

Learning Gabor Texture Features for Fine-Grained Recognition

Lanyun Zhu¹ Tianrun Chen² Jianxiong Yin³ Simon See³ Jun Liu^{1*}

Singapore University of Technology and Design¹ Zhejiang University² NVIDIA AI Tech Centre³

lanyun.zhu@mymail.sutd.edu.sg tianrun.chen@zju.edu.cn

{jianxiongy, ssee}@nvidia.com jun_liu@sutd.edu.sg

Abstract

Extracting and using class-discriminative features is critical for fine-grained recognition. Existing works have demonstrated the possibility of applying deep CNNs to exploit features that distinguish similar classes. However, CNNs suffer from problems including frequency bias and loss of detailed local information, which restricts the performance of recognizing fine-grained categories. To address the challenge, we propose a novel texture branch as complementary to the CNN branch for feature extraction. We innovatively utilize Gabor filters as a powerful extractor to exploit texture features, motivated by the capability of Gabor filters in effectively capturing multi-frequency features and detailed local information. We implement several designs to enhance the effectiveness of Gabor filters, including imposing constraints on parameter values and developing a learning method to determine the optimal parameters. Moreover, we introduce a statistical feature extractor to utilize informative statistical information from the signals captured by Gabor filters, and a gate selection mechanism to enable efficient computation by only considering qualified regions as input for texture extraction. Through the integration of features from the Gabor-filter-based texture branch and CNN-based semantic branch, we achieve comprehensive information extraction. We demonstrate the efficacy of our method on multiple datasets, including CUB-200-2011, NA-bird, Stanford Dogs, and GTOS-mobile. State-of-the-art performance is achieved using our approach.

1. Introduction

Recognizing fine-grained categories (e.g., different flower types or bird species) is challenging due to the subtle differences in their visual appearance. Recently, deep neural network technology [17, 8, 31, 57, 52, 12] has made remarkable advancement, allowing for significant progress in extracting deep class-discriminative features for fine-grained recognition [10, 38, 41, 23, 3, 56]. For example, [54] develops a graph-based module to extract cross-class relationship features. [13, 26] use bilinear pooling to ex-

ploit high-order features from different networks. [55, 38] leverage multi-scale information to enhance feature effectiveness. The success of these methods highlights the crucial role of extracting powerful features in distinguishing fine-grained categories.

To exploit class-discriminative features, previous methods [38, 41, 23, 3, 56] have primarily relied on convolutional neural networks (CNNs) as the model backbone. CNNs have achieved significant success in many computer vision tasks. However, due to the following limitations, they may not be sufficient for capturing comprehensive information in fine-grained recognition. First, deep CNNs generally extract information with high receptive fields and resolution-reduced feature maps, leading to the loss of local detailed information that is critical for fine-grained recognition. Second, as discussed in [33, 45, 46], due to the regularity of the commonly used activation [45] and loss functions [33], CNNs prioritize learning low-frequency components while ignoring high-frequency components. This bias is detrimental to fine-grained classification, as the ignored high-frequency information, such as ripples and spots in bird feathers, can be important in distinguishing between similar classes. Considering these limitations of CNNs, we suggest that a more powerful feature extractor is required, which should be able to capture local detailed information and multi-frequency comprehensive features that are lost by CNNs but crucial for fine-grained recognition.

Witnessing the challenge, in this work, we introduce a novel branch that extracts texture features as complementary to the vanilla CNN features to facilitate fine-grained recognition. We employ Gabor filters as the texture extractor in the novel branch. The Gabor filter is a windowed Fourier transform that combines a sinusoidal signal with a Gaussian wave. It is widely used in image processing for extracting low-level texture features [30, 35, 29, 1, 40]. The motivations of using Gabor filters in fine-grained recognition can be illustrated in two-folds. (1) First, Gabor filters are effective in capturing local-detailed information. Previous research [5, 6] has found that Gabor functions are similar to the receptive field profiles in the mammalian cortical sim-

*Corresponding Author

ple cells. This indicates that Gabor filters can perform the similar function as cortical cells to capture effective and locally detailed texture information, which is critical for recognizing confusing samples. (2) Second, Gabor filters can extract sufficient high-frequency information. A Gabor filter contains a parameter that directly controls the frequency of the extracted information. It is therefore possible to exploit sufficient high-frequency components by constraining the parameters of some filters into a high-value range, overcoming the frequency-bias limitations of traditional CNNs. Benefiting from these advantages, our introduction of Gabor filters enables texture features to be extracted as an effective supplement to the typical CNN features for improving fine-grained recognition.

Utilizing Gabor filters in an effective manner is non-trivial. Previous methods with Gabor filters generally rely on hand-crafted designs for filter parameter setting, which cannot guarantee to be optimal. To overcome this limitation, we propose an approach to automatically learn parameters in a supervised manner, in which a value constraint strategy is used to stabilize the learning and enhance high-frequency information extraction. To extract effective information with Gabor filters, we further propose a statistical feature extractor to capture texture features from the Gabor response maps. We also introduce a gate mechanism that selects a limited number of appropriate regions as inputs to the texture branch, thereby reducing computational costs and avoiding information redundancy. These careful designs allow Gabor filters to benefit fine-grained recognition in an effective and efficient manner.

To the best of our knowledge, our work pioneers to learn texture features using Gabor filters in the task of fine-grained recognition. We perform extensive experiments on four datasets, which demonstrate the effectiveness of our method. Benefiting from the powerful texture features from the Gabor filters, our network can leverage the informative and comprehensive information for recognition, resulting in state-of-the-art (SOTA) performance on all datasets. To summarize, our contributions are as follows. (1) First, we pioneer learning Gabor filters in fine-grained recognition. The texture features extracted from Gabor filters serve as an effective supplement to the vanilla CNN features, significantly improving recognition performance. (2) Second, to enhance model effectiveness and efficiency, we implement several key designs based on Gabor filters, including the parameter learning method, a powerful Gabor feature extractor and a region selection gate. (3) Third, extensive experiments demonstrate the effectiveness of our method with state-of-the-art (SOTA) performance achieved on four datasets.

2. Related Work

Fine-Grained Recognition. The goal of fine-grained recognition is to classify classes that are visually similar.

Recent methods for the task can be categorized into two major types: feature-encoding methods [13, 49, 26, 9, 10, 38, 41, 23, 3, 56] and localization methods [20, 43, 2, 21, 44, 7, 36, 11]. Feature-encoding methods emphasize on extracting effective features for the better classification results. For example, some methods [13, 49, 26] use bilinear pooling to obtain informative high-level features. Other works [9, 10] use feature constraints to enhance feature effectiveness. [38, 41, 23, 3] extract multi-level features to get comprehensive information. [54] exploits cross-class relationships to capture class-discriminative features. Localization methods emphasize on selecting the most effective regions to capture features. [20, 43, 2] use extra annotations such as bounding boxes and key points to promote the localization ability of the network. [21, 44, 7, 36] select informative regions through activation maps. Different from these works, our method innovatively learns Gabor filter to effectively capture texture information. Additionally, the region selection gate in our method can be end-to-end trained without the need for additional annotation, giving it a unique advantage.

Gabor Filters. The Gabor filter is an effective tool for extracting texture features, which makes it widely used in various computer vision applications such as face recognition [30, 35, 29], vehicle verification [1], object detection [42, 22, 39] and gait recognition [40]. Recently, some works [18, 32, 50, 53] introduce Gabor filters into deep neural networks. However, these approaches have suffered from problems like unstable training, frequency bias and high computation costs, leading to low effectiveness and efficiency. Our method addresses the challenges faced by previous approaches through our careful designs, making Gabor filters to be effectively applied in fine-grained recognition.

3. Method

3.1. Overall Structure

Fig. 1 provides an overview of our approach. A semantic branch takes the entire image as input and generates a semantic feature \mathbf{M} through a vanilla CNN. Using intermediate features from the semantic branch, a gate mechanism selects impactful regions from the image. After that, each of the selected regions i is zoomed into the size $S \times S$ and fed into a texture branch to produce a texture feature \mathbf{T}_i . The texture branch comprises N learnable Gabor filters followed by a statistical feature extractor that performs the following two-step workflow. First, each learnable Gabor filter is convolved with the input region to generate an intensity map \mathbf{I} . In this way, for all the N learnable Gabor filters, a set of maps $\{\mathbf{I}^n\}_{n=1}^N$ are generated. Next, the statistical feature extractor processes each $\mathbf{I}^n \in \{\mathbf{I}^n\}_{n=1}^N$ through a Learnable Histogram Operator (LHO) and then passes the resulting $\{\mathbf{S}^n\}_{n=1}^N$ through a Filter Correlation Module (FCM) to generate the texture feature \mathbf{T} . Finally, the semantic feature \mathbf{M} and texture features $\{\mathbf{T}_i\}$ from all

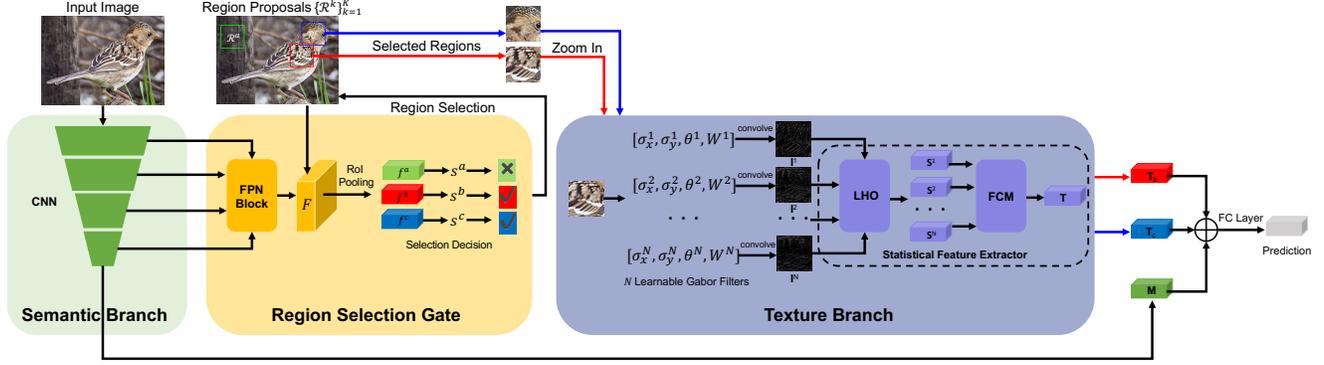


Figure 1. **Overall structure of our network.** The semantic branch produces a semantic feature M for the input image. The region selection gate selects impactful regions from a set of pre-defined region proposals. The texture branch generates a texture feature T for each of the selected regions that is zoomed into the size $S \times S$. Finally, semantic feature M and texture features T from all selected regions are added for prediction. Note that, for simplicity of illustration, the figure only presents three region proposals for an image. In actual practice, the region proposals follow the same setting as RCNN [15], comprising multiple regions with different positions and scales.

selected regions are added and fed into a fully connected layer to obtain the prediction.

