 7132–7141, 2018. 3, 5

[33] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017. 14

[34] Jiaxing Huang, Dayan Guan, Aoran Xiao, and Shijian Lu. Fsdr: Frequency space domain randomization for domain generalization. In *CVPR*, pages 6891–6902, 2021. 2

[35] Zeyi Huang, Haohan Wang, Eric P Xing, and Dong Huang. Self-challenging improves cross-domain generalization. In *ECCV*, pages 124–140. Springer, 2020. 7

[36] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. *NeurIPS*, 28, 2015. 3

[37] Xin Jin, Cuiling Lan, Wenjun Zeng, Zhibo Chen, and Li Zhang. Style normalization and restitution for generalizable person re-identification. In *CVPR*, pages 3143–3152, 2020. 2, 7, 8, 13, 14

[38] Yitzhak Katznelson. *An introduction to harmonic analysis*. Cambridge University Press, 2004. 2

[39] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 14

[40] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy Hospedales. Learning to generalize: Meta-learning for domain generalization. In *AAAI*, volume 32, 2018. 2, 14

[41] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain generalization. In *ICCV*, pages 5542–5550, 2017. 5, 13, 14

[42] Da Li, Jianshu Zhang, Yongxin Yang, Cong Liu, Yi-Zhe Song, and Timothy M Hospedales. Episodic training for domain generalization. In *ICCV*, pages 1446–1455, 2019. 14

[43] Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain generalization with adversarial feature learning. In *CVPR*, pages 5400–5409, 2018. 7, 13, 14

[44] Hanjun Li, Gaojie Wu, and Wei-Shi Zheng. Combined depth space based architecture search for person re-identification. In *CVPR*, pages 6729–6738, 2021. 15, 16

[45] Shaohua Li, Kaiping Xue, Bin Zhu, Chenkai Ding, Xindi Gao, David Wei, and Tao Wan. Falcon: A fourier transform based approach for fast and secure convolutional neural network predictions. In *CVPR*, pages 8705–8714, 2020. 2

[46] Wei Li and Xiaogang Wang. Locally aligned feature transforms across views. In *CVPR*, pages 3594–3601, 2013. 5, 13

[47] Wei Li, Rui Zhao, Tong Xiao, and Xiaogang Wang. Deepreid: Deep filter pairing neural network for person re-identification. In *CVPR*, pages 152–159, 2014. 5, 13

[48] Xin Li, Zhizheng Zhang, Guoqiang Wei, Cuiling Lan, Wenjun Zeng, Xin Jin, and Zhibo Chen. Confounder identification-free causal visual feature learning. *CoRR*, abs/2111.13420, 2021. 2

[49] Ya Li, Xinmei Tian, Mingming Gong, Yajing Liu, Tongliang Liu, Kun Zhang, and Dacheng Tao. Deep domain generalization via conditional invariant adversarial networks. In *ECCV*, pages 624–639, 2018. 2

[50] Yiying Li, Yongxin Yang, Wei Zhou, and Timothy Hospedales. Feature-critic networks for heterogeneous domain generalization. In *ICML*, pages 3915–3924. PMLR, 2019. 2

[51] Shengcai Liao and Ling Shao. Interpretable and generalizable person re-identification with query-adaptive convolution and temporal lifting. In *ECCV*, pages 456–474. Springer, 2020. 7, 8, 13

[52] Ci-Siang Lin, Yuan-Chia Cheng, and Yu-Chiang Frank Wang. Domain generalized person re-identification via cross-domain episodic learning. In *ICPR*, pages 6758–6763. IEEE, 2021. 7

[53] Jiawei Liu, Zhipeng Huang, Liang Li, Kecheng Zheng, and Zheng-Jun Zha. Debiased batch normalization via gaussian process for generalizable person re-identification. In *AAAI*, 2022. 2, 5, 13

[54] Quande Liu, Cheng Chen, Jing Qin, Qi Dou, and Pheng-Ann Heng. Feddg: Federated domain generalization on medical image segmentation via episodic learning in continuous frequency space. In *CVPR*, pages 1013–1023, 2021. 2, 3

[55] Chen Change Loy, Tao Xiang, and Shaogang Gong. Multi-camera activity correlation analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1988–1995. IEEE, 2009. 5, 13

[56] Hao Luo, Youzhi Gu, Xingyu Liao, Shenqi Lai, and Wei Jiang. Bag of tricks and a strong baseline for deep person re-identification. In *CVPR Workshops*, pages 0–0, 2019. 13, 14, 15

