

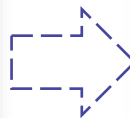


# 卷积神经网络压缩

吴建鑫

南京大学

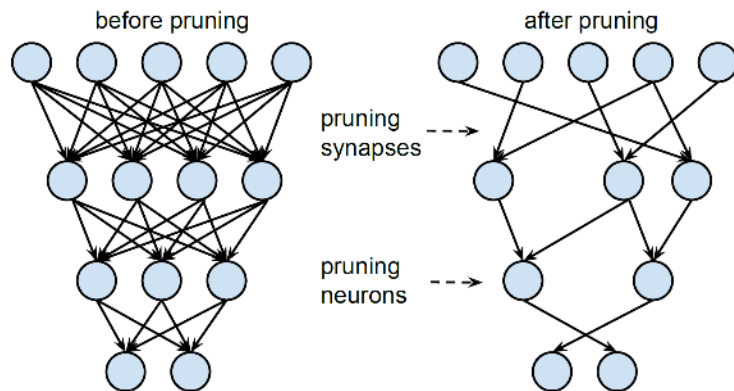
- 深度神经网络会消耗大量的计算资源



Q: 如何在资源受限的小型设备上运行深度网络?

A: 轻量级网络设计、参数量化、网络压缩 ...

- 深度学习时代的网络剪枝



[Han et al., NeurIPS 2015]

准则：所有权重低于某阈值的连接都会被移除

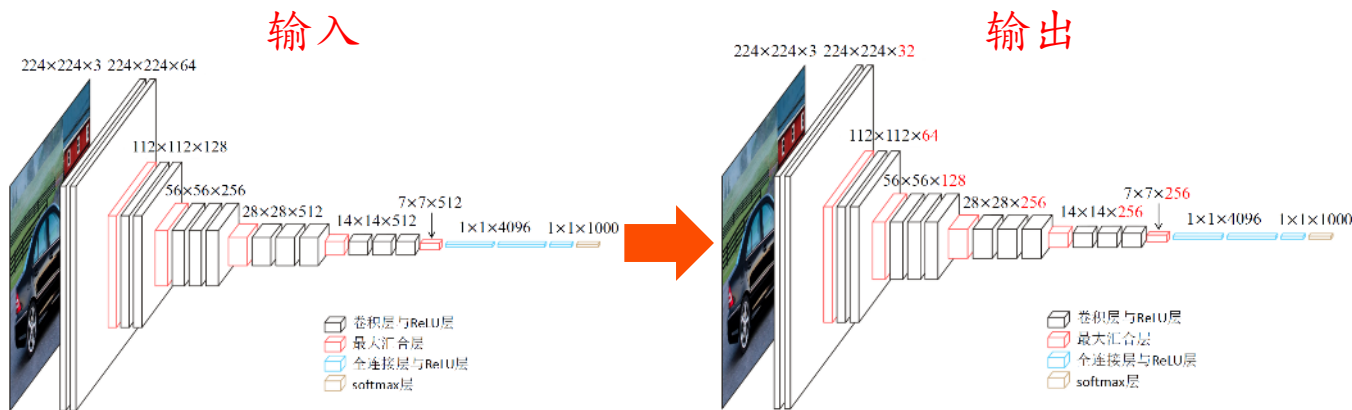
缺点：非结构化的模型很难获得实际加速效果

# 问题1：怎么剪？

---

- 简单答案：结构化剪枝！

## • 滤波器级别的深度网络剪枝



## • 结构化剪枝的优点

- ✓ 高实用性：可同时减少**计算量与内存占用**
- ✓ 低依赖性：结构化的剪枝**不会破坏模型结构**，无缝部署
- ✓ 可扩展性：模型剪枝可与**其他网络压缩技术**相结合