The subsequent sections are arranged as follows. In Sec. 3.2, we provide an introduction to Gabor filters, which are the key to our method to capture texture information. In Sec. 3.3, we illustrate the strategy to train our Gabor filters in a stable and effective manner. In Sec. 3.4, we elaborate on the statistical feature extractor that exploits informative statistical information from the output of Gabor filters. In Sec. 3.5, we introduce the gate mechanism that selects impactful regions into the texture branch for texture extraction.

3.2. Gabor Filter

The Gabor filter is a widely-used feature extractor in traditional image processing, owing to its biological relevance and powerful texture extraction ability. It has optimal joint localization in both spatial and frequency domains, which makes it possible to be effectively applied in both spatial- and frequency-wise analysis. In the spatial domain, a Gabor Filter is mathematically defined as follows:

$$g(x, y) = \frac{1}{2\pi\sigma_x\sigma_y} \exp\left\{-\frac{1}{2}\left(\frac{\tilde{x}^2}{\sigma_x^2} + \frac{\tilde{y}^2}{\sigma_y^2}\right)\right\} \exp[2\pi jW\tilde{x}]. \quad (1)$$

$$\begin{cases} \tilde{x} = x \cos \theta + y \sin \theta \\ \tilde{y} = -x \sin \theta + y \cos \theta \end{cases} \quad (2)$$

This representation can be regarded as a windowed Fourier transform, i.e., a Gaussian kernel function modulated by a sinusoidal plane wave. Specifically, $[\sigma_x, \sigma_y]$ defines the scale of the Gabor filter, and it can determine the effective size of a pixel's neighborhood where weighted summation occurs. θ defines the orientation of the filter. W denotes the radial frequency of the sinusoid. These four parameters $\{\sigma_x, \sigma_y, \theta, W\}$ can fully define a Gabor filter. By altering the parameters, different filters can be generated to capture different signal components.

3.3. Learnable Gabor Filter Parameters

An effective Gabor filter requires the appropriate combination of parameters to match a given task. Previous works [40, 30, 42] have generally used a hand-crafted approach for parameter selection, relying on heuristic rules to inform manual parameter settings. However, this kind of hand-crafted approach is highly dependent on the expertise of the individual and cannot guarantee effectiveness for the task at hand. To address this limitation, we propose an alternative method that automatically learns the optimal parameters in a supervised manner, ensuring that the parameter combination can better fit the task. However, learning the parameters presents significant challenges. Firstly, during experimentation, erratic loss fluctuations are observed in the training process, indicating a lack of training stability (see the experimental section for details). Secondly, the features extracted by the trained network exhibit a substantial bias towards low-frequency information. The high-frequency information, on the other hand, occupies only a small fraction of all features. Specifically, only 6% of the trained Gabor filters have a frequency parameter W higher than 0.5. This tendency is mainly due to the regularity of commonly used activation and loss functions in deep learning (refer to [45, 33] for details), which is detrimental since the high-frequency information can be vital for fine-grained recognition.

To ensure training stability and filter effectiveness, we propose two solutions to tackle the above problems. To address the first challenge, we propose constraining parameter values into a valid range, which is theoretically derived to avoid potential problems such as frequency aliasing that could hinder training stability. To address the second challenge, we propose a dual-range strategy to further regularize the frequency parameter W , ensuring that sufficient high-frequency components can be extracted. These solutions enable our Gabor filters to be learned in a stable and ef-

fective manner. A detailed illustration of these solutions is presented in the subsequent sections.

Learning Parameters under Constraints. We first introduce a novel parameter value constraint to enhance the stability of model training. Our approach is motivated by the analysis that training instability could be attributed to the mismatch between image properties and the ideal application scenario of Gabor filters. Specifically, Gabor filters are mathematically defined for infinite length signals [37]; the digital image signals, on the other hand, are finite in both spatial and frequency domains due to the limited length of image width and sampling frequency. As a result, directly applying Gabor filters to finite-length images may cause mathematical deficiency and distortion, leading to frequency aliasing and other issues that can impair model training [39]. To address this problem, we propose to concentrate most of the Gabor filter energy within the finite signal zone. This ensures that only a small amount of filter energy spills over the finite signal, minimizing the negative effects of using infinite-length-defined filters on finite-length images. Such a constraint could be implemented by limiting the value ranges of filter parameters, which determine the effective zone and energy concentration area of a Gabor filter. Specifically, according to the properties of Gabor filters and images mentioned above, we first derive $[0, \pi]$, $[\frac{5}{2\pi(1-2W)}, \frac{S}{5}]$, $[\frac{5}{2\pi}, \frac{S}{5}]$ and $[0, \frac{2\pi S-25}{4\pi S}]$ to be the valid ranges for parameters θ , σ_x , σ_y and W respectively (please refer to **supplementary materials** for a detailed derivation of these valid ranges). Next, we propose a constraint function to ensure each parameter to fall within its valid range. Specifically for a parameter p with lower bound l_p and upper bound u_p , we propose to obtain p through a learnable parameter \mathbf{p} as follows:

$$p = l_p + (u_p - l_p) \text{Sigmoid}(\mathbf{p}) = \frac{u_p + l_p e^{-\mathbf{p}}}{1 + e^{-\mathbf{p}}}. \quad (3)$$

The parameter \mathbf{p} can be learned to have any value ranging from $-\infty$ to $+\infty$. By applying Eq. 3, parameter p can be guaranteed to fall within its valid range $[l_p, u_p]$. By doing so, parameter values are constrained to enhance the stability of training.

High Frequency Enhancement. To address the frequency bias problem and improve high-frequency information extraction, we further propose a dual-constraint strategy for the frequency parameter W . Specifically, we constrain half of Gabor filters to only exploit high-frequency components, making the network to capture sufficient high-frequency information for facilitating fine-grained recognition. This is enabled by setting two value intervals for W : $[0, \frac{2\pi S-25}{8\pi S}]$ and $[\frac{2\pi S-25}{8\pi S}, \frac{2\pi S-25}{4\pi S}]$, which are equally divided from the valid range $[0, \frac{2\pi S-25}{4\pi S}]$ of the frequency parameter. Let N be the total number of Gabor filters to be learned. Using the method in Eq. 3, we constrain W of $N/2$ fil-

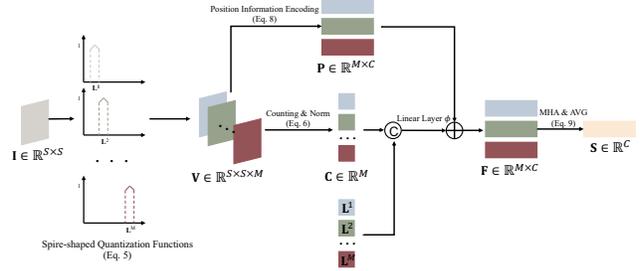


Figure 2. The structure of Learnable Histogram Operator (LHO).

ters to fall between $[0, \frac{2\pi S-25}{8\pi S}]$ and the other $N/2$ filters to fall between $[\frac{2\pi S-25}{8\pi S}, \frac{2\pi S-25}{4\pi S}]$, such that they serve as the low-frequency expert and high-frequency expert respectively. This strategy enhances the extraction ability of high-frequency information, alleviating the frequency bias problem in traditional DNNs where high-frequency features are difficult to extract.

3.4. Statistical Feature Extractor

The above method enables Gabor filters to be trained stably and effectively in a deep learning framework. These filters can capture texture signals from an input image. In the next step, we target at extracting effective features from these filters to facilitate fine-grained recognition. From the perspective of digital image processing, texture can be represented by the statistical properties of feature maps [16]. Based on this principle, in some previous research [40, 4, 1], textures were exploited by extracting hand-crafted statistics such as variance, contrast and smoothness from the output of Gabor filters. However, these statistics have limited effectiveness because they are hand-crafted designed. In this work, we step further and take advantage of deep learning to extract deep statistical features that are more comprehensive and effective. Specifically, we propose a novel statistical feature extractor that contains a **Learnable Histogram Operator (LHO)** followed by a **Filer Correlation Module (FCM)**. In the following sections, we provide a detailed illustration of QCO and FCM respectively.