[57] Massimiliano Mancini, Samuel Rota Bulo, Barbara Caputo, and Elisa Ricci. Best sources forward: domain generalization through source-specific nets. In *ICIP*, pages 1353–1357, 2018. 2

[58] Xintian Mao, Yiming Liu, Wei Shen, Qingli Li, and Yan Wang. Deep residual fourier transformation for single image deblurring. *arXiv preprint arXiv:2111.11745*, 2021. 2

[59] Michael Mathieu, Mikael Henaff, and Yann LeCun. Fast training of convolutional networks through ffts. *arXiv preprint arXiv:1312.5851*, 2013. 2

[60] Toshihiko Matsuura and Tatsuya Harada. Domain generalization using a mixture of multiple latent domains. In *AAAI*, volume 34, pages 11749–11756, 2020. 13, 14

[61] Diganta Misra, Trikay Nalamada, Ajay Uppili Arasanipalai, and Qibin Hou. Rotate to attend: Convolutional triplet attention module. In *WACV*, pages 3139–3148, 2021. 5

[62] Saeid Motiian, Marco Piccirilli, Donald A Adjeroh, and Gianfranco Doretto. Unified deep supervised domain adaptation and generalization. In *ICCV*, pages 5715–5725, 2017. 2, 7, 14

[63] Varsha Nair, Moitrayee Chatterjee, Neda Tavakoli, Akbar Siami Namin, and Craig Snoeyink. Fast fourier transformation for optimizing convolutional neural networks in object recognition. *arXiv preprint arXiv:2010.04257*, 2020. 2

[64] Xingang Pan, Ping Luo, Jianping Shi, and Xiaoou Tang. Two at once: Enhancing learning and generalization capacities via ibn-net. In *ECCV*, 2018. 2

[65] Mandela Patrick, Dylan Campbell, Yuki Asano, Ishan Misra, Florian Metze, Christoph Feichtenhofer, Andrea Vedaldi, and João F Henriques. Keeping your eye on the ball: Trajectory attention in video transformers. In *NeurIPS*, 2021. 3

[66] Ioannis Pitas. *Digital image processing algorithms and applications*. John Wiley & Sons, 2000. 2

[67] Anish Prabhu, Ali Farhadi, Mohammad Rastegari, et al. Butterfly transform: An efficient fft based neural architecture design. In *CVPR*, pages 12024–12033, 2020. 2

[68] Harry Pratt, Bryan Williams, Frans Coenen, and Yalin Zheng. Fcnn: Fourier convolutional neural networks. In *ECML PKDD*, pages 786–798, 2017. 2

[69] Zequn Qin, Pengyi Zhang, Fei Wu, and Xi Li. Fcanet: Frequency channel attention networks. In *ICCV*, pages 783–792, 2021. 3, 5

[70] Yongming Rao, Wenliang Zhao, Zheng Zhu, Jiwen Lu, and Jie Zhou. Global filter networks for image classification. In *NeurIPS*, volume 34, 2021. 2, 5

[71] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015. 13, 15

[72] Seonguk Seo, Yumin Suh, Dongwan Kim, Geeho Kim, Jongwoo Han, and Bohyung Han. Learning to optimize domain specific normalization for domain generalization. In *ECCV*, pages 68–83, 2020. 2

[73] Shiv Shankar, Vihari Piratla, Soumen Chakrabarti, Siddhartha Chaudhuri, Preethi Jyothi, and Sunita Sarawagi. Generalizing across domains via cross-gradient training. *ICLR*, 2018. 2, 7, 13, 14

[74] Jifei Song, Yongxin Yang, Yi-Zhe Song, Tao Xiang, and Timothy M Hospedales. Generalizable person re-identification by domain-invariant mapping network. In *CVPR*, pages 719–728, 2019. 7

[75] Yifan Sun, Changmao Cheng, Yuhan Zhang, Chi Zhang, Liang Zheng, Zhongdao Wang, and Yichen Wei. Circle loss: A unified perspective of pair similarity optimization. In *CVPR*, pages 6398–6407, 2020. 15

[76] Yifan Sun, Liang Zheng, Yi Yang, Qi Tian, and Shengjin Wang. Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline). In *ECCV*, pages 480–496, 2018. 15

[77] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. In *WACV*, pages 2149–2159, 2022. 2

[78] Jonathan Tremblay, Aayush Prakash, David Acuna, Mark Brophy, Varun Jampani, Cem Anil, Thang To, Eric Cameracci, Shaad Boochoon, and Stan Birchfield. Training deep networks with synthetic data: Bridging the reality gap by domain randomization. In *CVPR workshops*, pages 969–977, 2018. 2

