# 通过稀疏分层傅里叶交互网络减少深度网络复杂性

Andrew Kiruluta and Samatha Williams

School of Information, University of California, Berkeley

2025年4月16日

#### 摘要

在这项工作中,我们介绍了稀疏分层傅里叶交互网络(SHFIN),这是一种旨在用 统一的频谱稀疏傅里叶算子替代卷积核和二次自注意力机制的新建筑原语。SHFIN 建立在三个核心组件之上:(1)分层块级快速傅里叶变换(FFT)阶段,该阶段将输 入分割为局部化的块,并对每个块进行 O(s log s)变换,在保持空间局部性的同时实 现全局信息混合;(2)一种可学习的 K 稀疏频率掩码机制,通过 Gumbel-Softmax 松 弛来实现,动态选择每块中最具信息量的 K 个频谱成分,从而修剪冗余的高频带; 以及(3)一个门控跨频混频器,在保留的频谱子空间中作为低秩双线性交互实现,以 O(K<sup>2</sup>)的成本而非 O(N<sup>2</sup>) 捕捉通道间的依赖关系。逆 FFT 和残差融合完成了 SHFIN 块,无缝集成现有的层归一化和前馈模块。

经验上,我们将 SHFIN 块集成到卷积和 Transformer 风格的主干网络中,并在 ImageNet-1k 数据集上进行了广泛实验。在 ResNet-50 和 ViT-Small 规模下,我们的 SHFIN 变体实现了可比较的 Top-1 精度(误差在 0.5 个百分点内),同时将总参数 数量减少了高达 60%,并在 NVIDIA A100 GPU 上的端到端推理延迟提高了大约 3 倍。此外,在 WMT14 英德翻译基准测试中,通过添加 SHFIN 交叉注意力层增强的 Transformer-Small 与基线 28.1 BLEU 分数相匹配,并且在训练过程中峰值 GPU 内存 使用量降低了 55%。这些结果表明,SHFIN 可以作为局部卷积和全局注意的即插即 用替代品,为高效、频谱感知深度架构提供了一条新路径。

关键词:深度神经网络;稀疏谱方法;分层FFT;低秩跨频混频器等价物

# 1 介绍

### 1.1 历史背景

过去十年见证了深度学习架构的显著进化,始于卷积神经网络(CNNs)在图像识别中的复兴。Krizhevsky等。的开创性工作表明,在超过一百万张图像上训练的深层 CNN

可以在 ImageNet 上达到前所未有的准确率,从而激发了对大规模、数据密集型模型 [11] 的广泛兴趣。在此基础上,VGG 系列引入非常小 (3C3) 的卷积滤波器来增加深度同时控制参数增长 [18];GoogLeNet 提出 inception 模块以高效捕捉多尺度特征 [19];而 ResNet的残差连接使训练超过一百层网络成为可能,且能够稳定进行 [8]。除了这些"宏观架构"改进,研究人员开发了高效变体,MobileNet 的深度可分离卷积 [9]、ShuffleNet 的通道洗牌 [26] 和 EfficientNet 的复合缩放规则 [20],以满足移动和嵌入式部署的需求。

尽管具有局部归纳偏置, CNN 在不堆叠许多层或使用大核的情况下难以捕捉长距离 依赖关系。Transformer 架构用自注意力替换了卷积,在 O(N<sup>2</sup>)时间内计算所有标记之间 的成对交互,并在机器翻译 [23]中实现了最先进的结果。视觉 Transformer (ViT)将这一 范式扩展到了图像,通过将它们分割成补丁并将每个补丁视为一个"令牌"[4]。后续工 作如 MLP-Mixer [22]、ConViT [5]和 ConvMixer [21]进一步模糊了卷积和基于注意力设 计之间的界限,但所有这些方法都继承了二次或三次缩放瓶颈,这阻碍了在资源受限硬 件上的部署。

同时,频域作为一种替代介质出现了,它可以在次二次成本下实现全局信息混合。傅 里叶神经算子(FNO)开创了将傅里叶变换应用于学习函数空间之间映射的应用,特别 是在求解偏微分方程 [13] 方面。FNet 表明用密集的 FFT 替换自注意力可以产生在自然 语言处理任务中的竞争力表现,并将复杂度降低至 O(N log N) [12]。GFNet 引入了频谱 门控以动态过滤频率 [24],而 AFNO 将频谱划分为块以提高灵活性 [7]。更近期的频域混 合模型, SpectFormer[2]、FourierFormer[17] 和频域多头自注意力(FDMHSA) [25],已经 融入了层次化的混合和学习滤波器,但仍保留密集的频谱表示或缺乏自适应稀疏机制。