Learnable Histogram Operator. We first employ a Learnable Histogram Operator (LHO) to exploit statistical information, which is improved from [58]. Fig. 2 illustrates the structure of LHO. LHO receives the output from each Gabor filter as input and produces a statistical feature $\mathbf{S} \in \mathbb{R}^C$, where C refers to the channel number. Let $\mathbf{I} \in \mathbb{R}^{S \times S}$ be the intensity map generated by convolving a Gabor filter with the processed image, where S refers to image size. LHO first defines M levels by equally dividing the minimal and maximal values for intensities of all pixels. Specifically, the m -th level \mathbf{L}^m is computed as:

$$\mathbf{L}^m = \min(\mathbf{I}) + \frac{\max(\mathbf{I}) - \min(\mathbf{I})}{M} \cdot m \quad (4)$$

Next, each pixel $\mathbf{I}^{i,j}$ on \mathbf{I} is quantized to a vector $\mathbf{V}^{i,j} \in$

\mathbb{R}^M . Concretely, let I be the interval between adjacent levels ($I = \mathbf{L}^{m+1} - \mathbf{L}^m$), the m -th channel $\mathbf{V}^{i,j,m}$ on $\mathbf{V}^{i,j}$ is computed as follows:

$$\mathbf{V}^{i,j,m} = \begin{cases} 1 - |\mathbf{L}^m - \mathbf{I}^{i,j}| & \text{if } |\mathbf{L}^m - \mathbf{I}^{i,j}| < \frac{I}{2} \\ 0 & \text{else} \end{cases} \quad (5)$$

In this way, \mathbf{I} is quantized to $\mathbf{V} \in \mathbb{R}^{S \times S \times M}$. Note a spire-shaped function is used for quantization instead of a binary one to ensure that the operation is differentiable. After the quantization, the number of pixels that fall into each level is counted and normalized to obtain a counting feature $\mathbf{C} \in \mathbb{R}^M$. Specifically, \mathbf{C}^m for the m -th level is computed as:

$$\mathbf{C}^m = \frac{\sum_{i=1}^S \sum_{j=1}^S \mathbf{V}^{i,j,m}}{\sum_{m=1}^M \sum_{i=1}^S \sum_{j=1}^S \mathbf{V}^{i,j,m}}. \quad (6)$$

We concatenate \mathbf{C}^m with \mathbf{L}^m , and obtain the feature \mathbf{F}^m for the m -th quantization level as follows:

$$\mathbf{F}^m = \phi(\text{CAT}(\mathbf{C}^m, \mathbf{L}^m)) + \mathbf{P}^m, \quad (7)$$

where ϕ refers to a linear layer. $\mathbf{P}^m \in \mathbb{R}^C$ encodes the position information for pixels quantized to the m -th level, which is computed as the weighted sum of position features from all pixels as follows:

$$\mathbf{P}^m = \sum_{i=1}^S \sum_{j=1}^S \mathbf{V}^{i,j,m} \text{PE}(i, j), \quad (8)$$

where PE refers to the position encoding operation as in vit [8]. Finally, the statistical texture \mathbf{S} is produced by correlating features from different levels:

$$\mathbf{S} = \text{Avg}(\text{MHA}(\{\mathbf{F}^m\}_{m=1}^M)), \quad (9)$$

where MHA denotes a multi-head attention and Avg refers to the average over all levels. Using LHO has two advantages. Firstly, it can be considered as a way to extract deep features from a differentiable and learnable histogram, which can produce effective and appropriate statistical features that can benefit the task [16, 58]. Secondly, conventional statistical features are typically pixel-relationship-invariant. For example, rearranging pixel positions does not alter the mean and variance values. In contrast, our method is pixel-relationship-aware as it encodes \mathbf{P}^m that represents the position information of pixels quantized to each level, which provides additional information that can be useful for fine-grained recognition. These advantages enable the extracted features from LHO to be effective in facilitating fine-grained recognition.

Filter Correlation Module. The proposed LHO can produce a statistical feature from the output of each Gabor filter, resulting in a set of features $\{\mathbf{S}^n\}_{n=1}^N$ for all filters. We

further propose a Filter Correlation Module (FCM) to aggregate statistical information from all Gabor filters as follows:

$$\mathbf{T} = \text{Avg}\left(\text{MHA}\left(\{\mathbf{S}^n + \phi(\sigma_x^n, \sigma_y^n, \theta^n, W^n)\}_{n=1}^N\right)\right), \quad (10)$$

in which \mathbf{S}^n refers to the statistical feature for the n -th Gabor filter. $\phi(\sigma_x^n, \sigma_y^n, \theta^n, W^n)$ uses a linear layer to encode filter parameter information, which can indicate the type of texture extracted by the Gabor filter. Eq. 10 exploits the correlations among different texture components from different Gabor filters, generating a powerful feature with long-range correlations that can facilitate fine-grained recognition.

3.5. Region Selection Gate

Through the methods outlined above, we can build a powerful feature extractor that captures texture information using Gabor filters. As a common practice used by previous methods, the extractor can generate features in an image-based [40] manner or a patch-based manner [39]. Image-based methods directly feed the whole image into the extractor, which might lead to the blurring of local detailed information. Patch-based methods address the problem by equally dividing the image into several patches and extracting a feature from each patch. However, it is computationally expensive as each image position must be processed by the feature extractor. Intuitively, the texture information in different regions of an image is not equally important. For example, some regions may contain critical textures for recognition, such as feather of the bird, while others may be less informative, consisting of background textures such as sky, road and lawn. Therefore, extracting textures from all positions is redundant and computationally inefficient, which can impede recognition accuracy and slow down algorithmic speed. To address the issue, we propose a novel approach to solely extract textures from informative regions that are automatically selected from a gate mechanism. More specifically, the gate mechanism for region selection involves three steps. Firstly, we introduce a shallow FPN block that takes the intermediate features of the CNN-based semantic branch as inputs and generates a feature map F . Please refer to supplementary materials for a detailed description for the structure of this block. Since features from the CNN branch have a high receptive field and contain rich semantic information, F can identify key parts of an image and assist in selecting class-discriminative regions in appropriate positions. Secondly, we introduce a set of regions proposals, denoted as $\{\mathcal{R}^k\}_{k=1}^K$, which is the same as that in RCNN [15], comprising regions with different scales and positions. For each region \mathcal{R}^k , a feature vector f^k is generated by performing an ROI pooling on the corresponding region of the feature map F . This feature vector is then fed into a linear layer to generate a score s^k .

Lastly, regions that require texture extraction are selected based on their scores. This mechanism enables only the regions that contain class-discriminative information to be selected for texture extraction, thereby reducing computation costs and avoiding information redundancy.

To reduce computation for capturing texture features, we expect ‘hard’ selections instead of the ‘soft’ attention scores to be generated so that only a limited number of regions need texture extraction. However, such a hard selection may prohibit the gradients to be back-propagated to the earlier layers in optimization. To address this issue, we employ the Improved Semantic Hashing [24, 25] for region selection. In training, to encourage more decision space to be explored randomly, we add a standard Gaussian noise ϵ to s^k , resulting in the noise-injected score \hat{s}^k . Then two vectors are generated as follows:

$$c^k = \sigma'(\hat{s}^k); d^k = \mathbb{1}(\hat{s}^k > 0), \quad (11)$$

where σ' refers to the saturating Sigmoid function:

$$\sigma'(\hat{s}^k) = \max(0, \min(1, 1.2\sigma(\hat{s}^k) - 0.1)), \quad (12)$$

where σ denotes the original Sigmoid function. In this way, d^k is a binarized discrete feature with value $\{0, 1\}$, which is used for making the selection decision. c^k is a continuously differentiable vector that can be used to approximate gradients in back-propagation, whose details are presented in supplementary materials.

Discussion with Previous Methods. Note that the proposed region gate mechanism differs from other key part localization methods used for fine-grained recognition. Unlike existing methods, our method can be trained end-to-end and generate hard decisions without the need for extra annotations. These properties enable our method to be highly effective and practical for region selection. Conversely, existing key part localization techniques have several limitations. For example, some methods [20, 43] require extra annotations such as bounding box and key points for joint training; [14] uses an indifferentiable selection process, making the network unable to be trained end-to-end; other methods [7, 21] typically produce soft attention rather than hard decisions, which cannot reduce computation costs. Our approach can be more practical and effective benefiting from our careful designs.

3.6. Optimization

The proposed network is optimized by the following loss function:

$$\mathcal{L} = L_{ce}(p, l) + \lambda \sum_{k=1}^K d^k, \quad (13)$$

where L_{ce} denotes the cross-entropy loss. p and l are the prediction logits and ground truth label. d^k denotes the binarized discrete decision as in Eq. 11. The second item

Method	Backbone	Accuracy
ResNet50 [17]	ResNet50	84.5
ResNet101 [17]	ResNet101	85.5
DenseNet161 [19]	DenseNet161	85.5
PC-Dense161 [9]	DenseNet161	86.9
NTSNet [48]	ResNet50	87.5
Cross-X [34]	ResNet50	87.7
DCL [3]	ResNet50	87.8
S3N [7]	ResNet50	88.5
iSQRT-COV [27]	ResNet101	88.7
GaRD [54]	ResNet50	89.6
APINet [59]	DenseNet161	90.0
Ours	ResNet50	90.8
Ours	ResNet101	91.3
Ours	DenseNet161	91.5

Table 1. Comparison results on CUB-200-2011 dataset.

in Eq. 13 encourages fewer regions to be selected, thus avoiding the redundant information and reducing computation costs. λ is a hyper-parameter to control the trade-off between the two loss items.