[79] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008. 15

[80] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *CVPR*, pages 5018–5027, 2017. 14

[81] Riccardo Volpi, Hongseok Namkoong, Ozan Sener, John C Duchi, Vittorio Murino, and Silvio Savarese. Generalizing to unseen domains via adversarial data augmentation. *NeurIPS*, 31, 2018. 2

[82] Guanshuo Wang, Yufeng Yuan, Xiong Chen, Jiwei Li, and Xi Zhou. Learning discriminative features with multiple granularities for person re-identification. In *ACMMM*, pages 274–282, 2018. 15

[83] Haohan Wang, Xindi Wu, Zeyi Huang, and Eric P Xing. High-frequency component helps explain the generalization of convolutional neural networks. In *CVPR*, 2020. 1, 2, 3, 8

[84] Qilong Wang, Banggu Wu, Pengfei Zhu, Peihua Li, Wangmeng Zuo, and Qinghua Hu. Eca-net: Efficient channel attention for deep convolutional neural networks. In *CVPR*, pages 11531–11539, 2020. 3

[85] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, pages 7794–7803, 2018. 3, 5

[86] Guoqiang Wei, Cuiling Lan, Wenjun Zeng, Zhizheng Zhang, and Zhibo Chen. Toalign: task-oriented alignment for unsupervised domain adaptation. *NeurIPS*, 34:13834–13846, 2021. 2

[87] Longhui Wei, Shiliang Zhang, Wen Gao, and Qi Tian. Person transfer gan to bridge domain gap for person re-identification. In *CVPR*, pages 79–88, 2018. 5, 13, 16

[88] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *ECCV*, pages 3–19, 2018. 3, 4, 5

[89] Tong Xiao, Shuang Li, Bochao Wang, Liang Lin, and Xiaogang Wang. Joint detection and identification feature learning for person search. In *CVPR*, pages 3415–3424, 2017. 5, 13

[90] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017. 14

[91] Qinwei Xu, Ruipeng Zhang, Ya Zhang, Yanfeng Wang, and Qi Tian. A fourier-based framework for domain generalization. In *CVPR*, pages 14383–14392, 2021. 2, 3

[92] Zhiqin John Xu. Understanding training and generalization in deep learning by fourier analysis. *arXiv preprint arXiv:1808.04295*, 2018. 1, 2, 3, 8

[93] Zhi-Qin John Xu, Yaoyu Zhang, and Yanyang Xiao. Training behavior of deep neural network in frequency domain. In *ICONIP*, 2019. 1, 2, 3, 8

[94] Cheng Yan, Guansong Pang, Xiao Bai, Changhong Liu, Ning Xin, Lin Gu, and Jun Zhou. Beyond triplet loss: person re-identification with fine-grained difference-aware pairwise loss. *IEEE Trans Multimedia*, 2021. 15

[95] Yanchao Yang and Stefano Soatto. Fda: Fourier domain adaptation for semantic segmentation. In *CVPR*, pages 4085–4095, 2020. 2, 3

[96] Qiaosi Yi, Jinhao Liu, Le Hu, Faming Fang, and Guixu Zhang. Contrastive learning for local and global learning mri reconstruction. *arXiv preprint arXiv:2111.15200*, 2021. 2

[97] Dong Yin, Raphael Gontijo Lopes, Jon Shlens, Ekin Dogus Cubuk, and Justin Gilmer. A fourier perspective on model robustness in computer vision. *NeurIPS*, 2019. 1, 2

[98] Xiangyu Yue, Yang Zhang, Sicheng Zhao, Alberto Sangiovanni-Vincentelli, Kurt Keutzer, and Boqing Gong. Domain randomization and pyramid consistency: Simulation-to-real generalization without accessing target domain data. In *ICCV*, pages 2100–2110, 2019. 2

[99] Hang Zhang, Chongruo Wu, Zhongyue Zhang, Yi Zhu, Haibin Lin, Zhi Zhang, Yue Sun, Tong He, Jonas Mueller, R Manmatha, et al. Resnest: Split-attention networks. *arXiv preprint arXiv:2004.08955*, 2020. 14

[100] Yi-Fan Zhang, Hanlin Zhang, Zhang Zhang, Da Li, Zhen Jia, Liang Wang, and Tieniu Tan. Learning domain invariant representations for generalizable person re-identification. *arXiv preprint arXiv:2103.15890*, 2021. 7

[101] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, and Zhibo Chen. Densely semantically aligned person re-identification. In *CVPR*, pages 667–676, 2019. 13