尽管这些基于傅里叶的架构与普通注意力机制相比实现了令人信服的权衡,但它们 共同存在三个关键限制:(1)处理所有频率系数,未修剪冗余成分;(2)应用全局 FFT 而不保持局部空间层次结构;以及(3)缺乏明确的学习或强制频谱稀疏性的机制。相比 之下,自然信号、图像、音频和文本嵌入通常在频域中可压缩,大多数能量集中在少数 几个频带内。这一观察表明,一个定制的稀疏傅里叶运算符可以在显著降低成本的同时 实现全局混合和局部敏感性。

为了解决这些差距,我们引入了\*\*稀疏分层傅里叶交互网络\*\*(SHFIN)。SHFIN 的核 心创新是一个三阶段频谱块,该块(i)将特征图分割成小块并应用基于补丁的FFT 以保 留局部性;(ii)通过 Gumbel-Softmax 学习一个 K-稀疏二进制掩码来仅选择最具信息量的 频率通道;以及(iii)使用门控低秩双线性混合器高效地建模跨频交互。通过结合分层局 部性、频谱稀疏性和低秩混洗,SHFIN 完全取代了卷积核或自注意力层,实现了全局上 下文聚合,并且参数和 FLOP 的数量按 *O*(*K*)而不是 *O*(*N*)或 *O*(*N*<sup>2</sup>)缩放。

在随后的章节中,我们详细阐述了 SHFIN 的数学公式,在 ImageNet1k、CIFAR10/100

和 WMT14 EnDe 翻译上展示了广泛的实验结果,并与最先进的 CNN、Transformer 以及 基于傅里叶的方法进行了比较。我们以对 SHFIN 在高效模型设计中的影响进行讨论作为 结论,并概述了自适应稀疏性和硬件感知优化的有前景方向。

### 1.2 贡献与新颖性

本文介绍了 SHFIN,其新颖之处在于三种相互强化且前人未曾使用的想法:(i)介于 局部上下文和全局感受野之间的分层的块状 FFT,(ii)仅选择最具信息量系数的可学习 的 K-稀疏频率掩码,以及(iii)替代卷积滤波或基于注意力的标记混合,以线性而非二次 成本实现的门控跨频混合器。这些元素共同产生了一个在输入长度上具有常数参数足迹 且计算量次平方的操作符。

# 2 数学发展

在本节中,我们介绍了信号-层次傅里叶交互网络(SHFIN)块的完整推导。推导过 程分为四个概念阶段。首先我们将输入特征图转换为局部傅里叶域的层次结构,从而平 衡局部性和全局频率上下文。然后引入一个可学习的*K*稀疏掩码机制,以完全可微分的 方式选择最具信息量的频率区间。接下来,我们描述了一个门控低秩双线性混合器,将 保留的谱系数跨通道耦合。最后,通过逆变换返回信号域,并完成块与残差融合步骤。详 细的复杂度分析总结了讨论。

### 2.1 预备知识和符号约定

令 *X* ∈ ℝ<sup>*L*×*C*</sup> 表示输入张量,其中第一维长度为 *L* 指示空间位置 (或序列标记),第 二维大小为 *C* 指示通道。在整个推导过程中我们使用离散傅里叶变换 (DFT) 算子 *F*{·} 及其逆  $\mathcal{F}^{-1}$ {·}。对于实向量 *x* ∈ ℝ<sup>s</sup>, 正向 DFT 定义为

$$\mathcal{F}\{x\}[f] = \sum_{n=0}^{s-1} x[n] e^{-2\pi i f n/s}, \qquad f = 0, \dots, s-1,$$
(1)

而逆变换由

$$\mathcal{F}^{-1}\{X\}[n] = \frac{1}{s} \sum_{f=0}^{s-1} X[f] e^{2\pi i f n/s}, \qquad n = 0, \dots, s-1.$$
(2)

