
Global Filter Networks for Image Classification

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent advances in self-attention and pure multi-layer perceptrons (MLP) models
2 for vision have shown great potential in achieving promising performance with
3 fewer inductive biases. These models are generally based on learning interaction
4 among spatial locations from raw data. The complexity of self-attention and MLP
5 grows quadratically as the image size increases, which makes these models hard
6 to scale up when high-resolution features are required. In this paper, we present
7 the Global Filter Network, a conceptually simple yet computationally efficient
8 architecture, that learns long-term spatial dependencies in the frequency domain
9 with log-linear complexity. Our architecture replaces the self-attention layer in
10 vision transformers with three key operations: a 2D discrete Fourier transform,
11 an element-wise multiplication between frequency-domain features and learnable
12 global filters, and a 2D inverse Fourier transform. We exhibit favorable accuracy/
13 complexity trade-offs of our models on ImageNet, which achieve competitive
14 performance with vision transformers while maintaining the high efficiency of
15 MLP models.

16 1 Introduction

17 The transformer architecture, originally designed for the natural language processing (NLP) tasks [33],
18 has shown promising performance on various vision problems recently [7, 31, 21, 40, 3]. Different
19 from convolutional neural networks (CNNs), vision transformer models use self-attention layers to
20 capture long-term dependencies, which are able to learn more diverse interactions between spatial
21 locations. The pure multi-layer perceptrons (MLP) models [29, 30] further simplify the vision
22 transformers by replacing the self-attention layers with MLPs that are applied across spatial locations.
23 Since fewer inductive biases are introduced, these two kinds of models have the potential to learn
24 more generic and flexible interactions among spatial locations from raw data.

25 One primary challenge of applying self-attention and pure MLP models to vision tasks is the
26 considerable computational complexity that grows quadratically as the number of tokens increases.
27 Therefore, typical vision transformer style models usually consider a relatively small resolution for
28 the intermediate features (*e.g.* 14×14 tokens are extracted from the input images in both ViT [7]
29 and MLP-Mixer [29]). This design may limit the applications of downstream dense prediction tasks
30 like detection and segmentation. A possible solution is to replace the global self-attention with
31 several local self-attention like Swin transformer [21]. Despite the effectiveness in practice, local
32 self-attention brings quite a few hand-made choices (*e.g.*, window size, padding strategy, *etc.*) and
33 limits the receptive field of each layer, which may hurt the generality of this kind of models.

34 In this paper, we present a new conceptually simple yet computationally efficient architecture
35 called Global Filter Network (*GFNet*), which follows the trend of removing inductive biases from
36 vision models while enjoying the log-linear complexity in computation. The basic idea behind our
37 architecture is to learn the interactions among spatial locations in the frequency domain. Different
38 from the self-attention mechanism in vision transformers and the fully connected layers in MLP

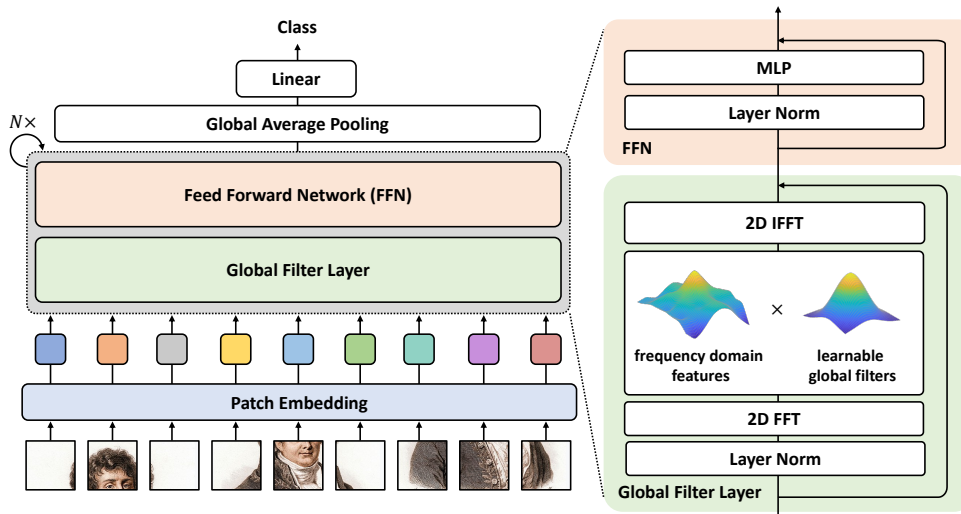


Figure 1: **The overall architecture of the Global Filter Network.** Our architecture is based on Vision Transformer (ViT) models with some minimal modifications. We replace the self-attention sub-layer with the proposed *global filter layer*, which consists of three key operations: a 2D discrete Fourier transform to convert the input spatial features to the frequency domain, an element-wise multiplication between frequency-domain features and the global filters, and a 2D inverse Fourier transform to map the features back to the spatial domain. The efficient fast Fourier transform (FFT) enables us to learn arbitrary interactions among spatial locations with log-linear complexity.

	Complexity (FLOPs)	# Parameters
Depthwise Convolution	$k^2 HWD$	$k^2 D$
Self-Attention	$4HWD^2 + 2H^2W^2D$	$4D^2$
Spatial MLP	H^2W^2D	H^2W^2
Global Filter	$HWD[\log_2(HW)] + HWD$	HWD

Table 1: Comparisons of the proposed *Global Filter* with prevalent operations in deep vision models. H , W and D are the height, width and the number of channels of the feature maps. k is the kernel size of the convolution operation. The proposed global filter is much more efficient than self-attention and spatial MLP.