[102] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, and Zhibo Chen. Multi-granularity reference-aided attentive feature aggregation for video-based person re-identification. In *CVPR*, pages 10407–10416, 2020. 3

[103] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, Zhibo Chen, and Shih-Fu Chang. Beyond triplet loss: Meta prototypical n-tuple loss for person re-identification. *IEEE Trans Multimedia*, 2021. 15

[104] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, Xin Jin, and Zhibo Chen. Relation-aware global attention for person re-identification. In *CVPR*, pages 3186–3195, 2020. 3, 5, 15

[105] Cairong Zhao, Xinbi Lv, Zhang Zhang, Wangmeng Zuo, Jun Wu, and Duoqian Miao. Deep fusion feature representation learning with hard mining center-triplet loss for person re-identification. *IEEE Trans Multimedia*, pages 3180–3195, 2020. 15, 16

[106] Yuyang Zhao, Zhun Zhong, Fengxiang Yang, Zhiming Luo, Yaojin Lin, Shaozi Li, and Nicu Sebe. Learning to generalize unseen domains via memory-based multi-source meta-learning for person re-identification. In *CVPR*, 2021. 5, 13

[107] Yuyang Zhao, Zhun Zhong, Fengxiang Yang, Zhiming Luo, Yaojin Lin, Shaozi Li, and Nicu Sebe. Learning to generalize unseen domains via memory-based multi-source meta-learning for person re-identification. In *CVPR*, pages 6277–6286, 2021. 7, 8, 14

[108] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *ICCV*, pages 1116–1124, 2015. 5, 13, 16

[109] Wei-Shi Zheng, Shaogang Gong, and Tao Xiang. Associating groups of people. In *BMVC*, volume 2, pages 1–11, 2009. 5, 13

[110] Zhedong Zheng, Xiaodong Yang, Zhiding Yu, Liang Zheng, Yi Yang, and Jan Kautz. Joint discriminative and generative learning for person re-identification. In *CVPR*, pages 2138–2147, 2019. 15, 16

[111] Zhedong Zheng, Liang Zheng, and Yi Yang. Unlabeled samples generated by gan improve the person re-identification baseline in vitro. In *Proceedings of the IEEE international conference on computer vision*, pages 3754–3762, 2017. 13, 16

[112] Kaiyang Zhou, Yongxin Yang, Andrea Cavallaro, and Tao Xiang. Omni-scale feature learning for person re-identification. In *ICCV*, pages 3702–3712, 2019. 15, 16

[113] Kaiyang Zhou, Yongxin Yang, Timothy Hospedales, and Tao Xiang. Deep domain-adversarial image generation for domain generalisation. In *AAAI*, pages 13025–13032, 2020. 2

[114] Kaiyang Zhou, Yongxin Yang, Yu Qiao, and Tao Xiang. Domain adaptive ensemble learning. *TIP*, 2021. 2, 7, 13

[115] Ying Zhu, Runwei Ding, Weibo Huang, Peng Wei, Ge Yang, and Yong Wang. Hmfca-net: Hierarchical multi-frequency based channel attention net for mobile phone surface defect detection. *Pattern Recognit Lett*, 153:118–125, 2022. 3

[116] Zijie Zhuang, Longhui Wei, Lingxi Xie, Tianyu Zhang, Hengheng Zhang, Haozhe Wu, Haizhou Ai, and Qi Tian. Rethinking the distribution gap of person re-identification with camera-based batch normalization. In *ECCV*, pages 140–157. Springer, 2020. 7

# Supplementary Material

## 6. More Datasets and Implementation Details

We evaluate the effectiveness of our proposed Deep Frequency Filtering (DFF) for Domain Generalization (DG) on **Task-1**: the close-set classification task and **Task-2**: the open-set retrieval task, *i.e.*, person re-identification (ReID). More details about the datasets and our experiment configurations are introduced in this section.

### 6.1. Datasets for Task-1

We use the most commonly used Office-Home [41] and PACS dataset [41]. Specifically, Office-Home consists of 4 domains (Art (Ar), Clip Art (Cl), Product (Pr), Real-World (Rw)), each consisting of 65 categories, with an average of 70 images per category, for a total of 15,500 images. PACS consists of 9991 samples in total from 4 domains (*i.e.*, Photo (P), Art Painting (A), Cartoon (C) and Sketch (S)). All these 4 domains share 7 object categories. They are commonly used domain generalization (DG) benchmark on the task of classification. We validate the effectiveness of our proposed method for generalization in close-set classification task on Office-Home and PACS. Following the typical setting, we conduct experiments on this dataset under the leave-one-out protocol (see Table 8 Protocol-1 and Protocol-2), where three domains are used for training and the remaining one is considered as the unknown target domain.