4. Experiments

4.1. Implementation Details

Following [59], we adopt different networks as the semantic branch, including ResNet50, ResNet101 and DenseNet161. We use SGD as the optimizer with the momentum of 0.9. The initial learning rate is set to 1e-4, which decays by 0.1 for every 20 epochs (overall 100 epochs). Batch size is 16 for all datasets. Following previous works [54, 28], we use random cropping of 448×448 in training and center crop at inference. We apply the commonly-used data augmentation strategies for training, including random scaling, left-right flipping and random cropping. S indicating the size of regions that are zoomed in is set to 112. N indicating the number of learnable Gabor filters is set to 128. M indicating the number of quantization levels in QCO is set to 8. λ in Eq. 13 is set to 0.2. Experiments are conducted on 2 NVIDIA TITAN V100 GPUs.

4.2. Comparison with State-of-the-art

We evaluate our method on four datasets, including three fine-grained datasets CUB-200-2011, NA-bird, Stanford Dogs and a terrain recognition dataset GTOS-mobile. The comparison results on these datasets are presented in Tables 1, 2, 3, and 4, respectively. Our method achieves the best performance when using DenseNet161 as the semantic branch, reaching accuracy scores of 91.5, 90.4, 92.1 and 87.5, outperforming the second-best method by 1.5, 2.3, 1.8 and 4.8, respectively. It is worth noting that even when using ResNet50, the least effective backbone among the three, our method remains superior to all other compared methods. This highlights the high effectiveness of our method. The advantage of our method is especially significant on the challenging terrain dataset GTOS-mobile. This is due to our

Method	Backbone	Accuracy
ResNet50 [17]	ResNet50	82.2
ResNet101 [17]	ResNet101	82.9
DenseNet161 [19]	DenseNet161	83.1
Cross-X [34]	ResNet50	86.4
GaRD [54]	ResNet50	88.0
APINet [59]	DenseNet101	88.1
Ours	ResNet50	89.5
Ours	ResNet101	90.0
Ours	DenseNet161	90.4

Table 2. Comparison results on NA-bird dataset.

Method	Backbone	Accuracy
ResNet50 [17]	ResNet50	80.1
ResNet101 [17]	ResNet101	80.6
DenseNet161 [19]	DenseNet161	80.9
PC-Dense161 [9]	DenseNet161	82.8
Cross-X [34]	ResNet50	88.9
APINet [59]	ResNet101	90.3
Ours	ResNet50	91.0
Ours	ResNet101	91.7
Ours	DenseNet161	92.1

Table 3. Comparison results on Stanford Dogs dataset.

Method	Backbone	Accuracy
ResNet50 [17]	ResNet50	69.4
ResNet101 [17]	ResNet101	72.0
DenseNet161 [19]	DenseNet161	72.9
B-CNN [28]	ResNet50	75.8
Deep-TEN [51]	ResNet50	76.1
DEP [47]	ResNet50	82.7
Ours	ResNet50	86.4
Ours	ResNet101	87.0
Ours	DenseNet161	87.5

Table 4. Comparison results on GTOS-mobile dataset.

Method	Accuracy	Flops(G)	Params(M)
Baseline (ResNet50)	84.5	16.26	24.69
+ TB & RSG	90.8	20.72	32.41
+ TB (image-based manner)	89.0	36.88	29.82
+ TB (patch-based manner)	89.9	37.04	29.82

Table 5. Ablation results of different components in our method. TB and RSG refer to the texture branch and region selection gate, respectively.

method’s ability to extract texture features, which are essential for the classification of visually-similar and closely-related terrain classes.

4.3. Ablation Study

We perform ablation study to verify the effectiveness of our designs. Experiments in this part are conducted on CUB-200-2011 dataset with ResNet50 as the semantic branch. Due to paper length limitation, more ablation study results are presented in **supplementary materials**.

Ablation of Different Components. As shown in Fig. 1, in addition to the basic CNN-based semantic branch, our method consists of two key components: a Gabor-filter-based texture branch (TB) and a region selection gate (RSG). We conduct experiments to validate the effect of these components and present the results in Table. 5, which includes the validation accuracy, computation costs and parameter numbers of different settings. We set the network that only has a ResNet50 semantic branch to be the baseline. Our method with adding both texture branch & region selection gate increases accuracy from 84.5 to 90.8. When the region selection gate is removed, extracting textures by feeding the whole image into texture branch (image-based manner) reduces accuracy to 89.0 and increases computation by 78%; equally dividing image into 8×8 patches and getting a feature from each patch (patch-based manner) reduces accuracy to 89.0 and increases computation by 79%. The results demonstrate that the region selection gate can reduce computation and improve performance.

Ablation of Learnable Gabor Filters. Results in this part are presented in Table. 6. (1) First, we validate the effectiveness of Gabor filters. Replacing all Gabor filters by the 3×3 convolution layers decreases accuracy by 4.9. This indicates that the improvement does not simply come from

Method	Accuracy	Flops(G)	Params(M)
Ours	90.8	20.72	32.41
Ours replacing Gabor filters with conv	85.9	20.03	32.95
Ours w/o parameter learning	86.5	20.72	32.41
Ours w/o parameter value constraints	75.4	20.40	32.41
Ours w/o high frequency enhancement	88.1	20.40	32.41

Table 6. Ablation results of learnable Gabor filters.

an increase in the number of parameters, but from the effectiveness of using Gabor filters for texture extraction. (2) Second, we evaluate the effectiveness of learning Gabor filter parameters in a supervised manner. We replace the learnable parameters with the manually set ones, where the setting rules are the same as [40]. This reduces accuracy by 4.3, showing the effectiveness of supervised learning. (3) Third, we verify the necessity of using parameter value constraints. In Fig. 3, we present the loss curves and validation accuracy curves for methods w/ and w/o using value constraints. Using value constraints results in a smoother loss decline curve and accuracy improvement curve. In contrast, free training causes violent loss and accuracy fluctuations, especially at the beginning stage of training, indicating that the training process is unstable. Compared to free training, using value constraints significantly improves validation accuracy by 15.4. (4) Last, we evaluate the high frequency enhancement strategy by setting two constraints for frequency parameter W . Without the strategy and directly constraining W into its valid range reduces accuracy by 2.7. Results demonstrate the effectiveness of our designs for the learnable Gabor filter parameters. Also note, using the proposed learning methods requires no extra parameters and only increases computation very slightly.

Ablation of Statistical Feature Extractor. We propose a statistical feature extractor to leverage texture information from the output of Gabor filters. The extractor consists of two parts: a Learnable Histogram Operator (LHO) and a Filter Correlation Module (FCM). We conduct experiments to validate the effect of these components and present the results in Table. 7. Removing LHO and FCM reduces accuracy from 90.8 to 87.0 and 89.2 respectively. Also note that, both LHO and FCM are lightweight, only requiring extra parameters of 4.26M and 3.56M, increasing flops by

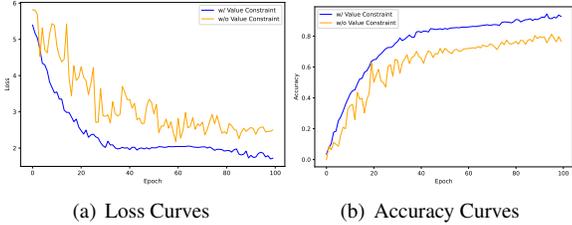


Figure 3. Loss curves (a) and accuracy curves (b) in different settings. Blue and orange lines indicate curves for training Gabor filters w/ and w/o using parameter value constraints, respectively.

Method	Accuracy	Flops(G)	Params(M)
Ours	90.8	20.72	32.41
Ours w/o LHO	87.0	18.03	28.25
Ours w/o FCM	89.2	19.55	28.85

Table 7. Ablation results of statistical feature extractor.

15% and 6% respectively. Results show the high effectiveness and efficiency of LHO and FCM.

Ablation of LHO. In Table. 8, we further evaluate the designs in LHO. First, we replace LHO with a hand-crafted statistical feature composed of mean, variance, maximum and minimum values, which reduces accuracy by 2.9. This demonstrates that LHO can produce statistical features that are more effective than hand-crafted ones widely used by previous methods. Next, we remove position information \mathbf{P}^m in Eq. 7 and multi-head attention in Eq. 9, which reduces accuracy by 1.5 and 2.1 respectively, showing the importance of encoding position features and correlating different levels in LHO. Finally, we move LHO from texture branch to CNN-based semantic branch. This reduces accuracy by 4.6, showing that the effectiveness of LHO relies on Gabor filters to extract texture-related signals.

4.4. Visualization

In Fig. C, we provide a visualization of the outputs obtained from learned Gabor filters. More specifically, Fig. C (c) and Fig. C (d) show the average output of all high-frequency and low-frequency Gabor filters, respectively. As can be observed, the high-frequency filters primarily capture information of undulating areas such as speckles and ripples, whereas the low-frequency filters primarily capture information related to smooth changing areas. Both these two kinds of information are critical for recognition. By exploiting sufficient and balanced multi-frequency features through the carefully-designed learnable Gabor filters, our method can leverage comprehensive information for effective fine-grained recognition. Due to space limitation, more visualization results are presented in **supplementary**.

4.5. Discussion of Computation and Parameter

Compared to the ResNet50 baseline, our method only increases computation and parameter usage slightly. Specifically, on CUB-200-2011 dataset, ResNet50 consumes

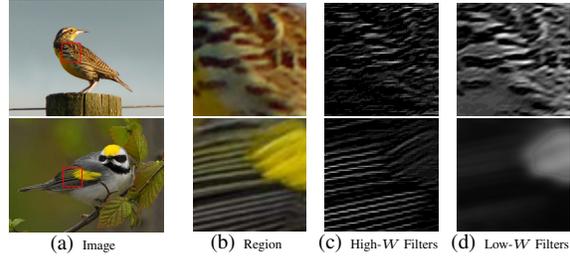


Figure 4. Visualization of Gabor filters. (a), (b), (c) and (d) present the original images, the selected regions, average output of all high-frequency and low-frequency Gabor filters, respectively.