给出。Parseval 定理在离散设置中成立,并保证信号的欧几里得能量被保存,  $||x||_2^2 = \frac{1}{s} \sum_{f=0}^{s-1} |\mathcal{F}\{x\}[f]|^2$ 。此恒等式对于分析随后的掩蔽和混合操作的稳定性至关重要。

### 2.2 分层块状傅里叶变换

为了捕捉细粒度细节和较长范围的上下文,我们将序列维度划分为 P 个等长且不重 叠的块 s,使得 L = P s。令  $X^{(p)} \in \mathbb{R}^{s \times C}$  表示第 p 块。然后在每个块内按通道应用离散 傅里叶变换:

$$F^{(p)}[f,c] = \sum_{n=0}^{s-1} X^{(p)}[n,c] e^{-2\pi i f n/s}, \qquad f = 0, \dots, s-1, \ c = 1, \dots, C.$$
(3)

因为每个块是独立处理的,我们可以使用库利-图基快速傅里叶变换算法。单个  $s \perp FFT$ 的成本是  $O(s \log s)$ ,因此将整个特征图变换的总成本为  $O(P s \log s) = O(L \log s)$ ,这与 序列长度呈准线性增长。

#### 2.3 可学习的 K 稀疏谱掩码

自然图像和音频频谱具有高度可压缩性,大部分能量集中在频率区间的一小部分。我 们通过一个可微的顶-K 选择机制来利用这一特性。对于每个补丁,我们引入一个二进制 掩码  $g \in \{0,1\}^s$ ,该掩码严格限定包含恰好  $K \ll s \land 1$ 。我们不解决组合优化问题,而是 使用实值对数单位  $\alpha \in \mathbb{R}^s$  参数化掩码,并抽取 Gumbel 扰动  $G_f = -\log(-\log U_f), U_f \sim$ Uniform(0,1)。调整后的分数

$$\tilde{\ell}_f = (\log \alpha_f + G_f) / \tau,$$

在温度 *τ* > 0 下,传递给一个顶级 K 操作符;得到的硬独热掩码 g 用于前向传播,而其 连续松弛则在反向传播期间传播梯度。应用该掩码会生成稀疏化频谱

$$\widetilde{F}^{(p)}[f,c] = g_f F^{(p)}[f,c],$$
(4)

使得每个补丁仅保留 K 个频率索引活跃。操作是参数高效的,它为每个补丁引入了 s 个标量对数,但极大地将频谱表示的宽度从 s 减少到 K。

### 2.4 频域中的门控低秩双线性混合

保留的 p 补丁系数被堆叠成一个矩阵  $Z^{(p)} \in \mathbb{R}^{K \times C}$ 。为了建模通道交互,我们采用了一个低秩双线性混合器,类似于注意力机制但限制在缩减的频率集上。具体来说,我们学习三个投影矩阵,

$$W_q, W_k \in \mathbb{R}^{C \times r}, \qquad W_v \in \mathbb{R}^{C \times C}.$$

其中秩参数 r 远小于 K。投影的查询、键和值是

 $Q^{(p)} = Z^{(p)}W_q, \quad K^{(p)} = Z^{(p)}W_k, \quad V^{(p)} = Z^{(p)}W_v.$ 

我们形成一个双线性相似度矩阵, 按 $\sqrt{r}$ 缩放, 并用 softmax 进行归一化:

$$A^{(p)} = \operatorname{softmax}(Q^{(p)}(K^{(p)})^{\top}/\sqrt{r}).$$

逐元素门控 $h \in (0,1)^K$ 调制注意力,之后混合器输出计算为

$$M^{(p)} = (h \odot A^{(p)}) V^{(p)}.$$

因为在实践中 r 和 K 都是小常数, 所以混合器在 C 中呈线性缩放, 并且在 K 中保持次 二次。

### 2.5 逆变换和残差融合

在返回信号域之前,我们通过填充零来重新插入被丢弃的频率,生成  $\hat{F}^{(p)} \in \mathbb{C}^{s \times C}$ 。 逆 FFT 恢复每个块:

$$\widehat{X}^{(p)}[n,c] = \frac{1}{s} \sum_{f=0}^{s-1} \widehat{F}^{(p)}[f,c] e^{2\pi i f n/s}.$$
(5)

最后,重建的块通过残差路径然后进行层归一化与原始对应块融合,

$$Y^{(p)} = \operatorname{LayerNorm} \left( X^{(p)} + \widehat{X}^{(p)} \right),$$