### 6.2. Datasets for Task-2

Person re-identification (ReID) is a representative open-set retrieval task, where different domains and datasets do not share their label space. We employ existing ReID protocols to evaluate the generalization ability of our method. *i*) For Protocol-3 and Protocol-4, we also follow the leave-one-out protocols as in [53, 106]. Among the four datasets (CUHK-SYSU (CS) [89], MSMT17 (MS) [87], CUHK03 (C3) [47] and Market-1501 (MA) [108]), three are selected as the seen domain for training and the remaining one is selected the unseen domain data for testing. Differently, Protocol-3 only adopts the training set of seen domains for model training while in Protocol-3, the testing set of the seen domains are also included for training model. *ii*) For Protocol-5 in Table 8, several large-scale ReID datasets *e.g.*, CUHK02 (C2) [46], CUHK03 (C3) [47], Market-1501 (MA) [108] and CUHK-SYSU (CS) [89], are viewed as multiple source domains. Each small-scale ReID dataset including VIPeR [25], PRID [30], GRID [55] and iLIDS [109] is used as an unseen target domain, respectively. To comply with the General Ethical Conduct, we exclude DukeMTMC from the source domains. The final performance is obtained by averaging 10 repeated experiments with random splits of training and testing sets.

Table 8. The evaluation protocols. "Com-" refers to combining the train and test sets of source domains for training. "Pr", "Ar", "Cl", "Rw" are short for the Product, Art, Clip Art and Real-World domains in Office-Home dataset [41], respectively. "P", "A", "C", "S" are short for the Photo, Painting, Cartoon, Sketch domains in PACS dataset [41], respectively. "MA", "CS", "C3", "MS" denote Market-1501 [108], CUHK-SYSU [89], CUHK03 [47], MSMT17 [87], respectively. Note that for person ReID, the commonly used DukeMTMC [111] has been withdrawn by its publisher, is thus no longer used.

| Task | Setting | Training Data | Testing Data |
|---|---|---|---|
| Close-set classification | Protocol-1 | Cl,Pr,Rw | Ar |
| | | Ar,Pr,Rw | Cl |
| | | Ar,Cl,Rw | Pr |
| | | Ar,Cl,Pr | Rw |
| | Protocol-2 | C,P,S | A |
| | | A,P,S | C |
| | | A,C,P | S |
| | | A,C,S | P |
| Open-set retrieval | Protocol-3 | CS+C3+MS | MA |
| | | MA+CS+MS | C3 |
| | | MA+CS+C3 | MS |
| | Protocol-4 | Com-(CS+C3+MT) | MA |
| | | Com-(MA+CS+MS) | C3 |
| | | Com-(MA+CS+C3) | MS |
| | Protocol-5 | Com-(MA+C2+C3+CS) | PRID |
| | | | GRID |
| | | | VIPeR |
| | | | iLIDs |

### 6.3. Networks

Following the common practices of domain generalizable classification (Task-1) [8, 43, 73, 114] and person ReID (Task-2) [9, 14, 37, 51], we build the networks equipped with our proposed Deep Frequency Filtering (DFF) for these two tasks on the basis of ResNet-18 and ResNet-50, respectively. As introduced in the Sec. 3.4 of our manuscript, we evaluate the effectiveness of our proposed DFF based on the two-branch architecture of Fast Fourier Convolution (FFC) in [11]. In particular, we adopt our DFF operation to the spectral transformer branch of this architecture. Unless otherwise stated, the ratio $r$ in splitting $\mathbf{X} \in \mathbb{R}^{C \times H \times W}$ into $\mathbf{X}^g \in \mathbb{R}^{rC \times H \times W}$ and $\mathbf{X}^l \in \mathbb{R}^{(1-r)C \times H \times W}$ is set to 0.5. We conduct an ablation study on this ratio in this Supplementary as follows.

### 6.4. Training

Following common practices [8, 27, 43, 56, 60, 101], we adopt ResNet-18 and ResNet-50 [26] as our backbone for Task-1 and Task-2, respectively. Each convolution layer of the backbone is replaced with our DDF module. Unless specially stated, we first pretrain the models on ImageNet [71] then fine-tune them on Task-1 or Task-2, referring to the common practices [8, 9, 14, 73, 114]. We introduce our training configurations with more details in the following.

