



用于高效视觉任务的多用卷积核

(Minor revision in T-PAMI)

韩凯 吴恩华

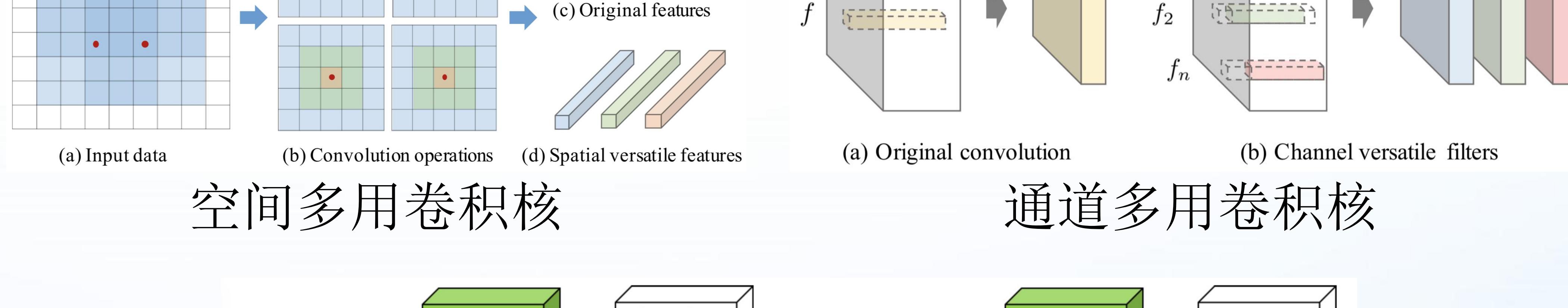
摘要

卷积神经网络模型在计算机视觉、图形学等领域应用广泛，但是卷积网络体积大速度慢。本文提出用于加速神经网络的多用卷积核，利用少量卷积核生成更多卷积核。所提出多用卷积核能够大幅降低卷积网络的计算量和参数量。

方法

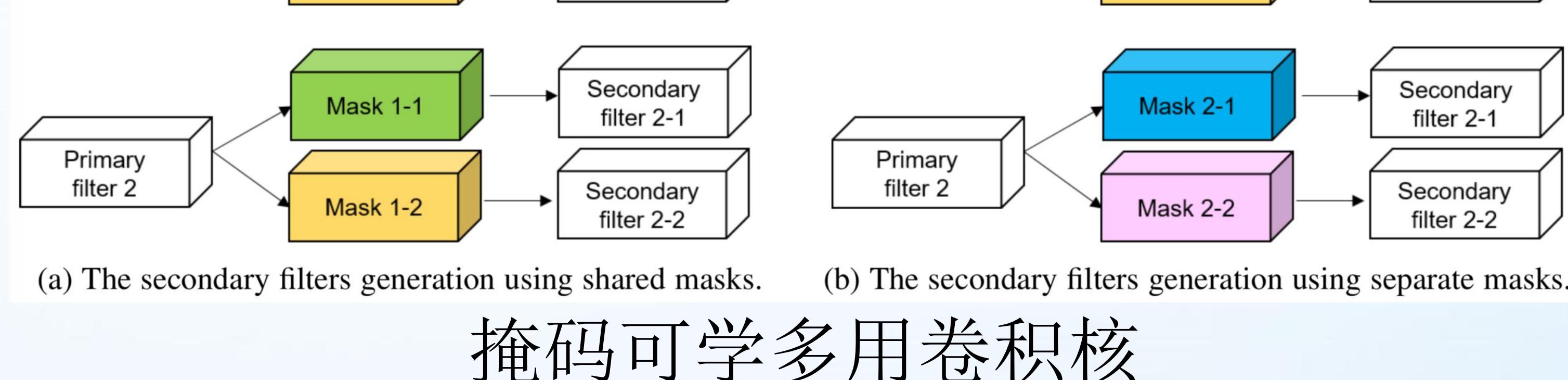
建立少量主卷积核，从主卷积核中产生更多子卷积核。

1. 空间多用卷积核：从高度、宽度维度“挖”出子卷积核。
2. 通道多用卷积核：从通道维度“挖”出子卷积核。
3. 掩码可学多用卷积核：不限制维度，让二值掩码是可学习的，和权重一起训练。



空间多用卷积核

通道多用卷积核



掩码可学多用卷积核

实验

ImageNet图像分类模型压缩

ResNet-50	#Param ($\times 10^7$)	Mem (MB)	#MUL ($\times 10^9$)	Top1err	Top5err
Vanilla [20]	2.6	97	4.1	24.7%	7.8%
S-Versatile	1.9	76	3.0	24.5%	7.6%
S+C-Versatile	1.1	42	1.5	25.5%	8.2%
Shared L-Versatile ($s=4$)	0.8	30	1.1	26.8%	8.9%
Separate L-Versatile ($s=4$)	0.9	33	1.1	25.5%	8.0%
Separate L-Versatile ($s=32$)	0.36	14	0.26	26.5%	8.7%

Top-1 Error (%) vs #MUL (10^9)

Legend: ResNet-50 (black circle), Ours (red circle), Winograd (yellow star), HRank (blue triangle), SSS (green diamond), ThiNet (purple cross), Slimmable (yellow triangle), GAL (cyan diamond), MetaPruning (pink circle)

Approximate data points from the graph:

#MUL (10^9)	ResNet-50	Ours	Winograd	HRank	SSS	ThiNet	Slimmable	GAL	MetaPruning
0.4	35.0	26.0	35.0	30.0	30.0	30.0	30.0	30.0	30.0
1.0	27.0	25.5	27.0	29.0	29.0	29.0	29.0	29.0	29.0
1.5	27.0	25.0	27.0	28.0	28.0	28.0	28.0	28.0	28.0
2.0	24.5	24.5	24.5	25.5	25.5	25.5	25.5	25.5	25.5
3.0	24.0	24.0	24.0	24.0	24.0	24.0	24.0	24.0	24.0
4.0	24.0	24.0	24.0	24.0	24.0	24.0	24.0	24.0	24.0

COCO目标检测模型压缩