39 models, the interactions among tokens are modeled as a set of learnable *global filters* that are applied
 40 to the spectrum of the input features. Since the global filters are able to cover all the frequencies,
 41 our model can capture both long-term and short-term interactions. The filters are directly learned
 42 from the raw data without introducing human priors. Our architecture is largely based on the vision
 43 transformers only with some minimal modifications. We replace the self-attention sub-layer in vision
 44 transformers with three key operations: a 2D discrete Fourier transform to convert the input spatial
 45 features to the frequency domain, an element-wise multiplication between frequency-domain features
 46 and the global filters, and a 2D inverse Fourier transform to map the features back to the spatial
 47 domain. Since the Fourier transform is used to mix the information of different tokens, the global filter
 48 is much more efficient compared to the self-attention and MLP thanks to the $\mathcal{O}(L \log L)$ complexity
 49 of the fast Fourier transform algorithm (FFT) [4]. Benefiting from this, the proposed global filter
 50 layer is less sensitive to the token length L and thus is compatible with larger feature maps and
 51 CNN-style hierarchical architectures. The overall architecture of GFNet is illustrated in Figure 1. We
 52 also compare the proposed global filter with prevalent operations in deep vision models in Table 1.

53 Our experiments on ImageNet verify the effectiveness of GFNet. Our method achieves competitive
 54 performance with DeiT [31] while maintaining the high efficiency of MLP-like models [30]. With
 55 a similar architecture, our model outperform the recent MLP models including ResMLP [30] and
 56 gMLP [20]. When using the hierarchical architecture, GFNet can further enlarge the gap. GFNet
 57 also works well on downstream transfer learning datasets. Our results demonstrate that GFNet can be
 58 a very competitive alternative to vision transformers and MLP-like models in accuracy/complexity
 59 trade-offs.

60 2 Related works

61 **Vision transformers.** Since Dosovitskiy *et al.* [7] introduce transformers to the image classifica-
62 tion and achieve a competitive performance compared to CNNs, transformers begin to exhibit their
63 potential in various vision tasks [2, 3, 40]. Recently, there are a large number of works which aim
64 to improve the transformers [31, 32, 21, 34, 13, 8, 36]. These works either seek for better training
65 strategies [31, 8] or design better architectures [21, 34, 36] or both [32, 8]. However, most of the
66 architecture modification of the transformers [34, 13, 21, 36] introduces additional inductive biases
67 similar to CNNs. In this work, we only focus on the standard transformer architecture [7, 31] and
68 our goal is to replace the heavy self-attention layer ($\mathcal{O}(L^2)$) to an more efficient operation which can
69 still model the interactions among different spatial locations without introducing the inductive biases
70 associated with CNNs.

71 **MLP-like models.** More recently, there are several works that question the importance of self-
72 attention in the vision transformers and propose to use MLP to replace the self-attention layer in the
73 transformers [29, 30, 20]. The MLP-Mixer [29] employs MLPs to perform token mixing and channel
74 mixing alternatively in each block. ResMLP [30] adopts a similar idea but substitutes the Layer
75 Normalization with an Affine transformation for acceleration. The recently proposed gMLP [20] uses
76 a spatial gating unit to re-weight tokens in the spatial dimension. However, all of the above models
77 include MLPs to mix the tokens spatially, which brings two drawbacks: (1) like the self-attention
78 in the transformers, the spatial MLP still requires computational complexity quadratic to the length
79 of tokens. (2) unlike transformers, MLP models are hard to scale up to higher resolution since the
80 weights of the spatial MLPs have fixed sizes. Our work follows this trend and successfully resolves
81 the above issues in MLP-like models. The proposed GFNet enjoys log-linear complexity and can be
82 easily scaled up to any resolution.

83 **Applications of Fourier transform in vision.** Fourier transform has been an important tool in
84 digital image processing for decades [25, 1]. With the breakthroughs of CNNs in vision [10, 9],
85 there are a variety of works that start to incorporate Fourier transform in some deep learning
86 method [19, 35, 6, 17]. Some of these works employ discrete Fourier transform to convert the images
87 to the frequency domain and leverage the frequency information to improve the performance in
88 certain tasks [17, 35], while others utilize the convolution theorem to accelerate the CNNs via fast
89 Fourier transform (FFT) [19, 6]. In this work, we propose to use learnable filters to interchange
90 information among the tokens in the Fourier domain, inspired by the frequency filters in the digital
91 image processing [25]. We also take advantage of some properties of FFT to reduce the computational
92 costs and the number of parameters.

93 3 Method

94 3.1 Preliminaries: discrete Fourier transform

95 We start by introducing the discrete Fourier transform (DFT), which plays an important role in the
96 area of digital signal processing and is a crucial component in our GFNet. For clarity, We first
97 consider the 1D DFT. Given a sequence of N complex numbers $x[n]$, $0 \leq n \leq N - 1$, the 1D DFT
98 converts the sequence into the frequency domain by:

$$X[k] = \sum_{n=0}^{N-1} x[n] e^{-j(2\pi/N)kn} := \sum_{n=0}^{N-1} x[n] W_N^{kn} \quad (3.1)$$

99 where j is the imaginary unit and $W_N = e^{-j(2\pi/N)}$. The formulation of DFT in Equation (3.1) can
100 be derived from the Fourier transform for continuous signal by sampling in both the time domain and
101 the frequency domain (see Appendix A for details). Since $X[k]$ repeats on intervals of length N , it
102 is suffice to take the value of $X[k]$ at N consecutive points $k = 0, 1, \dots, N - 1$. Specifically, $X[k]$
103 represents to the spectrum of the sequence $x[n]$ at the frequency $\omega_k = 2\pi k/N$.

104 It is also worth noting that DFT is a one-to-one transformation. Given the DFT $X[k]$, we can recover
105 the original signal $x[n]$ by the inverse DFT (IDFT):

$$x[n] = \frac{1}{N} \sum_{k=0}^{N-1} X[k] e^{j(2\pi/N)kn}. \quad (3.2)$$

Algorithm 1 Pseudocode of Global Filter Layer.

```
# x: the token features, B x H x W x D (where N = H * W)
# K: the filter kernel, H x W_hat x D (where W_hat = W // 2 + 1, see Section 3.2 for details)

X = rfft2(x, dim=(1, 2))
X_tilde = X * K
x = x + irfft2(X_tilde, dim=(1, 2))
```

rfft2/irfft2: 2D FFT/IFFT for real signal

106 For real input $x[n]$, it can be proved that (see Appendix A) its DFT is conjugate symmetric, i.e.,
107 $X[N - k] = X^*[k]$. The reverse is true as well: if we perform IDFT to $X[k]$ which is conjugate
108 symmetric, a real discrete signal can be recovered. This property implies that the half of the DFT
109 $\{X[k] : 0 \leq k \leq \lceil N/2 \rceil\}$ contains the full information about the frequency characteristics of $x[n]$.