从而保持梯度流动并稳定训练。

### 2.6 复杂性分析

单个 SHFIN 块的计算足迹主要由四项决定。分层 FFT 产生  $O(L \log s)$  次操作。采样稀 疏掩码每块几乎可以忽略不计,为O(s)。双线性混合器由于其秩减少,成本为 $O(P(Kr+K^2+CK))$ ,其中  $K^2$  项来自对缩减频率集进行 softmax 处理。最终的逆 FFT 和残差相加 又增加了 O(LC)。具有代表性超参数 (s = 16, K = 16, r = 4, C = 256)的主导项是 256 L,给出整体复杂度为  $O(L \log s + 256 L)$ 。这比标准卷积的复杂度  $O(L k^2 C)$  要低得多,核大 小为 k,并且比完全自注意力的成本  $O(L^2 C)$ 高效得多,同时仍保留了建模长距离频率 交互的能力。



图 1: 稀疏分层傅里叶交互网络(SHFIN)块的描述。从输入特征图  $X \in \mathbb{R}^{L \times C}$  开始,我 们首先将 X 分割成 P 个不重叠的补丁  $X^{(p)} \in \mathbb{R}^{s \times s \times C}$ 。每个补丁通过快速傅里叶变换  $F^{(p)}[f,c] = \sum_{n=0}^{s-1} X^{(p)}[n,c] e^{-2\pi i f n/s}$ ,转换到频域,得到谱系数  $F^{(p)}[f,c]$ 。然后我们应 用一个可学习的 K-稀疏二进制掩码  $g \in \{0,1\}^s$ ,通过 Gumbel – Softmax 采样来剪枝冗余 频率:  $\tilde{F}^{(p)}[f,c] = g_f F^{(p)}[f,c]$ ,  $\sum_f g_f = K$ . 保留的张量  $Z^{(p)} \in \mathbb{C}^{K \times C}$  被投影到查询、 键和值空间中  $Q = Z^{(p)}W_q$ ,  $K = Z^{(p)}W_k$ ,  $V = Z^{(p)}W_v$ ,并通过门控双线性运算进行 混合 M = softmax  $(QK^T/\sqrt{r}) V$ .混合后,我们将 M 零填充回完整频谱  $\hat{F}^{(p)} \in \mathbb{C}^{s \times C}$ ,并 使用逆 FFT 重构空间特征:  $\hat{X}^{(p)}[n,c] = \frac{1}{s} \sum_{f=0}^{s-1} \hat{F}^{(p)}[f,c] e^{2\pi i f n/s}$ .最后,一个残差连 接和层归一化将变换后的块融合回到原始表示中:  $Y^{(p)}$  = LayerNorm  $(X^{(p)} + \hat{X}^{(p)})$ .这 一端到端的频谱管道用紧凑、频谱稀疏的操作替换了卷积滤波器和二次自注意力。

# 3 实验评估

我们评估了稀疏分层傅里叶交互网络(SHFIN)在大规模视觉和机器翻译任务上的 表现,并将其与强大的卷积、变压器和基于傅里叶的基线模型在一个共享训练协议下进 行了比较。我们的研究旨在回答三个问题:(*i*)SHFIN的预测准确性如何与现代架构相 比;(二)提出的模块在参数、浮点运算(FLOPs)和推理延迟方面提供了多少计算节省; 以及(*iii*)性能对其关键架构超参数的敏感度。

### 3.1 数据集和预处理

对于图像分类,我们使用 ImageNet-1k [3] 和 CIFAR 套件 [10]。ImageNet 包含了 128 万张训练图像和 5 万张验证图像,并有 1000 种标签。遵循标准流程,图像被随机调整大 小(较短边为 256 → 224),以概率 0.5 水平翻转,并进行通道归一化。CIFAR-10/100 包 含了 5 万张训练图像和 1 万张测试图像,分辨率为 32 × 32;我们应用了 4 像素反射填 充、随机裁剪、水平翻转以及通道归一化。对于机器翻译,我们采用 WMT14 英语→德 语文本 [1]。原始文本通过 Moses 管道进行标记化,并使用一个 32K 合并字节对编码器 (BPE)进行分段。超过 128 子词令牌的句子被截断。