Method	Accuracy	Flops(G)	Params(M)
Ours	90.8	20.72	32.41
Ours w/o LHO w/ HSF	87.9	18.16	28.25
Ours w/o \mathbf{P}^m in Eq. 7	89.3	18.50	30.93
Ours w/o MHA in Eq.9	88.7	19.49	29.73
Ours w/o LHO in TB w/ LHO in CNN	86.2	22.85	32.90

Table 8. Ablation results of LHO. HSF refers to a hand-crafted statistical feature composed of mean, variance, maximum and minimum values. MHA denotes multi-head attention. TB denotes texture branch.

16.26G flops of computation, requiring 24.69M parameters, reaching accuracy of 84.5. Our method consumes 20.72G flops of computation, requiring 32.41M parameters, reaching accuracy of 90.8. Compared to ResNet50, our method can improve accuracy significantly (+6.3) with only increasing computation and parameter usage by 27% and 31% respectively. This demonstrates the high effectiveness and high efficiency of our method with outstanding performance and reasonable computation costs.

5. Conclusion

This paper presents a novel network that extracts texture features to facilitate fine-grained recognition using Gabor filters. We make the Gabor filter parameters to be automatically learned in a stable and effective manner through our carefully-designed value constraint strategies. We also propose a statistical feature extractor and a region selection gate, allowing informative and effective Gabor texture features to be extracted with very few extra computation costs. Our method is effective, as shown by the extensive experiments with state-of-the-art (SOTA) performance.

Acknowledgement This research is supported by the National Research Foundation, Singapore under its AI Singapore Programme (AISG Award No: AISG2-PhD-2021-08-006), Singapore Ministry of Education (MOE) AcRF Tier 2 under Grant MOE-T2EP20222-0009, the National Research Foundation Singapore through AI Singapore Programme under Grant AISG-100E-2020-065, and SUTD SKI Project under Grant SKI 2021_02_06.

References

- [1] Jon Arrospe and Luis Salgado. Log-gabor filters for image-based vehicle verification. *IEEE Transactions on Image Processing*, 22(6):2286–2295, 2013. 1, 2, 4
- [2] Steve Branson, Grant Van Horn, Serge Belongie, and Pietro Perona. Bird species categorization using pose normalized deep convolutional nets. *arXiv preprint arXiv:1406.2952*, 2014. 2
- [3] Yue Chen, Yalong Bai, Wei Zhang, and Tao Mei. Destruction and construction learning for fine-grained image recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5157–5166, 2019. 1, 2, 6
- [4] David A Clausi and Huang Deng. Design-based texture feature fusion using gabor filters and co-occurrence probabilities. *IEEE transactions on image processing*, 14(7):925–936, 2005. 4
- [5] John G Daugman. Two-dimensional spectral analysis of cortical receptive field profiles. *Vision research*, 20(10):847–856, 1980. 1
- [6] John G Daugman. Uncertainty relation for resolution in space, spatial frequency, and orientation optimized by two-dimensional visual cortical filters. *JOSA A*, 2(7):1160–1169, 1985. 1
- [7] Yao Ding, Yanzhao Zhou, Yi Zhu, Qixiang Ye, and Jianbin Jiao. Selective sparse sampling for fine-grained image recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6599–6608, 2019. 2, 6
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 1, 5
- [9] Abhimanyu Dubey, Otkrist Gupta, Pei Guo, Ramesh Raskar, Ryan Farrell, and Nikhil Naik. Pairwise confusion for fine-grained visual classification. In *Proceedings of the European conference on computer vision (ECCV)*, pages 70–86, 2018. 2, 6, 7
- [10] Abhimanyu Dubey, Otkrist Gupta, Ramesh Raskar, and Nikhil Naik. Maximum-entropy fine grained classification. *Advances in neural information processing systems*, 31, 2018. 1, 2
- [11] Jianlong Fu, Heliang Zheng, and Tao Mei. Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4438–4446, 2017. 2
- [12] Xiao Fu, Shangzhan Zhang, Tianrun Chen, Yichong Lu, Lanyun Zhu, Xiaowei Zhou, Andreas Geiger, and Yiyi Liao. Panoptic nerf: 3d-to-2d label transfer for panoptic urban scene segmentation. In *2022 International Conference on 3D Vision (3DV)*, pages 1–11. IEEE, 2022. 1
- [13] Yang Gao, Oscar Beijbom, Ning Zhang, and Trevor Darrell. Compact bilinear pooling. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 317–326, 2016. 1, 2
- [14] Weifeng Ge, Xiangru Lin, and Yizhou Yu. Weakly supervised complementary parts models for fine-grained image classification from the bottom up. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3034–3043, 2019. 6
- [15] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014. 3, 5
- [16] Rafael C Gonzales and Paul Wintz. *Digital image processing*. Addison-Wesley Longman Publishing Co., Inc., 1987. 4, 5
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 6, 7
- [18] Xiao-dong Hu, Xin-qing Wang, Fan-jie Meng, Xia Hua, Yujia Yan, Yu-yang Li, Jing Huang, and Xun-lin Jiang. Gabor-cnn for object detection based on small samples. *Defence Technology*, 16(6):1116–1129, 2020. 2
- [19] 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. 6, 7
- [20] Shaoli Huang, Zhe Xu, Dacheng Tao, and Ya Zhang. Part-stacked cnn for fine-grained visual categorization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1173–1182, 2016. 2, 6
- [21] Zixuan Huang and Yin Li. Interpretable and accurate fine-grained recognition via region grouping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8662–8672, 2020. 2, 6
- [22] Anil K Jain, Nalini K Ratha, and Sridhar Lakshmanan. Object detection using gabor filters. *Pattern recognition*, 30(2):295–309, 1997. 2
- [23] Ruyi Ji, Longyin Wen, Libo Zhang, Dawei Du, Yanjun Wu, Chen Zhao, Xianglong Liu, and Feiyue Huang. Attention convolutional binary neural tree for fine-grained visual categorization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10468–10477, 2020. 1, 2
- [24] Łukasz Kaiser and Samy Bengio. Discrete autoencoders for sequence models. *arXiv preprint arXiv:1801.09797*, 2018. 6
- [25] Łukasz Kaiser, Samy Bengio, Aurko Roy, Ashish Vaswani, Niki Parmar, Jakob Uszkoreit, and Noam Shazeer. Fast decoding in sequence models using discrete latent variables. In *International Conference on Machine Learning*, pages 2390–2399. PMLR, 2018. 6
- [26] Shu Kong and Charless Fowlkes. Low-rank bilinear pooling for fine-grained classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 365–374, 2017. 1, 2
- [27] Peihua Li, Jiangtao Xie, Qilong Wang, and Zilin Gao. Towards faster training of global covariance pooling networks by iterative matrix square root normalization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 947–955, 2018. 6