110 DFT is widely used in modern signal processing algorithms for mainly two reasons: (1) the input and
111 output of DFT are both discrete thus can be easily processed by computers; (2) there exist efficient
112 algorithms for computing the DFT. The *fast Fourier transform* (FFT) algorithms take advantage of
113 the symmetry and periodicity properties of W_N^{kn} and reduce the complexity to compute DFT from
114 $\mathcal{O}(N^2)$ to $\mathcal{O}(N \log N)$. The inverse DFT (3.2), which has a similar form to the DFT, can also be
115 computed efficiently using the inverse fast Fourier transform (IFFT).

116 The DFT described above can be extend to 2D signals. Given the 2D signal $X[m, n], 0 \leq m \leq$
117 $M - 1, 0 \leq n \leq N - 1$, the 2D DFT of $x[m, n]$ is given by:

$$X[u, v] = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} x[m, n] e^{-j2\pi(\frac{um}{M} + \frac{vn}{N})}. \quad (3.3)$$

118 The 2D DFT can be viewed as performing 1D DFT on the two dimensions alternatively. Similar to 1D
119 DFT, 2D DFT of real input $x[m, n]$ satisfied the conjugate symmetry property $X[M - u, N - v] =$
120 $X^*[u, v]$. The FFT algorithms can also be applied to 2D DFT to improve computational efficiency.

121 3.2 Global Filter Networks

122 **Overall architecture.** Recent advances in vision transformers [7, 31] demonstrate that models
123 based on self-attention can achieve competitive performance even without the inductive biases
124 associated with the convolutions. Henceforth, there are several works [30, 29] that exploit approaches
125 (e.g., MLPs) other than self-attention to mix the information among the tokens. The proposed Global
126 Filter Networks (GFNet) follows this line of work and aims to replace the heavy self-attention layer
127 ($\mathcal{O}(N^2)$) with a simpler and more efficient one.

128 The overall architecture of our model is depicted in Figure 1. Our model takes as an input $H \times W$
129 non-overlapping patches and projects the flattened patches into $L = HW$ tokens with dimension
130 D . The basic building block of GFNet consists of: 1) a *global filter layer* that can exchange spatial
131 information efficiently ($\mathcal{O}(L \log L)$); 2) a feedforward network (FFN) as in [7, 31]. The output tokens
132 of the last block are fed into a global average pooling layer followed by a linear classifier.

133 **Global filter layer.** We propose global filter layer as an alternative to the self-attention layer which
134 can mix tokens representing different spatial locations. Given the tokens $\mathbf{x} \in \mathbb{R}^{H \times W \times D}$, we first
135 perform 2D FFT (see Section 3.1) along the spatial dimensions to convert \mathbf{x} to the frequency domain:
136

$$\mathbf{X} = \mathcal{F}[\mathbf{x}] \in \mathbb{C}^{H \times W \times D}, \quad (3.4)$$

137 where $\mathcal{F}[\cdot]$ denotes the 2D FFT. Note that \mathbf{X} is a complex tensor and represents the spectrum of \mathbf{x} .
138 We can then modulate the spectrum by multiplying a learnable filter $\mathbf{K} \in \mathbb{C}^{H \times W \times D}$ to the \mathbf{X} :

$$\tilde{\mathbf{X}} = \mathbf{K} \odot \mathbf{X}, \quad (3.5)$$

139 where \odot is the element-wise multiplication (also known as the Hadamard product). The filter \mathbf{K} is
140 called the *global filter* since it has the same dimension with \mathbf{X} , which can represent an arbitrary filter
141 in the frequency domain. Finally, we adopt the inverse FFT to transform the modulated spectrum $\tilde{\mathbf{X}}$
142 back to the spatial domain and update the tokens via a residual connection [10],

$$\mathbf{x} \leftarrow \mathbf{x} + \mathcal{F}^{-1}[\tilde{\mathbf{X}}]. \quad (3.6)$$

143 The formulation of the global filter layer is motivated by the frequency filters in the digital image
 144 processing [25], where the global filter \mathbf{K} can be regarded as a set of learnable frequency filters
 145 for different hidden dimensions. It can be proved (see Appendix A) that the global filter layer is
 146 equivalent to a depthwise *global circular convolution* with the filter size $H \times W$. Therefore, the
 147 global filter layer is different from the standard convolutional layer which adopts a relatively small
 148 filter size to enforce the inductive biases of the locality. We also find although the proposed global
 149 filter can also be interpreted as a spatial domain operation, the filters learned in our networks exhibit
 150 more clear patterns in the frequency domain than the spatial domain, which indicates our models tend
 151 to capture relation in the frequency domain instead of spatial domain (see Figure 4). Note that the
 152 global filter implemented in the frequency domain is also much more efficient compared to the spatial
 153 domain, which enjoys a complexity of $\mathcal{O}(DL \log L)$ while the vanilla depthwise global circular
 154 convolution in the spatial domain has $\mathcal{O}(DL^2)$ complexity. We will also show that the global filter
 155 layer is better than its local convolution counterparts in the experiments.