## 问题2：剪哪些？

---

- 简单答案：数学优化 vs. 启发式！

## • 启发式剪枝算法的缺点

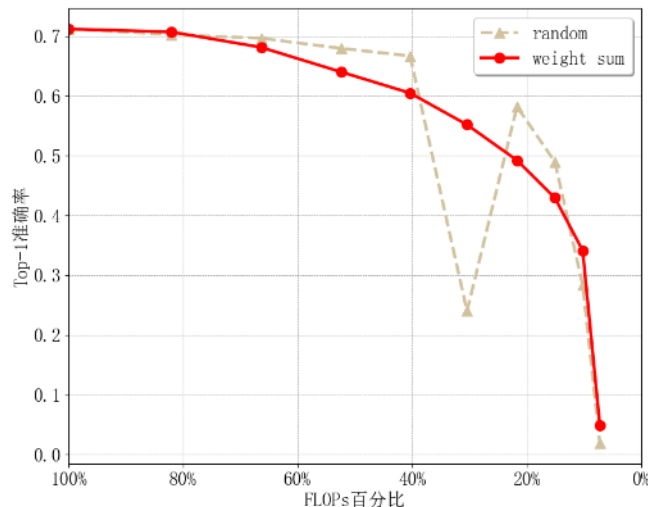
### □ 基线方法

- weight sum [ICLR'17]: 将每个滤波器权重的绝对值作为其重要性得分
- random: 随机丢弃一些滤波器

### □ 缺点

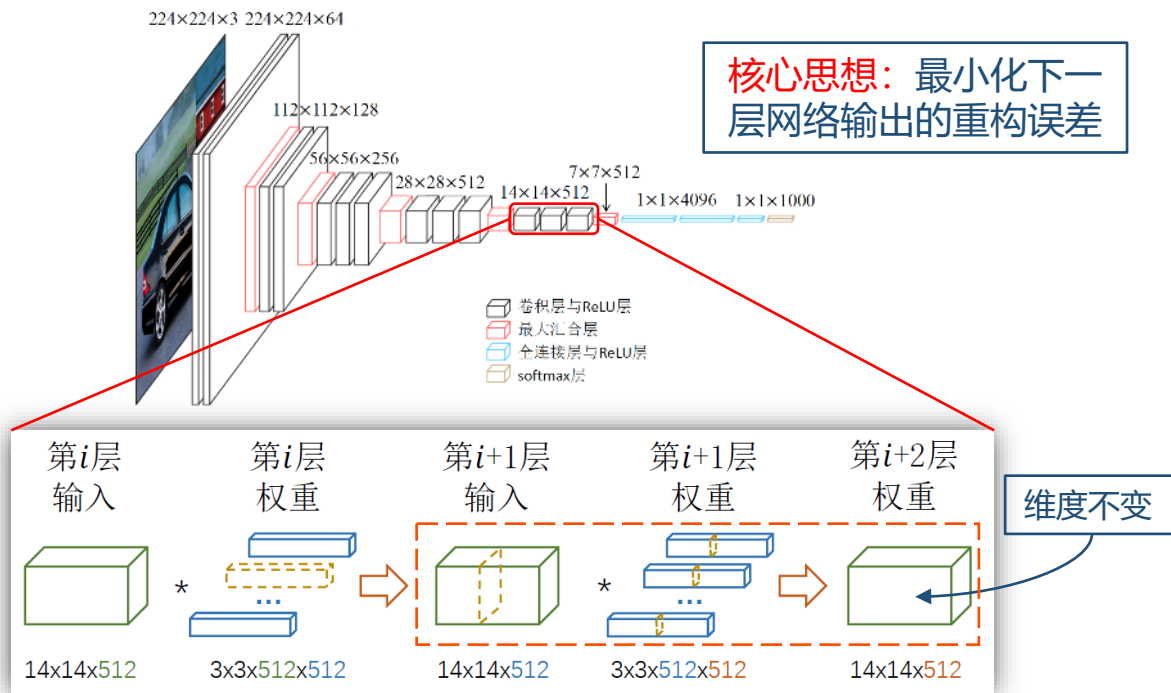
- 并不总是有效, 找不到成功或失败的原因
- 缺乏理论解释

启发: 能否对剪枝进行形式化建模,  
通过数学推导给出一定的保障?