- [28] Tsung-Yu Lin, Aruni RoyChowdhury, and Subhransu Maji. Bilinear cnn models for fine-grained visual recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 1449–1457, 2015. 6, 7
- [29] Chengjun Liu. Gabor-based kernel pca with fractional power polynomial models for face recognition. *IEEE transactions on pattern analysis and machine intelligence*, 26(5):572–581, 2004. 1, 2
- [30] Chengjun Liu and Harry Wechsler. Independent component analysis of gabor features for face recognition. *IEEE transactions on Neural Networks*, 14(4):919–928, 2003. 1, 2, 3
- [31] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 1
- [32] Shangzhen Luan, Chen Chen, Baochang Zhang, Jungong Han, and Jianzhuang Liu. Gabor convolutional networks. *IEEE Transactions on Image Processing*, 27(9):4357–4366, 2018. 2
- [33] Tao Luo, Zheng Ma, Zhi-Qin John Xu, and Yaoyu Zhang. Theory of the frequency principle for general deep neural networks. *arXiv preprint arXiv:1906.09235*, 2019. 1, 3
- [34] Wei Luo, Xitong Yang, Xianjie Mo, Yuheng Lu, Larry S Davis, Jun Li, Jian Yang, and Ser-Nam Lim. Cross-x learning for fine-grained visual categorization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8242–8251, 2019. 6, 7
- [35] Beom-Seok Oh, Kar-Ann Toh, Andrew Beng Jin Teoh, and Zhiping Lin. An analytic gabor feedforward network for single-sample and pose-invariant face recognition. *IEEE Transactions on Image Processing*, 27(6):2791–2805, 2018. 1, 2
- [36] Marcel Simon and Erik Rodner. Neural activation constellations: Unsupervised part model discovery with convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 1143–1151, 2015. 2
- [37] Kenneth Steiglitz. *Digital Signal Processing Primer*. Courier Dover Publications, 2020. 4, 12
- [38] Ming Sun, Yuchen Yuan, Feng Zhou, and Errui Ding. Multi-attention multi-class constraint for fine-grained image recognition. In *Proceedings of the european conference on computer vision (ECCV)*, pages 805–821, 2018. 1, 2
- [39] Zehang Sun, George Bebis, and Ronald Miller. On-road vehicle detection using evolutionary gabor filter optimization. *IEEE Transactions on Intelligent Transportation Systems*, 6(2):125–137, 2005. 2, 4, 5
- [40] Dacheng Tao, Xuelong Li, Xindong Wu, and Stephen J Maybank. General tensor discriminant analysis and gabor features for gait recognition. *IEEE transactions on pattern analysis and machine intelligence*, 29(10):1700–1715, 2007. 1, 2, 3, 4, 5, 7
- [41] Zhihui Wang, Shijie Wang, Shuhui Yang, Haojie Li, Jianjun Li, and Zezhou Li. Weakly supervised fine-grained image classification via gaussian mixture model oriented discriminative learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9749–9758, 2020. 1, 2
- [42] David M Weber and David P Casasent. Quadratic gabor filters for object detection. *IEEE Transactions on Image Processing*, 10(2):218–230, 2001. 2, 3
- [43] Xiu-Shen Wei, Chen-Wei Xie, Jianxin Wu, and Chunhua Shen. Mask-cnn: Localizing parts and selecting descriptors for fine-grained bird species categorization. *Pattern Recognition*, 76:704–714, 2018. 2, 6
- [44] Tianjun Xiao, Yichong Xu, Kuiyuan Yang, Jiaying Zhang, Yuxin Peng, and Zheng Zhang. The application of two-level attention models in deep convolutional neural network for fine-grained image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 842–850, 2015. 2
- [45] Zhi-Qin John Xu, Yaoyu Zhang, Tao Luo, Yanyang Xiao, and Zheng Ma. Frequency principle: Fourier analysis sheds light on deep neural networks. *arXiv preprint arXiv:1901.06523*, 2019. 1, 3
- [46] Zhi-Qin John Xu, Yaoyu Zhang, and Yanyang Xiao. Training behavior of deep neural network in frequency domain. In *Neural Information Processing: 26th International Conference, ICONIP 2019, Sydney, NSW, Australia, December 12–15, 2019, Proceedings, Part I 26*, pages 264–274. Springer, 2019. 1
- [47] Jia Xue, Hang Zhang, and Kristin Dana. Deep texture manifold for ground terrain recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 558–567, 2018. 7
- [48] Ze Yang, Tiange Luo, Dong Wang, Zhiqiang Hu, Jun Gao, and Liwei Wang. Learning to navigate for fine-grained classification. In *Proceedings of the European conference on computer vision (ECCV)*, pages 420–435, 2018. 6
- [49] Chaojian Yu, Xinyi Zhao, Qi Zheng, Peng Zhang, and Xinge You. Hierarchical bilinear pooling for fine-grained visual recognition. In *Proceedings of the European conference on computer vision (ECCV)*, pages 574–589, 2018. 2
- [50] Ye Yuan, Li-Na Wang, Guoqiang Zhong, Wei Gao, Wencong Jiao, Junyu Dong, Biao Shen, Dongdong Xia, and Wei Xiang. Adaptive gabor convolutional networks. *Pattern Recognition*, 124:108495, 2022. 2
- [51] Hang Zhang, Jia Xue, and Kristin Dana. Deep ten: Texture encoding network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 708–717, 2017. 7
- [52] Shangzhan Zhang, Sida Peng, Tianrun Chen, Linzhan Mou, Haotong Lin, Kaicheng Yu, Yiyi Liao, and Xiaowei Zhou. Painting 3d nature in 2d: View synthesis of natural scenes from a single semantic mask. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8518–8528, 2023. 1
- [53] Xudong Zhao, Ran Tao, Wei Li, Wilfried Philips, and Wenzhi Liao. Fractional gabor convolutional network for multi-source remote sensing data classification. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–18, 2021. 2
- [54] Yifan Zhao, Ke Yan, Feiyue Huang, and Jia Li. Graph-based high-order relation discovery for fine-grained recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15079–15088, 2021. 1, 2, 6, 7

- [55] Heliang Zheng, Jianlong Fu, Zheng-Jun Zha, and Jiebo Luo. Learning deep bilinear transformation for fine-grained image representation. *Advances in Neural Information Processing Systems*, 32, 2019. [1](#)
- [56] Heliang Zheng, Jianlong Fu, Zheng-Jun Zha, and Jiebo Luo. Looking for the devil in the details: Learning trilinear attention sampling network for fine-grained image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5012–5021, 2019. [1](#), [2](#)
- [57] Lanyun Zhu, Tianrun Chen, Jianxiong Yin, Simon See, and Jun Liu. Continual semantic segmentation with automatic memory sample selection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3082–3092, 2023. [1](#)
- [58] Lanyun Zhu, Deyi Ji, Shiping Zhu, Weihao Gan, Wei Wu, and Junjie Yan. Learning statistical texture for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12537–12546, 2021. [4](#), [5](#)
- [59] Peiqin Zhuang, Yali Wang, and Yu Qiao. Learning attentive pairwise interaction for fine-grained classification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 13130–13137, 2020. [6](#), [7](#)

The supplementary materials are arranged as follows. In Sec. A, additional details are presented to provide a more comprehensive understanding of our method. In Sec. B, further experimental results are included to validate the effectiveness of our approach. In Sec. C and D, additional visualization results and analysis are provided to offer deeper insights into the functioning of our method.

A. More Details of Method

A.1. Derivation of Parameter Valid Ranges

In Sec 3.3 of the text, we propose a constraint to ensure the stability of model training. This constraint limits the values of Gabor filter parameters into their valid ranges, which can be derived based on the properties of Gabor filters and digital images in different domains. According to the periodicity of angles, the orientation parameter has a valid range of $[0, \pi]$. As for the other parameters, we provide a detailed derivation of their valid ranges in the subsequent sections.

Scale Parameter $[\sigma_x, \sigma_y]$. $[\sigma_x, \sigma_y]$ are the scale parameters that determine the filter effective size in both spatial and frequency domains. In spatial domain, a Gaussian function modulates the sinusoidal plane wave, which is defined in the infinite signal space to satisfy its mathematical properties [37]. However, an image is a finite length signal in the spatial domain, with its valid signal zone determined by the image width S . Specifically for an image with size $S \times S$, the valid zone for each axis can be expressed as $[-0.5S, 0.5S]$ with the image center as the coordinate origin. Previous research [37] indicates that directly applying infinite-length-defined filters to the finite-length image zone would cause mathematical deficiency, which could limit the effectiveness of Gabor filters due to waveform distortion. To alleviate the problem, we propose a solution to concentrate most of the Gabor filter energy within the finite signal zone. This ensures that only a small amount of filter energy spills over the finite signal, minimizing the negative effects of using infinite-length-defined filters on finite-length images. Specifically for a Gaussian with mean μ and variance σ , we constrain $[\mu - \alpha\sigma, \mu + \alpha\sigma]$ to fall in the valid signal zone of the image. α is a hyper-parameter to control the energy concentration degree. According to the experimental results shown in Table. 9, we choose 2.5 to be the optimal value for α . By doing so, 98.76% of the Gaussian energy can be subtended, and only 1.24% of filter energy spills out of the image signal, whose negative effect is negligible. Based on the above analysis, we derive the spatial-wise constraints for parameters $[\sigma_x, \sigma_y]$ as follows:

$$\begin{cases} [-0.25\sigma_x, 0.25\sigma_x] \subseteq [-0.5S, 0.5S] \\ [-0.25\sigma_y, 0.25\sigma_y] \subseteq [-0.5S, 0.5S] \end{cases} \quad (14)$$

α	Subtended Energy	Accuracy
1.0	68.27%	85.8
1.5	86.64%	89.0
2.0	95.45%	90.2
2.5	98.76%	90.8
3.0	99.73%	90.3
3.5	99.95%	90.0

Table 9. Ablation results of hyper-parameter α for constraining Gabor filter parameters. When the value of α is too small, the percentage of subtended energy is also low. This leads to a large amount of filter energy spilling out of the effective signal zone, which in turn negatively impacts training stability. Conversely, if α is too large, the filter parameters may be constrained to a small range, leading to a loss of information across certain frequencies. Experimental results show that the optimal value for α is 2.5. This choice strikes a balance between training stability and the availability of sufficient multi-frequency information.

We further analyze the frequency-wise constraints for $[\sigma_x, \sigma_y]$. We perform a Fourier transform on Eq. 1 of text and get the frequency-wise expression of Gabor filters as follows:

$$G(u, v) = \exp \left[-\frac{1}{2} \left((4\pi^2\sigma_x^2(u - W)^2) + 4\pi^2\sigma_y^2v^2 \right) \right], \quad (15)$$

As can be observed from Eq. 15, in frequency domain, the Gabor filter also contains a Gaussian with mean $\{W, 0\}$ and variance $\{\frac{1}{2\pi\sigma_x}, \frac{1}{2\pi\sigma_y}\}$ to control its effective size. The valid signal zone in frequency domain can be derived according to Nyquist sampling theorem, which indicates that for a given sample rate f_s , perfect reconstruction is guaranteed possible when the frequency $|W| < (f_s/2)$, otherwise signal aliasing would happen. In an image, the sample rate equals 1 pixel, so any frequency component larger than 0.5 is distorted thus being invalid. This means the valid signal zone in frequency domain is $[-0.5, 0.5]$. Following the constraints in spatial domain, we subtend 98.76% of the frequency-wise Gaussian energy into the valid signal zone to avoid distortion, getting the constraints in frequency domains as follows:

$$\begin{cases} [W - \frac{2.5}{2\pi\sigma_x}, W + \frac{2.5}{2\pi\sigma_x}] \subseteq [-0.5, 0.5] \\ [-\frac{2.5}{2\pi\sigma_y}, \frac{2.5}{2\pi\sigma_y}] \subseteq [-0.5, 0.5] \end{cases} \quad (16)$$

Solving Eq. 14 and Eq. 16, we get $[\frac{5}{2\pi(1-2W)}, \frac{S}{5}]$ and $[\frac{5}{2\pi}, \frac{S}{5}]$ to be the valid ranges for σ_x and σ_y respectively.