156 It is also worth noting that in the implementation, we make use of the property of DFT to reduce the
 157 redundant computation. Since \mathbf{x} is a real tensor, its DFT \mathbf{X} is conjugate symmetric, *i.e.* $\mathbf{X}[H -$
 158 $u, W - v, :] = \mathbf{X}^*[H, W, :]$. Therefore, we can take only the half of the values in the \mathbf{X} but preserve
 159 the full information at the same time:

$$\mathbf{X}_r = \mathbf{X}[:, 0 : \widehat{W}] := \mathcal{F}_r[\mathbf{x}], \quad \widehat{W} = \lceil W/2 \rceil, \quad (3.7)$$

160 Where \mathcal{F}_r denotes the 2D FFT for real input. In this way, we can implement the global filter as
 161 $\mathbf{K}_r \in \mathbb{C}^{H \times \widehat{W} \times D}$, which can reduce half the parameters. This can also ensure $\mathcal{F}_r^{-1}[\mathbf{K}_r \odot \mathbf{X}_r]$ is a
 162 real tensor, thus it can be added directly to the input \mathbf{x} . The global filter layer can be easily in modern
 163 deep learning frameworks (*e.g.*, PyTorch [24]), as is shown in Algorithm 1. The FFT and ITTF are
 164 well supported by GPU and CPU thanks to the acceleration libraries like cuFFT and `mk1-fft`, which
 165 makes our models perform well on hardware.

166 **Relationship to vision transformers and pure MLP models.** The GFNet follows the line of
 167 research about the exploration of approaches to mix the tokens. Compared to existing architectures
 168 like vision transformers and pure MLP models, we exhibit that GFNet has several favorable properties:
 169 1) GFNet is more efficient. The complexity of both the vision transformers [7, 31, 32] and the MLP
 170 models [29, 30] is $\mathcal{O}(L^2)$. Different from them, global filter layer only consists an FFT ($\mathcal{O}(L \log L)$),
 171 an element-wise multiplication ($\mathcal{O}(L)$) and an IFFT ($\mathcal{O}(L)$), which means the total computational
 172 complexity is $\mathcal{O}(L \log L)$. 2) Although pure MLP models are simpler compared to transformers, it is
 173 hard to fine-tune them on higher resolution (*e.g.*, from 224×224 resolution to 384×384 resolution)
 174 since they can only process a fixed number of tokens. As opposed to pure MLP models, we will show
 175 that our GFNet can be easily scaled up to higher resolution. Our model is more flexible since both the
 176 FFT and the IFFT have no learnable parameters and can process sequences with arbitrary length. We
 177 can simply interpolate the global filter \mathbf{K} to $\mathbf{K}' \in \mathbb{C}^{H' \times W' \times D}$ for different inputs, where $H' \times W'$
 178 is the target size. The interpolation is reasonable due to the property of DFT. Each element of the
 179 global filter $\mathbf{K}[u, v]$ corresponds to the spectrum of the filter at $\omega_u = 2\pi u/H, \omega_v = 2\pi v/W$ and
 180 thus, the global filter \mathbf{K} can be viewed as a sampling of a continuous spectrum $\mathbf{K}(\omega_u, \omega_v)$, where
 181 $\omega_u, \omega_v \in [0, 2\pi]$. Hence, changing the resolution is equivalent to changing the sampling interval of
 182 $\mathbf{K}(\omega_u, \omega_v)$. Therefore, we only need to perform interpolation to shift from one resolution to another.

183 We also notice recently a concurrent work FNet [18] leverages Fourier transform to mix tokens. Our
 184 work is distinct from FNet in three aspects: (1) FNet performs FFT to the input and directly adds the
 185 real part of the spectrum to the input tokens, which blends the information from different domains
 186 (spatial/frequency) together. On the other hand, GFNet draws motivation from the frequency filters,
 187 which is more reasonable. (2) FNet only keeps the real part of the spectrum. Note that the spectrum of
 188 real input is conjugate symmetric, which means the real part is exactly symmetric and thus contains
 189 redundant information. Our GFNet, however, utilizes this property to simplify the computation. (3)
 190 FNet is designed for NLP tasks, while our GFNet focuses on vision tasks. In our experiments, we
 191 also implement the FNet and show that our model outperforms it.

192 **Architecture variants.** Due to the limitation from the quadratic complexity in the self-attention,
 193 vision transformers [7, 31] are usually designed to process a relatively small feature map (*e.g.*,
 194 14×14). However, our GFNet, which enjoys log-linear complexity, avoids that problem. Since in
 195 our GFNet the computational costs do not grow significantly when the feature map size increases, we
 196 can adopt a hierarchical architecture inspired by the success of CNNs [16, 10]. Generally speaking,

197 we can start from a large feature map (*e.g.*, 56×56) and gradually perform downsampling after
198 a few blocks. In this paper, we mainly investigate three variants of GFNet: GFNet-12 that has
199 a similar architecture to DeiT-S [31] and ResMLP-12 [30], GFNet-H18 and GFNet-H24 that are
200 two hierarchical architectures with higher efficiency and performance respectively. The detailed
201 architecture can be found in Appendix B.

202 **Training techniques.** In our experiments, we find the default training procedure introduced for
203 DeiT [31] can already yield good performance for our GFNet. However, as is analyzed in [8],
204 transformers-like architectures may suffer from over-smoothing due to the global receptive field in
205 every layer. Inspired by [38], we employ the token-labeling [13] technique to provide each token a
206 supervision signal. For the sake of simplicity, we implement it by using a well-trained model (LV-
207 ViT [13]) as the teacher model and perform token-wise distillation (see Appendix B). Experiments
208 show that this can overcome over-smoothing and further enhance the performance of GFNet.

209 4 Experiments

210 We conduct extensive experiments to verify the effectiveness of our GFNet architecture on the
211 image classification task. We present the main results on ImageNet [5] and compare them with
212 various architectures. We also test our models on the downstream transfer learning datasets including
213 CIFAR-10/100 [15], Stanford Cars [14] and Flowers-102 [23]. Finally, we investigate the effects of
214 the proposed global filters and provide visualization to have an intuitive understanding of our method.