**Pre-training on ImageNet.** Following the common practices [26, 33, 90, 99] in this field, we adopt the commonly used data augmentation strategies including color jittering, random flipping and center cropping. The input image size is $224 \times 224$. We use the SGD optimizer with the base initial learning rate of 0.4, the momentum of 0.9 and the weight decay of 0.0001, and perform learning rate decay by a factor of 0.1 after 30, 60 and 80 epochs. A linear warm-up strategy is adopted in the first 5 epochs, where the learning rate is increased from 0.0 to 0.4 linearly. All models are trained for 90 epochs with the batchsize of 1024.

**Fine-tuning on Task-1.** The initial learning rate in this stage is set to 0.001. We train all models on this task using the SGD optimizer with the momentum of 0.9 and the weight decay of 0.0001. The batch size is set to 32. Following prior works [8, 43, 62, 73], we adopt the data augmentation strategies including random cropping, horizontal flipping, and random grayscale. The input images are resized to $224 \times 224$. On the PACS dataset [41], we train the models for 3,500 iterations; and on the Office-Home [80] dataset, we train the models for 3,000 iterations. The experiment results on Office-Home have been presented in the main paper while the results on PACS are placed in this Supplementary due to the length limitation.

**Fine-tuning on Task-2.** Following the common practices for domain generalizable person ReID [14, 16, 37, 107], we adopt the widely used data augmentation strategies, including cropping, random flipping, and color jittering. We use Adam [39] optimizer with the momentum of 0.9 and weight decay of 0.0005. The learning rate is initialized by $3.5 \times 10^{-4}$ and decayed using a cosine annealing schedule. The batch size is set to 128, including 8 identities and 16 images per identity. For the Protocol-3, Protocol-4 and Protocol-5, the models are trained for 60 epochs on their corresponding source datasets. Similar to previous work [56], the last spatial down-sampling in the "conv5_x" block is removed. The input images are resized to $384 \times 128$. Following [27], we use task-related loss including cross-entropy loss, arcface loss, circle loss and triplet loss. And we adopt a gradient reversal layer [22] encouraging the learning of domain-invariant features.

## 7. More Experiment Results

In this section, we present more experiment results to further evaluate the effectiveness of our proposed DFF.

### 7.1. More Experiments on the Task-1 (PACS)

We further evaluate the effectiveness of our proposed DFF on another commonly used dataset, *i.e.*, Office-Home [80], for investigating the domain generalization on the

Table 9. Performance (classification accuracy %) comparison with the state-of-the-art methods under Protocol-2 (*i.e.*, on PACS dataset) on close-set classification task. We use ResNet-18 as backbone. Best in bold.

| Method | Source→Target | | | | Avg |
| --- | --- | --- | --- | --- | --- |
| | C,P,S→A | A,P,S→C | A,C,P→S | A,C,S→P | |
| Baseline | 77.6 | 73.9 | 70.3 | 94.4 | 79.1 |
| MMD-AAE [43] | 75.2 | 72.7 | 64.2 | 96.0 | 77.0 |
| CrossGrad [73] | 78.7 | 73.3 | 65.1 | 94.0 | 77.8 |
| MetaReg [3] | 79.5 | 75.4 | 72.2 | 94.3 | 80.4 |
| JiGen [8] | 79.4 | 77.3 | 71.4 | 96.0 | 81.0 |
| MLDG [40] | 79.5 | 75.3 | 71.5 | 94.3 | 80.7 |
| MASF [19] | 80.3 | 77.2 | 71.7 | 95.0 | 81.1 |
| Epi-FCR [42] | 82.1 | 77.0 | **73.0** | 93.9 | 81.5 |
| MMLD [60] | 81.3 | 77.2 | 72.3 | **96.1** | 81.7 |
| Ours | **82.2** | **78.5** | 72.5 | 95.5 | **82.2** |

Table 10. Performance comparisons of different frequency transformations. In *"Baseline"*, we take vanilla ResNet-18/-50 as the backbone models. *"Wavelet (db3)"* and *"Wavelet (Haar)"* denote the wavelet transforms with the Daubechies3 and Haar filters, respectively.

| Method | Source→Target | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | MS+CS+C3→MA | | MS+MA+CS→C3 | | MA+CS+C3→MS | |
| | mAP | R1 | mAP | R1 | mAP | R1 |
| Base | 59.4 | 83.1 | 30.3 | 29.1 | 18.0 | 41.9 |
| Wavelet (db3) | 61.5 | 83.7 | 30.7 | 29.8 | 18.3 | 42.2 |
| Wavelet (Haar) | 61.1 | 83.6 | 30.5 | 29.7 | 18.5 | 42.3 |
| FFT (Ours) | **71.1** | **87.1** | **41.3** | **41.1** | **25.1** | **50.5** |