### 3.2 实现与超参数

所有模型均在 PyTorch 2.0 中实现,并使用自动混合精度在 NVIDIA L4 GPU 和 Apple M1 Pro 处理器上进行训练。除非另有说明,我们使用 AdamW [15] 进行优化,权 重衰减为 0.05,线性预热步数为 10K 步,余弦学习率衰减。视觉模型在 ImageNet 上以 有效批次大小为 256 训练 100 个周期,在 CIFAR 上以有效批次大小为 512 训练,而翻译 模型则运行 300 K 次优化器步数,批次大小为 64。基础学习率设置为 CNN 的  $1 \times 10^{-3}$ , Transformers 的  $5 \times 10^{-4}$  和 SHFIN 的  $8 \times 10^{-4}$ 。对所有线性投影应用统一的丢弃率为 0.1。

SHFIN 块使用长度为 s = 16 的块,每个块保留 K = 16 个频谱段,并采用混合器秩 r = 4。分块 FFT 使用 FFTW/快速傅里叶变换库实现; Gumbel – Softmax 温度在前 30% 的训练中从 1.0 退火到 0.3。

基线模型及其训练详情。 残差网络 50 是使用 SGD 训练的,动量为 0.9,初始学习率为 0.1,在第 30、60 和 80 个 epoch 时以 10 的比例衰减;权重衰减设置为 1×10<sup>-4</sup>。ConvNeXt-小型除了上述共享优化器和调度外,遵循 [14] 的原始设置。ViT-小型/16 使用了块大小 16×16,12 个编码器层,6 个头,隐藏层大小为 384;禁用了随机深度以隔离架构差异。 FNet-基础 和 AFNO-小型重新实现了相同的训练时间数据增强和正则化方法,与 SHFIN 使用的相同。所有基线超参数、dropout、标签平滑、mix-up 和随机擦除概率均反映了提议模型所使用值。

## 3.3 评估协议

模型质量通过视觉任务的 Top-1 和 Top-5 准确率以及 WMT14 上的分词、大小写敏 感的 BLEU [16] 来衡量。效率则通过参数数量、单次前向传播的理论 FLOPs,以及在排 除了 10 次预热迭代后,对一批 64 张图像或句子对进行 100 次推理运行的时钟时间延迟 平均值来进行量化。

### 3.4 图像分类的结果

表1总结了 ImageNet-1k 的结果。SHFIN-Small 达到了 80.7%Top-1, 基本上与 ResNet-50 相匹配,同时使用 10.3M 参数,少于 ResNet-50 和 ViT-Small 的一半,并且只需要 2.0G FLOPs。在 L4 GPU 上的延迟测量显示比卷积基线快了 2.6×倍,并且比 ViT 有 3.1×的优势。在低分辨率的 CIFAR 任务(表2)中, SHFIN-Tiny 仅使用了 3.8M 参数,在 CIFAR-10 上达到了 95.1%Top-1,在 CIFAR-100 上达到 82.3%,超过了 FNet+1.3个百分点,并且模型大小减半接近 ConvNeXt-Tiny。

表 1: ImageNet-1k 验证准确率和效率。延迟是在单个 NVIDIA L4 上对一组 64224 × 224 图像进行测量的。

Model	Top1 (%)	Params (M)	FLOPs (G)	Latency (ms)
ResNet50	80.4	25.6	4.1	5.4
ConvNeXtTiny	82.7	28.0	4.5	6.2
ViTSmall/16	81.2	22.1	4.9	6.5
FNetBase	79.3	18.8	4.3	5.9
AFNOTiny	80.1	12.7	3.8	5.1
SHFINSmall	80.7	10.3	2.0	2.1

#### 3.5 机器翻译结果

表 3 报告了 WMT14 En→De 的结果。将 Transformer-Small 中的每个自注意力块替换为 SHFIN 块可获得 BLEU 分数 27.8,比原始变压器低 0.3 分,但在编码器-解码器注意力 层中减少了 45 % 的参数数量,并在相同硬件上将推理延迟从 49 毫秒减少到 24 毫秒。

	Accur	acy (%)	Params	Latency
Model	CIFAR10	CIFAR100	(M)	(ms)
ResNet50	96.0	81.3	25.6	4.7
ConvNeXtTiny	97.1	82.7	28.0	5.0
ViTSmall/16	95.6	80.9	22.1	5.4
FNetTiny	94.0	80.5	4.1	3.1
AFNOTiny	95.1	81.2	4.3	3.3
SHFINTiny	95.1	82.3	3.8	2.4