215 4.1 ImageNet results

216 **Setups.** We conduct our main experiments on ImageNet [5], which is a widely used large-scale
217 benchmark for image classification. ImageNet contains roughly 1.2M images from 1,000 categories.
218 Following common practice [10, 31], we train our models on the training set of ImageNet and
219 report the single-crop top-1 accuracy on 50,000 validation images. To fairly compare with previous
220 works [31, 30], we follow the most training details for our models. We train our models for 300
221 epochs using the AdamW optimizer [22]. We set the initial learning rate as $\frac{\text{batch size}}{1024} \times 0.001$ and decay
222 the learning rate to $1e^{-5}$ using the cosine schedule. We use a linear warm-up learning rate in the first
223 5 epochs and apply gradient clipping to stabilize the training process. For GFNet-12, we faithfully
224 follow the regularization strategies used in DeiT [31] including Mixup [39], CutMix [37], EMA
225 model [26], RandomErase [41] and repeated augmentation [11]. For GFNet-H18 and GFNet-24,
226 we follow most settings of [31] except for the EMA model and repeated augmentation as suggested
227 in [21]. We set the stochastic depth coefficient [12] to 0, 0.1 and 0.3 for GFNet-12, GFNet-H18 and
228 GFNet-H24 respectively. During finetuning at the higher resolution, we use the hyper-parameters
229 suggested by the implementation of [31] and train the model for 30 epochs. All of our models are
230 trained on a single machine with 8 GPUs.

231 **Main results.** The main results are presented in Table 2. We compare our method with different
232 architectures for image classification including convolutional networks, vision transformers, and MLP-
233 like models that have similar complexity and number of parameters. One can find that our method
234 can clearly outperform recent MLP-like models like ResMLP [30] and gMLP [20]. Specifically,
235 GFNet-12 outperforms ResMLP-12 by 1.9% while having slightly fewer FLOPs. GFNet-12 also
236 achieves competitive results compared to gMLP-S and has 1.7 GFLOPs less computation. Our
237 hierarchical models even have more significant advantages. Benefiting from the log-linear complexity,
238 GFNet-H18 is 2.6% more accurate than ResMLP-12 while having 36% fewer FLOPs. These results
239 show the high efficiency of the global filter.

240 **Fine-tuning at higher resolution.** One prominent problem of MLP-like models is that the feature
241 resolution is not adjustable. On the contrary, the proposed global filter is more flexible. We
242 demonstrate the advantage of GFNet by finetuning the model trained at 224×224 resolution to
243 higher resolution following the practice in vision transformers [31]. As shown in Table 2, our model
244 can easily adapt to higher resolution with only 30 epoch finetuning and achieve better performance.

245 **Token-wise distillation.** We address the over-smoothing problem by introducing token-wise dis-
246 tillation supervision. This strategy improves the performance of our models by around 1% without

Table 2: **Main results on ImageNet.** We compare different architectures for image classification including convolutional networks, vision transformers, MLP-like models and our method that have comparable FLOPs and number of parameters. We report the top-1 accuracy on the validation set of ImageNet as well as the number of parameters and FLOPs. All of our models are trained with 224×224 images. We use “ $\uparrow 384$ ” to represent models finetuned on 384×384 images for 30 epochs. \dagger indicates the models trained with token-wise distillation.

	Model	Params (M)	FLOPs (G)	Res.	Top-1 Acc. (%)	Top-5 Acc. (%)
ConvNets	EfficientNet-B3 [28]	12	1.8	300	81.6	95.7
	EfficientNet-B4 [28]	19	4.2	380	82.9	96.4
	EfficientNet-B5 [28]	30	9.9	456	83.6	96.7
	RegNetY-4GF [27]	21	4.0	224	80.0	-
	RegNetY-8GF [27]	39	8.0	224	81.7	-
	RegNetY-16GF [27]	84	16.0	224	82.9	-
Transformers	ViT-B/16 [7]	86	55.4	384	77.9	-
	ViT-L/16 [7]	307	190.7	384	76.5	-
	DeiT-Ti [31]	5	1.2	224	72.2	91.1
	DeiT-S [31]	22	4.6	224	79.8	95.0
	DeiT-B [31]	86	17.5	224	81.8	95.6
MLP-like Models	Mixer-B/16	59	12.7	224	76.4	-
	ResMLP-12 [30]	15	3.0	224	76.6	-
	ResMLP-24 [30]	30	6.0	224	79.4	-
	ResMLP-36 [30]	45	8.9	224	79.7	-
	gMLP-Ti [20]	6	1.4	224	72.0	-
	gMLP-S [20]	20	4.5	224	79.4	-
	gMLP-B [20]	73	15.8	224	81.6	-
Ours	GFNet-12	16	2.8	224	78.5	94.0
	GFNet-H18	17	1.9	224	79.2	94.5
	GFNet-H24	37	5.6	224	81.0	95.2
	GFNet-12 $\uparrow 384$	18	8.4	384	80.5	95.1
	GFNet-H18 $\uparrow 384$	20	5.5	384	80.8	95.3
	GFNet-H24 $\uparrow 384$	42	16.5	384	82.0	95.6
	GFNet-12 \dagger	16	2.9	224	79.7	94.7
	GFNet-H18 \dagger	17	1.9	224	79.9	94.9
	GFNet-H24 \dagger	37	5.7	224	82.0	95.8
	GFNet-12 $\uparrow 384$ \dagger	18	8.6	384	81.5	95.7
	GFNet-H18 $\uparrow 384$ \dagger	20	5.6	384	81.5	95.7
	GFNet-H24 $\uparrow 384$ \dagger	43	16.6	384	83.1	96.4

247 bringing significantly large computation during inference. By incorporating both token-wise distilla-
 248 tion and larger image resolution, the performance can be further improved.