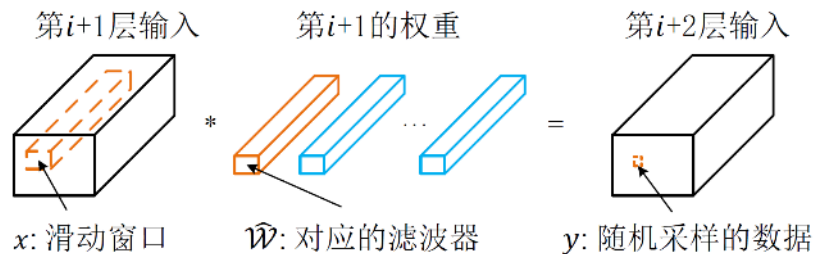


# 方法介绍

## • ThiNet



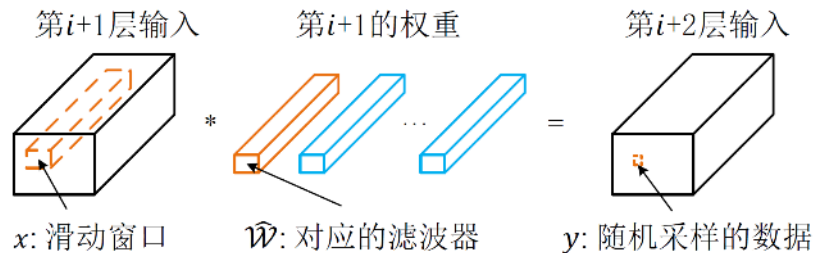




正常卷积: 
$$y = \sum_{c=1}^C \sum_{k_1=1}^K \sum_{k_2=1}^K \hat{\mathcal{W}}_{c,k_1,k_2} \times x_{c,k_1,k_2} + b.$$

定义: 
$$\hat{x}_c = \sum_{k_1=1}^K \sum_{k_2=1}^K \hat{\mathcal{W}}_{c,k_1,k_2} \times x_{c,k_1,k_2}, \quad \hat{y} = y - b.$$

代入可得: 
$$\hat{y} = \sum_{c=1}^C \hat{x}_c, \quad \text{其中, } C \text{ 是前一层滤波器的数量}$$



**剪枝:** 在 $C$ 个通道里选择一个子集 $S$ , 并移除那些不属于 $S$ 的滤波器

用数学语言来描述

**目标:** 找到子集  $S \subset \{1, 2, \dots, C\}$  使得下式尽可能成立

$$\hat{y} = \sum_{c \in S} \hat{x}_c$$

**方法:** 收集 $m$ 个训练样本 $\{(\hat{x}_i, \hat{y}_i)\}$ , 最小化 $\hat{y}$ 的重构误差

## 1. 转化为优化目标:

$$\arg \min_S \sum_{i=1}^m \left( \mathbf{y}_i - \sum_{j \in S} \mathbf{X}_{i,j} \right)^2$$

s.t.  $|S| = C \times r, \quad S \subset \{1, 2, \dots, C\}$

## 2. 通道放缩:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{i=1}^m (\mathbf{y}_i - \widehat{\mathbf{X}}_i \mathbf{w})^2$$

$$\hat{\mathbf{w}} = (\widehat{\mathbf{X}}^T \widehat{\mathbf{X}})^{-1} \widehat{\mathbf{X}}^T \mathbf{y}$$

## 3. 最终的优化目标:

$$\arg \min_S \sum_{i=1}^m \left( \mathbf{y}_i - \sum_{j \in S} \hat{\mathbf{w}}_j \mathbf{X}_{i,j} \right)^2$$

s.t.  $|S| = C \times r, \quad S \subset \{1, 2, \dots, C\}$

- 在数据集上的剪枝

- 基线方法:

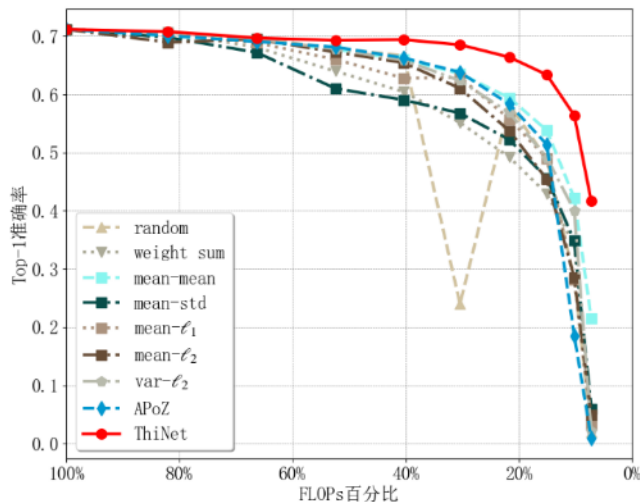
- random: 随机丢弃一些滤波器
- weight sum: 将每个滤波器权重的绝对值作为其重要性得分
- APoZ: 计算输出张量中每个通道的稀疏度作为其重要性得分
- mean-mean:  $s_i = \frac{1}{N} \sum \text{mean}(I(i, :, :))$
- mean-std:  $s_i = \frac{1}{N} \sum \text{std}(I(i, :, :))$
- mean- $\ell_1$ :  $s_i = \frac{1}{N} \sum \|I(i, :, :)\|_1$
- mean- $\ell_2$ :  $s_i = \frac{1}{N} \sum \|I(i, :, :)\|_2$
- var- $\ell_2$ :  $s_i = \frac{1}{N} \sum \text{var}(\|I(i, :, :)\|_2)$

- 数据集:

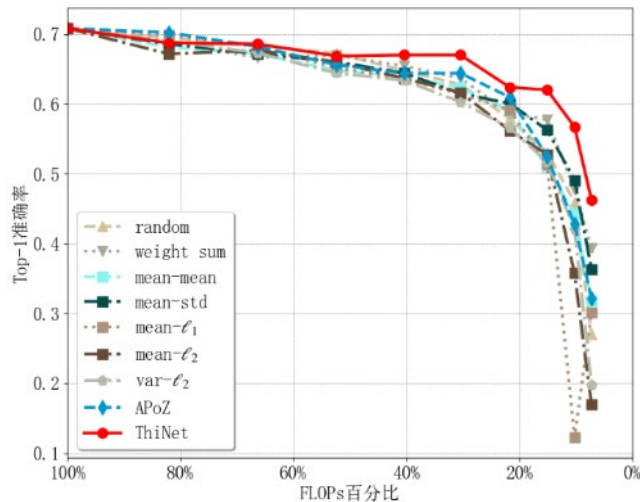
- CUB200: 鸟类数据集, 200类, 5994张训练图像, 5794张测试图像
- Indoor67: 室内场景数据集, 67类, 5360 张训练图像, 1340 张测试图像

- 在小数据集上的剪枝

CUB200数据集



Indoor67数据集



ThiNet算法在不同压缩率指标下均取得了良好的压缩效果

- 在ImageNet数据集上的剪枝 – VGG

算法	Top-1	Top-5	# 参数.	#FLOPs
VGG-16 <sup>[8] 1</sup>	71.50%	90.01%	138.34M	30.94B
连接级剪枝 <sup>[20]</sup>	68.66%	89.12%	10.30M	6.50B
APoZ-1 <sup>[25]</sup>	66.20%	87.60%	67.81M	–
APoZ-2 <sup>[25]</sup>	70.17%	89.69%	51.24M	–
Taylor-1 <sup>[26]</sup>	–	87.00%	–	11.50B
Taylor-2 <sup>[26]</sup>	–	84.50%	–	8.00B
weight sum <sup>[24] 2</sup>	69.35%	89.13%	131.44M	9.58B
Channel-Pruning (5×) <sup>[27] 3</sup>	67.80%	88.10%	130.87M	7.03B
从头训练 <sup>4</sup>	67.00%	87.45%	131.44M	9.58B
ThiNet-Conv-1	69.74%	89.41%	131.44M	9.58B
ThiNet-Conv-2	69.11%	88.86%	130.50M	6.91B
ThiNet-GAP	67.69%	88.13%	8.32M	9.34B