Frequency Parameter W . We further analyze the valid range for the frequency parameter W . Due to the symmetry of image frequencies, any W less than 0 is mirrored with its opposite number $-W$, so frequency components less than 0 are not considered and the lower bound for W is set to

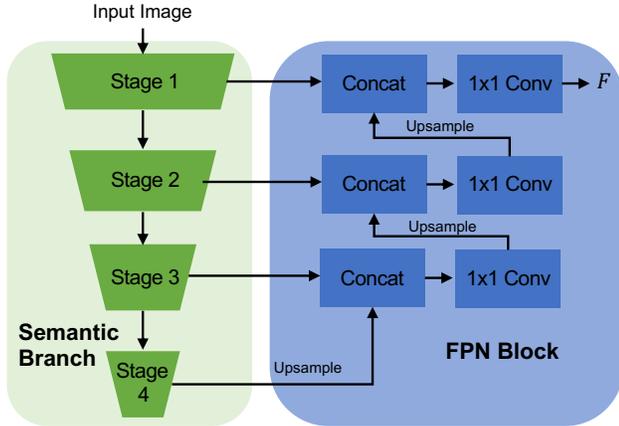


Figure 5. Structure of the FPN block in the proposed Region Selection Gate.

0. The upper bound can be derived from two constraints. First, according to Nyquist sampling theorem, frequency should be lower than $1/2$ to avoid aliasing. Second, the upper bound of σ_x should be higher than its lower bound. Mathematically, these constraints are formulated as follows:

$$\begin{cases} W < 0.5 \\ \frac{S}{5} > \frac{5}{2\pi(1-2W)} \end{cases} \quad (17)$$

Solving Eq.17, we get $\frac{2\pi S-25}{4\pi S}$ to be the upper bound and $[0, \frac{2\pi S-25}{4\pi S}]$ to be the valid range of W .

A.2. FPN Block in Region Selection Gate

The proposed region selection gate employs a FPN block to generate a feature F using the intermediate features of the CNN-based semantic branch, which is then used to assist in selecting informative regions for texture extraction (see Sec. 3.5 of the main paper for details). In Fig. 5, we show the detailed structure of the FPN block. The output channel number of all 1×1 convolution layers in the block is 128. This block integrates multi-level information from different layers. As a result, the generated F contains comprehensive information for the effective key part localization.

A.3. Back-Propagation of Improved Semantic Hashing.

In the proposed region selection gate, we employ the improved semantic hashing technique to make the selection operation differentiable. Specifically, for the k -th region proposal, a standard Gaussian noise is first added to its score s^k to produce \hat{s}^k . Then two vectors are generated from \hat{s}^k , including a binary discrete feature d^k and a continuously differentiable vector c^k (see Eq. 11 of the main paper for details). In forward-propagation, d^k is used to make region selection decisions. In back-propagation, we consider the

Method	Top-1 (%)
ResNet50	76.1
Ours (with ResNet50 backbone)	77.9

Table 10. Validation results on ImageNet.

gradient of c^k with respect to \hat{s}^k an approximation of the gradients for updating the parameters from the discrete gate d^k . This gradient replacement operation could be realized by $d^k = d^k + c^k - c^k.detach()$ in PyTorch. During inference, we skip the Gaussian noise sampling step and directly use the discrete output from its original score as the selection decision, i.e., $\mathbb{1}(s^k > 0)$.

B. More Experimental Results

B.1. Experiments on ImageNet

In addition to fine-grained recognition datasets, we also validate our method on ImageNet, which is a widely-used dataset for general image classification. The results are presented in Table. 10. As a baseline, ResNet50 achieves Top-1 accuracy of 76.1%. By using our method with ResNet50 as the semantic branch, we obtain Top-1 of 77.9%, which outperforms ResNet50 by 1.8%. Despite achieving higher accuracy, we observe that our method can bring greater improvement on fine-grained datasets than ImageNet. This can be explained by the different types of features required for different datasets. Specifically, the visual appearances and semantic meanings of different categories in ImageNet are significantly different, allowing us to classify different classes from a global perspective without explicitly exploiting local details. As a result, the high-level semantic information captured by CNN is already sufficient to distinguish different classes, while local detailed textures captured by our method can be less crucial. In contrast, fine-grained recognition datasets often include categories with very similar visual appearances and high-level semantic meanings (e.g., different bird species). These categories are very similar from the global view, only having subtle differences in some local areas. In this scenario, features from deep CNNs are insufficient for classification due to their lack of local detailed features and high-frequency information, as discussed in the Introduction section of our paper. Texture information extracted from our method can serve as an effective supplement to CNN features, significantly improving fine-grained recognition.

B.2. Ablation Study of Hyper-Parameters

In this section, we present ablation results of hyper-parameters used in our method, including the number of Gabor filters, the number of quantization levels in LHO, λ in the loss function, and the size that each region is zoomed into. Experiments in this section are conducted on CUB-200-2011 with ResNet50 as the semantic branch. We

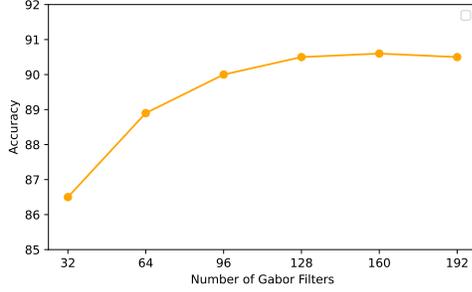


Figure 6. Ablation results of the Gabor filter number.

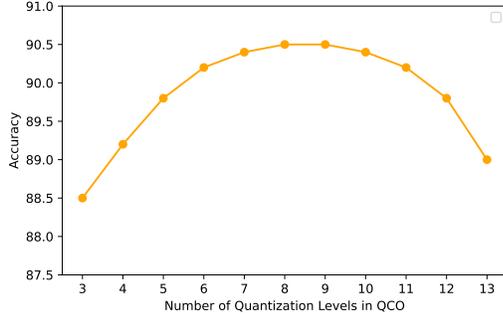


Figure 7. Ablation results for the number of quantization levels in LHO.

report the average results of 5 repeated experiments

Ablation of Gabor Filter Number The texture branch in our approach utilizes N learnable Gabor filters to process input regions. In Fig. 6, we present the validation results of our method using varying numbers of learnable Gabor filters. As shown in Fig. 6, increasing N from 32 to 128 results in an improvement in validation accuracy from 86.5 to 90.5. However, performance improvement becomes insignificant when N exceeds 128. Therefore, we choose $N = 128$ as the optimal number of Gabor filters.

Ablation for the Number of Quantization Levels. The proposed LHO involves a step that quantizes the intensity map into M levels in order to extract statistical information (refer to Eq. 4 and Eq. 5 of the main paper for further details). In Fig. 7, we present the validation results of using different numbers of quantization levels. It is observed that when M is greater than 7 and less than 11, the accuracy remains stable and near 90.5. Conversely, when M is too small, the quantization is coarse, resulting in less effective statistical feature extraction and lower validation accuracy. Furthermore, when M is too large, overfitting may occur to hinder the model’s effectiveness. Based on experimental results, we chose 8 as the setting for M . It is worth noting that our proposed method consistently outperforms the baseline ResNet50 significantly when M ranges from 3 to 13, thus demonstrating the high effectiveness of our approach.

λ	Accuracy	Flops(G)
0.01	90.0	25.57
0.05	90.2	23.26
0.1	90.7	21.05
0.2	90.8	20.72
0.3	90.5	20.65
0.5	90.2	20.46
1	90.0	20.20

Table 11. Ablation results of λ in Eq. 13 of the main paper.

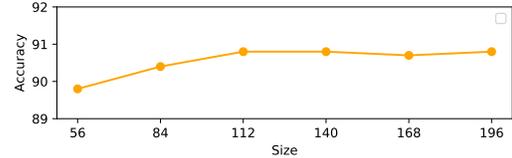


Figure 8. Ablation results for S indicating the size that each region is zoomed into.

Low-frequency Interval	High-frequency Interval	Accuracy
$[0, \frac{u_w}{5}]$	$[\frac{u_w}{5}, u_w]$	90.1
$[0, \frac{u_w}{4}]$	$[\frac{u_w}{4}, u_w]$	90.5
$[0, \frac{u_w}{3}]$	$[\frac{u_w}{3}, u_w]$	90.5
$[0, \frac{u_w}{2}]$	$[\frac{u_w}{2}, u_w]$	90.8
$[0, \frac{2u_w}{3}]$	$[\frac{2u_w}{3}, u_w]$	90.4
$[0, \frac{3u_w}{4}]$	$[\frac{3u_w}{4}, u_w]$	90.0
$[0, \frac{4u_w}{5}]$	$[\frac{4u_w}{5}, u_w]$	89.4

Table 12. Ablation results of different divisions for the low-frequency interval and high-frequency interval. $u_w = \frac{2\pi S - 25}{4\pi S}$ denotes the upper bound of frequency parameter W .

Ablation of λ in Loss Function. As shown in Eq. 13 of the main paper, we use a hyper-parameter λ to control the trade-off between the two loss items. In Table. 11, we present the validation accuracy and average flops when λ is set to different values.

The fluctuation of accuracy is less than 0.8 when λ varies from 0.01 to 1, showing that our method is non-sensitive to the choice of hyper-parameter λ .

Ablation of Size that Each Region is zoomed in. Each of the selected regions is zoomed into the size $S \times S$ before being feeding into the texture branch for feature extraction. In Fig. 8, we show the validation accuracy when setting S to different values. The accuracy keeps stable when S is greater than or equal to 112. Setting a smaller value for S results in a reduced input size for the texture branch and lower computational costs. Thus, we choose 112 to be the setting of S .

B.3. Ablation of High Frequency Enhancement Strategies.