249 4.2 Transfer learning results

250 **Setups.** To test the generality of our architecture and the learned representation, we evaluate
 251 GFNet on a set of commonly used transfer learning benchmark datasets including CIFAR-10 [15],
 252 CIFAR-100 [15], Stanford Cars [14] and Flowers-102 [23]. We follow the setting of previous
 253 works [28, 7, 31, 30], where the model is initialized by the ImageNet pre-trained weights and
 254 finetuned on the new datasets. During finetuning, we use the SGD optimizer and set the weight
 255 decay to $1e^{-4}$. We use batch size 512 and a smaller initial learning rate of 0.0001 with cosine decay.
 256 Linear learning rate warm-up in the first 5 epochs and gradient clipping with a max norm of 1 are
 257 also applied to stabilize the training. We keep most of the regularization methods unchanged except
 258 for removing Random Erasing and stochastic depth following [31].

Table 3: **Results on transfer learning datasets.** We report the top-1 accuracy on the four datasets as well as the number of parameters and FLOPs.

Model	FLOPs	Params	CIFAR-10	CIFAR-100	Flowers-102	Cars-196
ResNet50 [10]	4.1G	26M	-	-	96.2	90.0
EfficientNet-B7 [28]	37G	66M	98.9	91.7	98.8	94.7
ViT-B/16 [7]	55.4G	86M	98.1	87.1	89.5	-
ViT-L/16 [7]	190.7G	307M	97.9	86.4	89.7	-
DeiT-B/16 [31]	17.5G	86M	99.1	90.8	98.4	92.1
ResMLP-12 [30]	3.0G	15M	98.1	87.0	97.4	84.6
ResMLP-24 [30]	6.0G	30M	98.7	89.5	97.9	89.5
GFNet-12	2.8G	16M	98.6	89.0	98.1	92.4
GFNet-H24	5.6G	37M	98.8	89.3	98.5	93.1

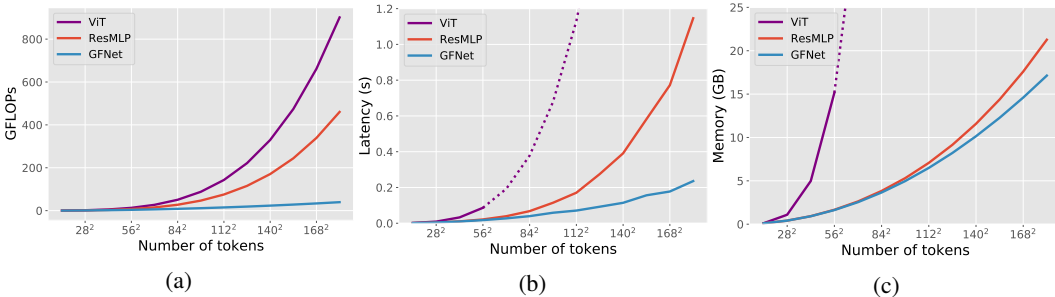


Figure 2: Comparisons among GFNet, ViT [7] and ResMLP [30] in (a) FLOPs (b) latency and (c) GPU memory with respect to the number of tokens (feature resolution). The dotted lines indicate the estimated values when the GPU memory has run out. The latency and GPU memory is measured using a single NVIDIA RTX 3090 GPU with batch size 32 and feature dimension 384.

259 **Results.** We evaluate the transfer learning performance of our basic model and best model. The
 260 results are presented in Table 3. The proposed models generally work well on downstream datasets.
 261 GFNet-H24 outperforms ResMLP-24 on three out of four datasets and achieves very close accuracy
 262 on the left one dataset. Our models also show competitive performance compared to state-of-the-art
 263 CNNs and vision transformers.

264 4.3 Analysis and visualization

265 **Efficiency of GFNet.** We demonstrate the efficiency of our GFNet in Figure 2, where the models
 266 are compared in theoretical FLOPs, actual latency and peak memory usage on GPU. We test a single
 267 building block of each model (including one token mixing layer and one FFN) with respect to the
 268 different numbers of tokens and set the feature dimension and batch size to 384 and 32 respectively.
 269 The self-attention model quickly runs out of memory when feature resolution exceeds 56², which
 270 is also the feature resolution of our hierarchical model. The advantage of the proposed architecture
 271 becomes larger as the resolution increases, which strongly shows the potential of our model in vision
 272 tasks requiring high-resolution feature maps.

273 **Complexity/accuracy trade-offs.** We show the complexity and accuracy trade-offs of various
 274 architectures in Figure 3. GFNet achieves the best trade-off among kinds of models.

275 **Ablation study on the global filter.** To more clearly show the effectiveness of the proposed global
 276 filters, we compare GFNet-12 with several baseline models that are equipped with different token
 277 mixing operations. The results are presented in Table 4. All models have a similar building block (
 278 token mixing layer + FFN) and the same feature dimension of $D = 384$. We also implement the
 279 recent FNet [18] for comparison, where a 1D FFT on feature dimension and a 2D FFT on spatial
 280 dimensions are used to mix tokens. As shown in Table 4, our method outperforms all baseline
 281 methods except DeiT-S that has 64% higher FLOPs.

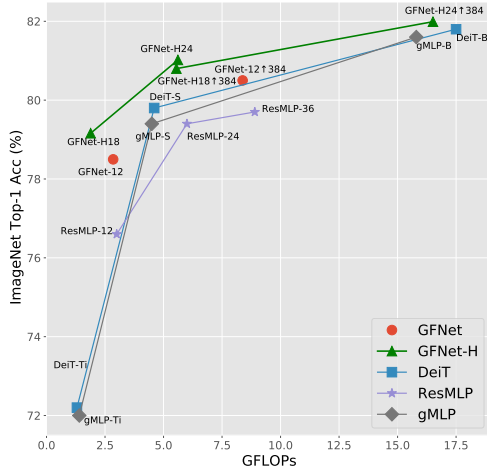


Figure 3: ImageNet acc. vs model complexity.

Table 4: Comparisons among the GFNet and other variants based on the transformer-like architecture on ImageNet. We show that GFNet outperforms the ResMLP [30], FNet [18] and models with local depth-wise convolutions. We also report the number of parameters and theoretical complexity in FLOPs.