Table 11. Performance comparisons of different dimensions on which the Fast Fourier Transform (FFT) is performed. *"FFT (CHW)"* refers to the models in which FFT is performed across the height (H), width (W) and channel (C) dimensions. In *"FFT (HW)"*, we just perform FFT across the height and width dimensions, *i.e.*, for each feature map independently, which is the default setting in this paper.

| Method | Source→Target | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | MS+CS+C3→MA | | MS+MA+CS→C3 | | MA+CS+C3→MS | |
| | mAP | R1 | mAP | R1 | mAP | R1 |
| Base | 59.4 | 83.1 | 30.3 | 29.1 | 18.0 | 41.9 |
| FFT (CHW) | 59.2 | 83.0 | 30.0 | 28.8 | 17.5 | 38.5 |
| FFT(HW) | **71.1** | **87.1** | **41.3** | **41.1** | **25.1** | **50.5** |

close-set classification. This dataset contains four domains (Aritistic, Clipart, Product and Real World) with 15,500 images of classes for home and office object recognition. Similar to the Protocol-1 on PACS dataset [41], we adopt a "Leave-One-Out" protocol for the evaluation on Office-Home where three domains are used for training while the remaining one is for testing. The experiment results are shown in Table 9. Our proposed DFF achieves significant improvements relative to the *Baseline* model, and outperforms the state-of-the-art methods on this dataset by a clear margin over all evaluation settings. This further demonstrates the effectiveness of DFF.
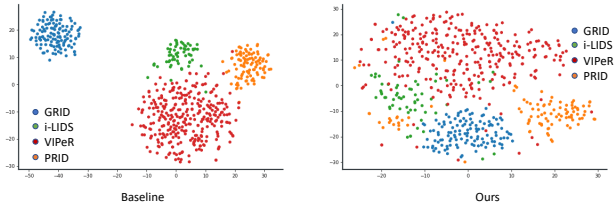
Figure 4. The t-SNE [79] visualization of ReID feature vectors learned by baseline (left) and our DFF (right) on four unseen target datasets (GRID, i-LIDS, VIPeR and GRID). Best viewed in color.

## 7.2. More Ablation Studies

**Experiments with other frequency transforms.** We preliminarily investigate the effectiveness of using other frequency transforms in implementing our conceptualized DFF. In particular, we replace the Fast Fourier Transform (FFT) in our proposed scheme by the wavelet transform with two widely used filters, *i.e.*, db3 and Haar. From the experiment results in Table 10, we observe that adopting the wavelet transform also delivers improvements compared to *Baseline*, but is inferior to adopting FFT. This is because the wavelet transform is a low frequency transformation such that our proposed filtering operation is performed in a local space, thus limiting the benefits of DFF.

**Design choices of performing FFT.** In our proposed scheme, for the given intermediate feature $\mathbf{X} \in \mathbb{R}^{C \times H \times W}$, we perform FFT for each channel independently to obtain the latent frequency representations, as described in the Sec. 3.2 of the main paper. Here, we investigate other design choices of perform FFT. In the Table 11, we find that performing FFT across H, W, C dimensions leads to performance drop compared to *Baseline*. For the intermediate feature $\mathbf{X} \in \mathbb{R}^{C \times H \times W}$, its different channels correspond to the outputs of different convolution kernels, which are independent in fact. Thus, we perform FFT on each channel of $\mathbf{X}$ independently.

**Ablation study on the ratio $r$.** We follow the overall architecture design of [11] to split the given intermediate feature $\mathbf{X} \in \mathbb{R}^{C \times H \times W}$ into $\mathbf{X}^g \in \mathbb{R}^{rC \times H \times W}$ and $\mathbf{X}^l \in \mathbb{R}^{(1-r)C \times H \times W}$ along its channel dimension with a ratio $r \in [0, 1]$. Our proposed filtering operation is only performed on $\mathbf{X}^g$. When setting $r = 0$, the models degenerate to the ResNet-18/-50 baselines. Setting $r = 1$ means that we perform DFF on the entire intermediate feature $\mathbf{X}$. As the experiment results in Table 12, we empirically find that the models with $r = 0.5$ achieve the best performance, exploiting the complementary of features in the frequency and original spaces.