To alleviate the frequency-bias problem and enhance the high-frequency texture extraction capability of Gabor filters, we propose a high frequency enhancement strategy

Low-frequency Filters	High-frequency Filters	Accuracy
$\frac{N}{5}$	$\frac{4N}{5}$	89.8
$\frac{N}{4}$	$\frac{3N}{4}$	90.3
$\frac{N}{3}$	$\frac{2N}{3}$	90.8
$\frac{N}{2}$	$\frac{N}{2}$	90.8
$\frac{2N}{3}$	$\frac{N}{3}$	90.7
$\frac{3N}{4}$	$\frac{N}{4}$	90.0
$\frac{4N}{5}$	$\frac{N}{5}$	89.4

Table 13. Ablation results of different amount allocations for the low-frequency filters and high-frequency filters.

by setting two value intervals for frequency parameter W : $[0, \frac{2\pi S - 25}{8\pi S}]$ and $[\frac{2\pi S - 25}{8\pi S}, \frac{2\pi S - 25}{4\pi S}]$, which are equally divided from the valid range $[0, \frac{2\pi S - 25}{4\pi S}]$ of W . We then constrain W of $N/2$ Gabor filters to fall between $[0, \frac{2\pi S - 25}{8\pi S}]$ and the other $N/2$ filters to fall between $[\frac{2\pi S - 25}{8\pi S}, \frac{2\pi S - 25}{4\pi S}]$, such that they serve as the low-frequency expert and high-frequency expert respectively. In Table. 12, we present the results obtained from applying different strategies with varying divisions for low-frequency and high-frequency intervals. We denote the upper bound of W as u_w , which equals $\frac{2\pi S - 25}{4\pi S}$. In Table. 13, we present the results of varying amounts of allocations for low-frequency and high-frequency filters. From the results in both tables, it can be observed that the performance remains stable when the division ranges from $\frac{u_w}{4}$ to $\frac{2u_w}{3}$ and the amount of low-frequency filters ranges from $\frac{N}{3}$ to $\frac{2N}{3}$. The results demonstrate that our proposed method is not sensitive to these specific settings.

C. More Visualization Results

C.1. Visualization of Selected Regions

In Fig. 9, we present some visualization results of the selected regions for texture extraction. These regions are marked by the red bounding boxes. The selected regions contain informative texture features that are difficult to be extracted by the vanilla CNNs. Using the proposed texture branch, we can extract effective texture features from these regions to facilitate fine-grained recognition.

C.2. Visualization for the Output of Gabor Filters

In Fig. 4 of the main paper, we have provided some visualization results of the outputs obtained from learned Gabor filters. In Fig. 12, we present more visualization results. More specifically, Fig. 12 (c) and Fig. 12 (d) show the average output of all high-frequency and low-frequency Gabor filters, respectively. As can be observed, the high-frequency filters primarily capture information of undulating areas such as speckles and ripples, whereas the low-frequency filters primarily capture information related to smooth chang-

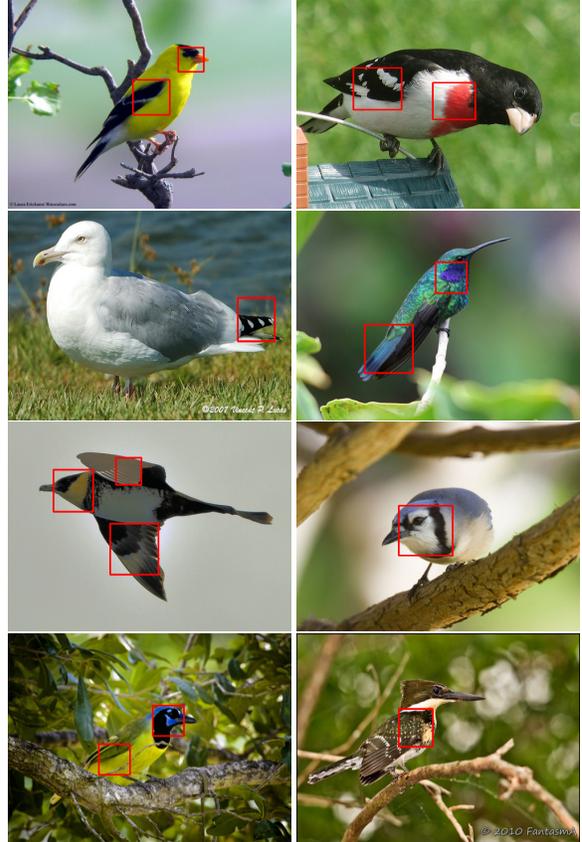


Figure 9. Visualization of selected regions for texture extraction. The selected regions are marked by the red bounding boxes.

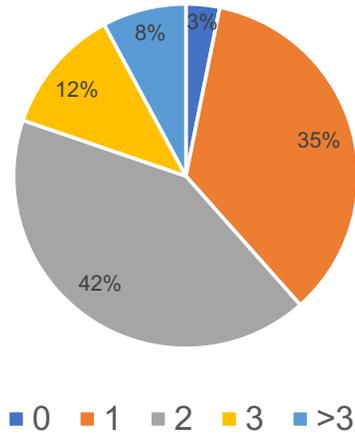


Figure 10. Percentage of images with different numbers of selected regions.

ing areas. Both kinds of information are critical for recognition. By exploiting sufficient and balanced multi-frequency features through the carefully-designed learnable Gabor filters, our method can leverage comprehensive information for effective fine-grained recognition.



Figure 11. Two image examples that have no region to be selected for texture extraction.

D. Statistical Analysis for the Number of Selected Regions.

Fig. 10 displays the percentage of images with various numbers of selected regions for texture extraction. The results indicate that, in general, only a few regions are selected for most images. This minimizes information redundancy and reduces computation costs. Specifically, 35% of all images have only one region selected, while 42% of all images have two regions selected for feature extraction. It is worth noting that a very small percentage of images have no regions selected for texture extraction. Fig. 11 illustrates two examples of such images. Typically, these images do not contain significant texture information that can facilitate recognition due to the low image quality or the category properties. Therefore, no region is selected for additional texture extraction.

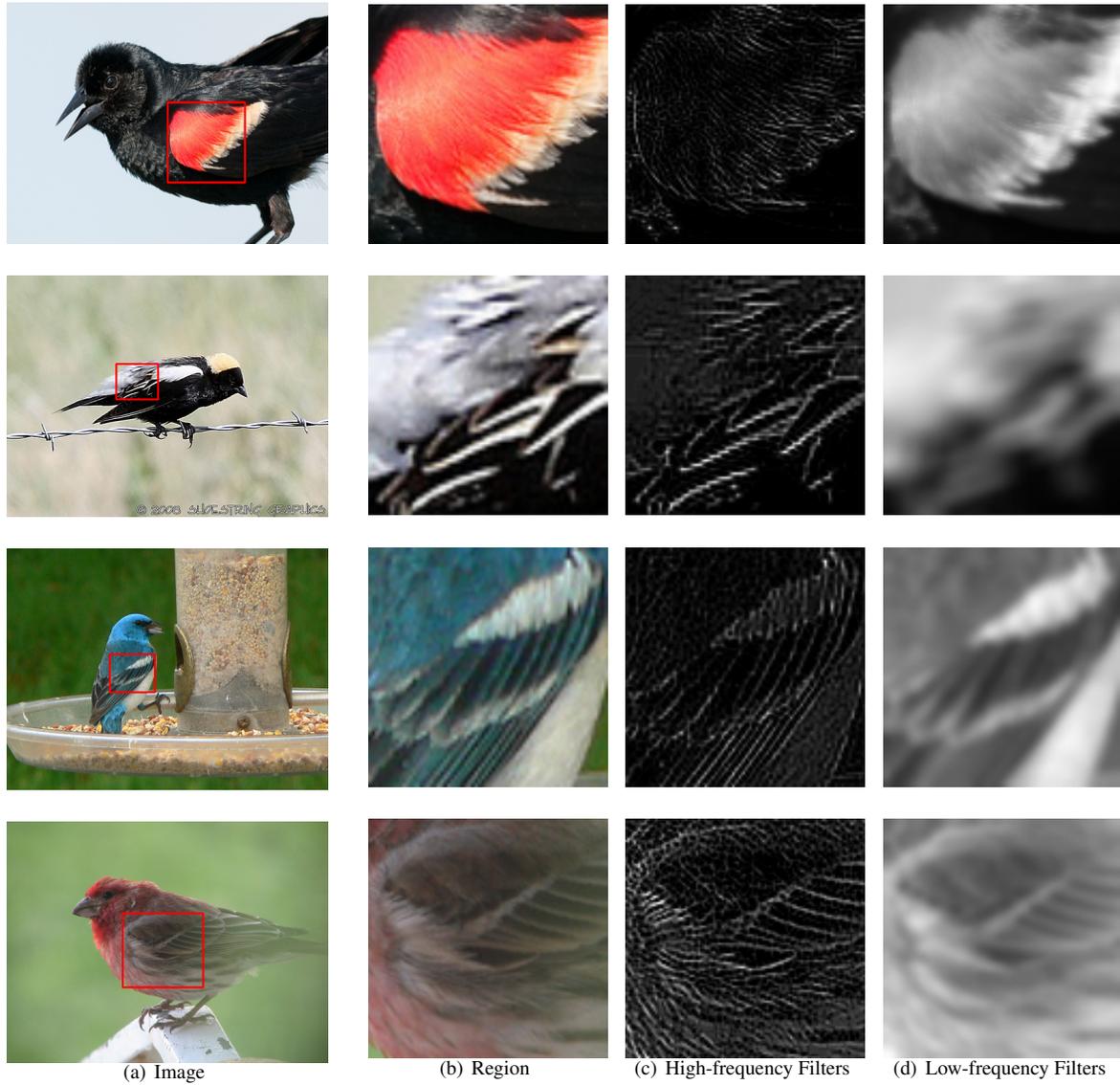


Figure 12. Visualization of output from Gabor filters. (a), (b), (c) and (d) present the original images, the selected regions, average output of all high-frequency and low-frequency Gabor filters, respectively.