Model	Acc (%)	Param (M)	FLOPs (G)
DeiT-S [31]	79.8	22	4.6
Local Conv (3×3)	77.7	15	2.8
Local Conv (5×5)	78.1	15	2.9
Local Conv (7×7)	78.2	15	2.9
ResMLP [30]	76.6	15	3.0
FNet [18]	71.2	15	2.9
GFNet-12	78.5	16	2.8

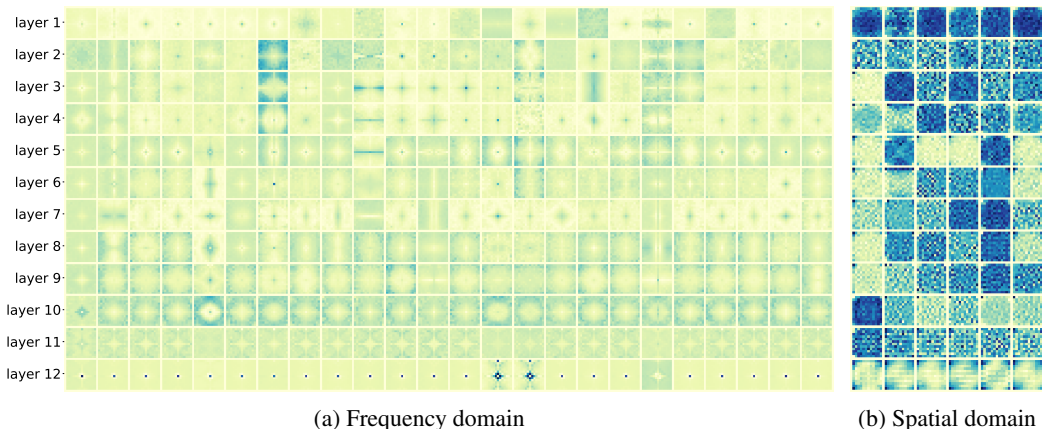


Figure 4: Visualization of the learned *global filters* in GFNet-12. We visualize the original frequency domain global filters in (a) and show the corresponding spatial domain filters for the first 6 columns in (b). There are more clear patterns in the frequency domain than spatial domain.

282 **Visualization.** Core operation in GFNet is the element-wise multiplication between frequency-
 283 domain features and the global filter. Therefore, it is easy to visualize and interpret. We visualize the
 284 frequency domain filters as well as their corresponding spatial domain filters in Figure 4. The learned
 285 global filters have more clear patterns in the frequency domain, where different layers have different
 286 characteristics. Interestingly, the filters in the last layer particularly focus on the low-frequency
 287 component. The corresponding filters in the spatial domain are less interpretable for humans.

288 5 Conclusion

289 We have presented the Global Filter Network (*GFNet*), which is a conceptually simple yet computa-
 290 tionally efficient architecture for image classification. Our model replaces the self-attention sub-layer
 291 in vision transformer with 2D FFT/IFFT and a set of learnable *global filters* in the frequency domain.
 292 Benefiting from the token mixing operation with log-linear complexity, our architecture is highly
 293 efficient. Our experimental results demonstrated that GFNet can be a very competitive alternative to
 294 vision transformers and MLP-like models in accuracy/complexity trade-offs. Since we only focus
 295 on vision tasks in this paper, extending our method to broader applications like NLP and speech
 296 understanding would be interesting future directions.

297 **References**

- 298 [1] Gregory A Baxes. *Digital image processing: principles and applications*. John Wiley & Sons,
299 Inc., 1994. 3
- 300 [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and
301 Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, pages 213–229.
302 Springer, 2020. 3
- 303 [3] Xin Chen, Bin Yan, Jiawen Zhu, Dong Wang, Xiaoyun Yang, and Huchuan Lu. Transformer
304 tracking. In *CVPR*, 2021. 1, 3
- 305 [4] James W Cooley and John W Tukey. An algorithm for the machine calculation of complex
306 fourier series. *Mathematics of computation*, 19(90):297–301, 1965. 2
- 307 [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
308 hierarchical image database. In *CVPR*, pages 248–255, 2009. 6
- 309 [6] Caiwen Ding, Siyu Liao, Yanzhi Wang, Zhe Li, Ning Liu, Youwei Zhuo, Chao Wang, Xuehai
310 Qian, Yu Bai, Geng Yuan, et al. Circnn: accelerating and compressing deep neural networks
311 using block-circulant weight matrices. In *MICRO*, pages 395–408, 2017. 3
- 312 [7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
313 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
314 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image
315 recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 1, 3, 4, 5, 7, 8
- 316 [8] Chengyue Gong, Dilin Wang, Meng Li, Vikas Chandra, and Qiang Liu. Improve vision
317 transformers training by suppressing over-smoothing. *arXiv preprint arXiv:2104.12753*, 2021.
318 3, 6
- 319 [9] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. 3
- 320 [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
321 recognition. In *CVPR*, pages 770–778, 2016. 3, 4, 5, 6, 8
- 322 [11] Elad Hoffer, Tal Ben-Nun, Itay Hubara, Niv Giladi, Torsten Hoefer, and Daniel Soudry.
323 Augment your batch: Improving generalization through instance repetition. In *CVPR*, pages
324 8129–8138, 2020. 6
- 325 [12] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with
326 stochastic depth. In *ECCV*, pages 646–661, 2016. 6
- 327 [13] Zihang Jiang, Qibin Hou, Li Yuan, Daquan Zhou, Xiaojie Jin, Anran Wang, and Jiashi Feng.
328 Token labeling: Training a 85.5% top-1 accuracy vision transformer with 56m parameters on
329 imagenet. *arXiv preprint arXiv:2104.10858*, 2021. 3, 6
- 330 [14] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for
331 fine-grained categorization. In *ICCVW*, pages 554–561, 2013. 6, 7
- 332 [15] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
333 2009. 6, 7
- 334 [16] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep
335 convolutional neural networks. *NeurIPS*, 25:1097–1105, 2012. 5
- 336 [17] Jae-Han Lee, Minhyeok Heo, Kyung-Rae Kim, and Chang-Su Kim. Single-image depth
337 estimation based on fourier domain analysis. In *CVPR*, pages 330–339, 2018. 3
- 338 [18] James Lee-Thorp, Joshua Ainslie, Ilya Eckstein, and Santiago Ontanon. Fnet: Mixing tokens
339 with fourier transforms. *arXiv preprint arXiv:2105.03824*, 2021. 5, 8, 9
- 340 [19] Shaohua Li, Kaiping Xue, Bin Zhu, Chenkai Ding, Xindi Gao, David Wei, and Tao Wan. Falcon:
341 A fourier transform based approach for fast and secure convolutional neural network predictions.
342 In *CVPR*, pages 8705–8714, 2020. 3