Table 12. Performance comparisons of our proposed DFF with different ratios. All models are built based on ResNet-18 for Task-1 while ResNet-50 for Task-2.

| Ratio | Source→Target | | | | | |
| | MS+CS+C3→MA | | MS+MA+CS→C3 | | MA+CS+C3→MS | |
| | mAP | R1 | mAP | R1 | mAP | R1 |
|---|---|---|---|---|---|---|
| 0.0 | 59.4 | 83.1 | 30.3 | 29.1 | 18.0 | 41.9 |
| 0.25 | 67.4 | 84.1 | 38.1 | 38.1 | 22.9 | 48.4 |
| **0.5(Ours)** | **71.1** | **87.1** | **41.3** | **41.1** | **25.1** | **50.5** |
| 0.75 | 70.8 | 86.8 | 40.7 | 40.6 | 21.0 | 44.9 |
| 1.0 | 64.2 | 83.4 | 29.3 | 28.1 | 17.6 | 40.4 |

Table 13. Performance comparisons of our proposed DFF with the corresponding ResNet baselines on ImageNet-1K classification. "*DFF-ResNet-18/-50*" denote the ResNet-18/-50 models equipped with our DFF.

| Method | Parameters | GFLOPs | Top-1 Acc. |
|---|---|---|---|
| ResNet-18 | 11.7M | 1.8 | 69.8 |
| DFF-ResNet-18 | 12.2M | 2.0 | 72.3 |
| ResNet-50 | 25.6M | 4.1 | 76.3 |
| DFF-ResNet-50 | 27.7M | 4.5 | 77.9 |

Table 14. Performance comparisons of our proposed DFF with the state-of-the-art methods on supervised person ReID. "*Base.*" refers to the baseline model.

| Model | Market-1501(MA) | | MSMT17(MT) | |
| | mAP | R1 | mAP | R1 |
|---|---|---|---|---|
| PCB [76] | 81.60 | 93.80 | - | - |
| BoT [56] | 85.90 | 94.50 | - | - |
| MGN [82] | 86.90 | 95.70 | - | - |
| JDGL [110] | 86.00 | 94.80 | 52.30 | 77.20 |
| GASM [28] | 84.70 | 95.30 | 52.50 | 79.50 |
| FPR [29] | 86.58 | 95.42 | - | - |
| HCTL [105] | 81.80 | 93.80 | 43.60 | 74.30 |
| OSNet [112] | 84.90 | 94.80 | 52.90 | 78.70 |
| RGA-SC [104] | 88.40 | 96.10 | 57.50 | 80.30 |
| CDNet [44] | 86.00 | 95.10 | 54.70 | 78.90 |
| Circle Loss [75] | 87.40 | 96.10 | 52.10 | 76.90 |
| AMD [10] | 87.15 | 94.74 | - | - |
| FIDI [94] | 86.80 | 94.50 | - | - |
| MPN-tuple [103] | 88.70 | 95.30 | 60.10 | 82.20 |
| ResNet-50 Base. | 81.63 | 93.89 | 50.84 | 76.78 |
| DFF-ResNet-50 | **90.21** | **96.17** | **60.21** | **82.95** |

## 7.3. More Visualization Results

We perform t-SNE visualization for the ReID feature vectors extracted by the baseline model and the model with our proposed DFF on four unseen datasets. As shown in Fig. 4, the four unseen target domains distribute more separately for the baseline model than that of ours. This indicates the domain gaps are effectively mitigated by our proposed Deep Frequency Filtering (DFF).

## 7.4. Effectiveness on ImageNet-1K Classification

ImageNet-1K [71] classification widely serves as a pre-training task, providing pre-trained weights as the model initialization for various downstream task. We present the effectiveness of our conceptualized DFF on ImageNet-

1K classification to showcase its potentials for more general purposes. As the results shown in Table 13, our DFF achieves 2.5%/1.6% improvements on the Top-1 classification accuracy compared to the corresponding baselines ResNet-18 and ResNet-50, respectively. Note that these improvements are achieved with the simple instantiation introduced in the Sec.3.3 of the main body. We believe more effective instantiations of DFF are worth exploring to make DFF contribute more in a wider range of fields.

### 7.5. Effectiveness on Supervised Person ReID

In the main body, we target domain generalization and present the effectiveness of our proposed DFF on domain generalizable person ReID. In this supplementary material, we also showcase its potential on improving supervised person ReID. Following the previous works [28, 44, 105, 110, 112] in this field, we evaluate our DFF on two most widely used datasets Market-1501 [108] and MSMT17 [87]. Note that another popular dataset DukeMTMC [111] has been taken down by its publisher. As shown in Table 14, the ResNet-50 equipped with DFF significantly outperforms the baseline model and reaches the SOTA performance on this task. This demonstrates the proposed DFF is also potentially beneficial for capturing discriminative features. We expect that it can contribute to more tasks.