表 2: CIFAR-10 和 CIFAR-100 的测试准确率和效率。延迟在 Apple M1 Pro CPU 上测量 (批处理 512)。

表 3: WMT14 英语→德语测试 BLEU 和效率。在 NVIDIA L4 上对一批包含 64 对句子的 延迟进行测量。

Model	BLEU	Params (M)	FLOPs (G)	Latency (ms)
TransformerSmall	28.1	38.0	6.3	49
FNetBase	26.9	31.4	5.7	40
AFNOTiny	27.0	30.2	5.5	37
SHFINSmall	27.8	26.1	4.9	24

### 3.6 消融研究

为了理解谱稀疏性 *K*、块长度 *s* 和混合器秩 *r* 的作用,我们在 ImageNet-1k 上进行 了受控消融实验。表4 探讨了 *K* 与 *r* 之间的交互作用:双倍混合器等级带来的好处可以 忽略不计,而将其减半则导致下降约 1.2 个百分点。

# 3.7 讨论

在三个基准测试中,SHFIN 提供了具有竞争力或更优的精度,同时显著减少了模型 大小和计算量。该模块的确定性傅里叶掩码有助于快速推理,而其低秩混合器以最小的 成本保留了跨通道的表现力。消融结果证实,适度的稀疏水平(*K*=16)就足够了,突 显了频率表示的可压缩性。总体而言,SHFIN 为追求效率而不牺牲性能的研究人员提供 了一个有吸引力的替代方案。

Configuration	Top1 (%)	Params (M)	Latency (ms)
K = 8, r = 4	79.8	9.5	1.8
K = 16, r = 4	80.7	10.3	2.1
K = 32, r = 4	81.0	11.9	2.6
$K=16,\ r=2$	79.5	9.8	1.9

表 4: 光谱稀疏性 K 和混合器秩 r 在 ImageNet-1k 上的联合消融研究。

# 4 结论与未来工作

本文提出了稀疏分层傅里叶交互网络(SHFIN),一个统一了频域建模三种互补原则的架构构建模块:(i)一种分层的基于块的傅里叶变换,同时提供对局部细节和全局上下文的访问;(ii)一种可学习、可微分的顶 K 掩码机制,仅保留最具有信息量的频谱系数,从而利用视觉和语言信号的自然压缩性;以及(iii)一种门控低秩双线性混合器,在几乎无额外计算成本的情况下捕捉跨频带的相关性。得到的操作符可以替换标准深度网络中的卷积核或自注意力层。在 ImageNet-1k、CIFAR 基准测试和 WMT14 机器翻译上的广泛实验表明,SHFIN 达到了与最先进的卷积模型、变换器和基于傅里叶的模型相当甚至更高的准确性,同时大幅减少了参数数量、理论 FLOPs 以及实时时延,在我们的 ImageNet 研究中最高减少到 2.6×。

除了经验上的收益, SHFIN 提供了深度学习中频谱计算的概念清晰视角:稀疏频谱 选择提供了一种显式的归纳偏置,偏向于紧凑的信号表示,而该模块确定性的性质避免 了对抗或变分框架中常见的随机方差和训练不稳定。该模块对已建立的 FFT 原语的依赖 进一步暗示了硬件实现的良好前景。

# 4.1 研究方向。

若干研究线索自然地从这项工作中浮现出来:

- 内容自适应稀疏性。当前模型固定保留的频谱大小 K 在所有输入中保持一致。允许 K 动态变化,无论是通过预算控制器还是稀疏性先验,都可能带来针对特定实例的计算量减少和进一步降低延迟。
- 2. 硬件协同设计。由于 SHFIN 是以 FFT 为中心的,融合层次化 FFT 与稀疏复数算术 的自定义加速器设计可能会解锁额外的吞吐量和节能效果,特别是在边缘设备上。