- 343 [20] Hanxiao Liu, Zihang Dai, David R So, and Quoc V Le. Pay attention to mlps. *arXiv preprint*
344 *arXiv:2105.08050*, 2021. 2, 3, 6, 7
- 345 [21] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining
346 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv preprint*
347 *arXiv:2103.14030*, 2021. 1, 3, 6
- 348 [22] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
349 *arXiv:1711.05101*, 2017. 6
- 350 [23] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large
351 number of classes. In *ICVGIP*, pages 722–729, 2008. 6, 7
- 352 [24] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
353 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative
354 style, high-performance deep learning library. *NeurIPS*, 2019. 5
- 355 [25] Ioannis Pitas. *Digital image processing algorithms and applications*. John Wiley & Sons, 2000.
356 3, 5
- 357 [26] Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging.
358 *SIAM journal on control and optimization*, 30(4):838–855, 1992. 6
- 359 [27] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Design-
360 ing network design spaces. In *CVPR*, 2020. 7
- 361 [28] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural
362 networks. In *ICML*, pages 6105–6114. PMLR, 2019. 7, 8
- 363 [29] Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas
364 Unterthiner, Jessica Yung, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, et al. Mlp-mixer: An
365 all-mlp architecture for vision. *arXiv preprint arXiv:2105.01601*, 2021. 1, 3, 4, 5
- 366 [30] Hugo Touvron, Piotr Bojanowski, Mathilde Caron, Matthieu Cord, Alaaeldin El-Nouby,
367 Edouard Grave, Armand Joulin, Gabriel Synnaeve, Jakob Verbeek, and Hervé Jégou. Resmlp:
368 Feedforward networks for image classification with data-efficient training. *arXiv preprint*
369 *arXiv:2105.03404*, 2021. 1, 2, 3, 4, 5, 6, 7, 8, 9
- 370 [31] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and
371 Hervé Jégou. Training data-efficient image transformers & distillation through attention. *arXiv*
372 *preprint arXiv:2012.12877*, 2020. 1, 2, 3, 4, 5, 6, 7, 8, 9
- 373 [32] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou.
374 Going deeper with image transformers. *arXiv preprint arXiv:2103.17239*, 2021. 3, 5
- 375 [33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
376 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008,
377 2017. 1
- 378 [34] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang.
379 Cvt: Introducing convolutions to vision transformers. *arXiv preprint arXiv:2103.15808*, 2021.
380 3
- 381 [35] Yanchao Yang and Stefano Soatto. Fda: Fourier domain adaptation for semantic segmentation.
382 In *CVPR*, pages 4085–4095, 2020. 3
- 383 [36] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zihang Jiang, Francis EH Tay,
384 Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch
385 on imagenet. *arXiv preprint arXiv:2101.11986*, 2021. 3
- 386 [37] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon
387 Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In
388 *ICCV*, pages 6023–6032, 2019. 6

- 389 [38] Sangdoon Yun, Seong Joon Oh, Byeongho Heo, Dongyoon Han, Junsuk Choe, and Sanghyuk
 390 Chun. Re-labeling imagenet: from single to multi-labels, from global to localized labels. *arXiv*
 391 *preprint arXiv:2101.05022*, 2021. 6
- 392 [39] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond
 393 empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017. 6
- 394 [40] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei
 395 Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a
 396 sequence-to-sequence perspective with transformers. *arXiv preprint arXiv:2012.15840*, 2020.
 397 1, 3
- 398 [41] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data
 399 augmentation. In *AAAI*, volume 34, pages 13001–13008, 2020. 6

400 Checklist

- 401 1. For all authors...
- 402 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 403 contributions and scope? [Yes]
- 404 (b) Did you describe the limitations of your work? [Yes] See Section 5.
- 405 (c) Did you discuss any potential negative societal impacts of your work? [N/A] We
 406 develop a general framework for image classification in this paper. Our model is not
 407 for specific applications.
- 408 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 409 them? [Yes]
- 410 2. If you are including theoretical results...
- 411 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 412 (b) Did you include complete proofs of all theoretical results? [N/A]
- 413 3. If you ran experiments...
- 414 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 415 mental results (either in the supplemental material or as a URL)? [Yes]
- 416 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 417 were chosen)? [Yes]
- 418 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 419 ments multiple times)? [No] We follow the common practice used in our baseline
 420 methods, where no error bars are reported.
- 421 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 422 of GPUs, internal cluster, or cloud provider)? [Yes] See our implementation details in
 423 the experiment part.
- 424 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 425 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 426 (b) Did you mention the license of the assets? [N/A]
- 427 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 428 (d) Did you discuss whether and how consent was obtained from people whose data you’re
 429 using/curating? [N/A]
- 430 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 431 information or offensive content? [N/A]
- 432 5. If you used crowdsourcing or conducted research with human subjects...
- 433 (a) Did you include the full text of instructions given to participants and screenshots, if
 434 applicable? [N/A]
- 435 (b) Did you describe any potential participant risks, with links to Institutional Review
 436 Board (IRB) approvals, if applicable? [N/A]
- 437 (c) Did you include the estimated hourly wage paid to participants and the total amount
 438 spent on participant compensation? [N/A]