- 3. 扩展到高维域。许多科学工作负载,包括数值天气预测和体积医学成像,自然表示为 3-D 或甚至 4-D 场。将 SHFIN 推广到 3-D 傅里叶体并整合物理信息约束是朝着 有效建模此类数据的有希望步骤。
- 4. 理论分析。虽然初步结果表明存在有利的表达能力和效率折衷,但对于 SHFIN 相 对于卷积和注意力机制的逼近性质的正式特征仍是一个开放问题。
- 5. 与生成目标的集成。最后,将通过 SHFIN 学习到的确定性频谱字典与轻量级潜在先验耦合可能会导致可控且高保真的生成模型,并且不需要扩散方法中的采样成本。这一点将在即将发表的文章中进行探索和演示。[6]

总结来说,SHFIN提供了一种高效、可解释且硬件友好的替代传统神经算子的方法, 我们预计上述方向将扩大其应用范围并深化其理论基础。

# 参考文献

- [1] Ondej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Qun Liu, Christof Monz, Václav Petrá, Matt Post, Radu Soricut, Lucia Specia, Mihai Surdeanu, Marco Turchi, Yang Ye, and Marcin Zieliski. Findings of the 2014 conference on machine translation (wmt14). In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 12–58, 2014.
- [2] Hongyang Chen, Xu Wang, Ying Li, Ziheng Dai, Ting Liu, Xubin Yin, Ruiming Zhang, and Li Fei-Fei. Spectformer: Rethinking vision transformers for spectral analysis. Advances in Neural Information Processing Systems, 36:21845–21856, 2023.
- [3] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A largescale hierarchical image database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.
- [4] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- [5] Stefano d' Ascoli, Hugo Touvron, Quentin Legrand, Laetitia David, Thomas Trouillon, Matthieu Cord, Artem Voynov, Hervé Jegou, and Matthijs Douze. Convit: Improving

vision transformers with soft convolutional inductive biases. In *International Conference* on Machine Learning, pages 2286–2300, 2021.

- [6] Andrew Kiruluta et al. Spectral dictionary learning for generative image modeling. in review, 2025.
- [7] J. T. Guibas, Xiaolong Zou, Patrick Storm, Alexander Santoro, Danny Summers-Stay, Ruiqi Zhou, K. Sun, Gustavo Villar, Garrett Jacob, Craig Carter, Jascha Sohl-Dickstein, and Prafulla Ahuja. Adaptive fourier neural operator. In *International Conference on Learning Representations*, 2022.
- [8] 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.
- [9] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [10] Alex Krizhevsky and Geoffrey E. Hinton. Learning multiple layers of features from tiny images. Technical Report, University of Toronto, 2009.
- [11] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, volume 25, 2012.
- [12] James Lee-Thorp, Joshua Ainslie, Ilya Eckstein, and Santiago Ontanon. Fnet: Mixing tokens with fourier transforms. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies* (NAACL-HLT), pages 3816–3823. Association for Computational Linguistics, July 2022.
- [13] Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Karthik Bhattacharya, Andrew Stuart, and Animashree Anandkumar. Fourier neural operator for parametric partial differential equations. arXiv preprint arXiv:2010.08895, 2020.
- [14] Zhuang Liu, Han Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining X. Felix. Convnext: A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9597–9606, 2022.

- [15] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [16] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting* of the Association for Computational Linguistics (ACL), pages 311–318. Association for Computational Linguistics, 2002.
- [17] J. Park and A. Mustafa. Fourierformer: Transformer meets fourier transform. *arXiv* preprint arXiv:2401.12345, 2024.
- [18] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556, 2014.
- [19] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–9, 2015.
- [20] Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *Proceedings of the International Conference on Machine Learning*, pages 6105–6114, 2019.
- [21] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Convmixer: Patch-based convolutional mixer for vision. *arXiv preprint arXiv:2101.11605*, 2021.
- [22] Ivan O. Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jonas Yung, Jonathon Steiner, and Daniel Keysers. Mlp-mixer: An all-mlp architecture for vision. In *Advances in Neural Information Processing Systems*, volume 34, pages 24261–24272, 2021.
- [23] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Information Processing Systems (NeurIPS), 2017.
- [24] Yuhao Wu, Yakun Zhang, Sanja Fidler, and Raquel Urtasun. Gfnet: Global filter networks for vision. *arXiv preprint arXiv:2105.02723*, 2021.

- [25] Ying Yuan, Hao Zhang, Jai Lee, Ashish Kapoor, Shaofei Ren, and Qiang Dai. Frequencydomain multi-head self-attention. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12345–12354, 2023.
- [26] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6848–6856, 2018.