- 在ImageNet数据集上的剪枝 – ResNet

模型	Top-1	Top-5	# 参数	#FLOPs	速度 (图片/秒)
原始模型 <sup>1</sup>	75.30%	92.20%	25.56M	7.72B	295.12
ThiNet-70	74.03%	92.11%	16.94M	4.88B	334.87
ThiNet-50	72.03%	90.99%	12.38M	3.41B	373.92
ThiNet-30	68.17%	88.86%	8.66M	2.20B	397.86

<sup>1</sup> ResNet 的精度是按照官方中心裁剪的设定进行测试的：原图保持长宽比短边缩放到 256，之后进行 224 × 224 中心裁剪 (<https://github.com/KaimingHe/deep-residual-networks>)。

- 针对以往剪枝算法多为启发式，缺少理论保障的缺点，对剪枝过程进行了形式化建模，提出了一种基于重构误差最小化的剪枝算法：ThiNet
- Jian-Hao Luo, Jianxin Wu and Weiyao Lin. ThiNet: A filter level pruning method for deep neural network compression. **In Proceedings of the IEEE International Conference on Computer Vision (ICCV)**, Venice, Italy, October 2017, pp. 3505-3513.
- Jian-Hao Luo, Hao Zhang, Hong-Yu Zhou, Chen-Wei Xie, Jianxin Wu and Weiyao Lin. ThiNet: Pruning CNN Filters for a Thinner Net. **IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)**, 41(10), 2019: 2525-2538.

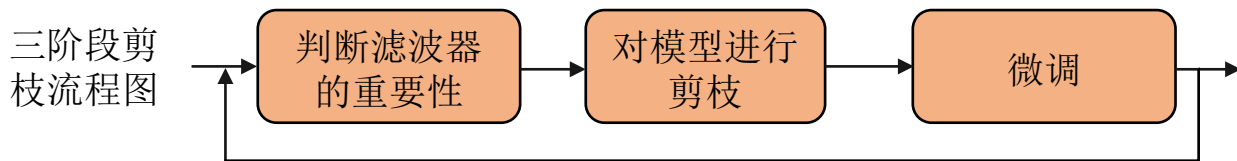


# 问题3：能不能自动化？

---

- 简单答案：能，且应该自动化！

- 三阶段剪枝的缺点

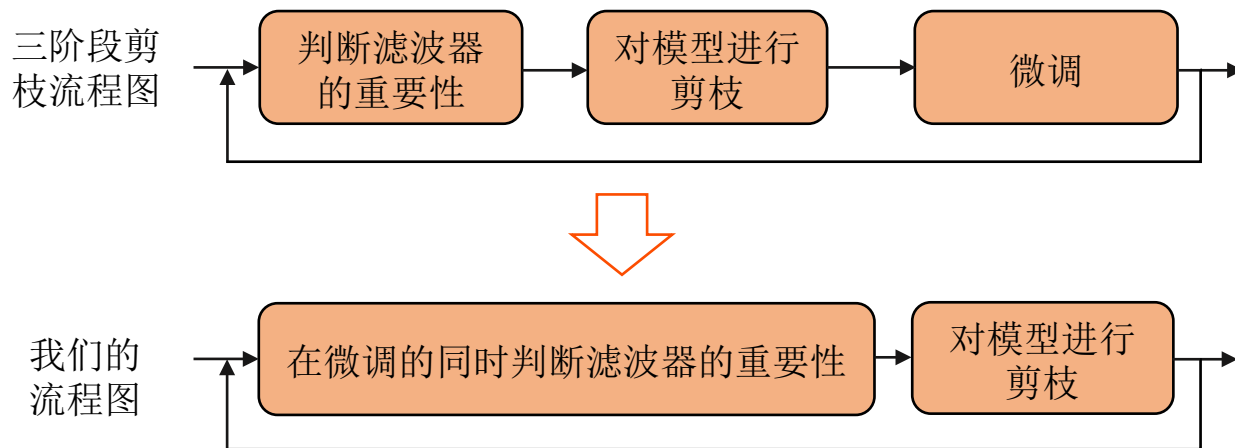


□ 三阶段剪枝策略：

- 衡量重要性——移除弱滤波器——微调恢复精度
- 过于依赖人工手动设计的重要性衡量标准
- 重要性衡量标准无法适用于所有的情况
- 每个层的压缩比都相等

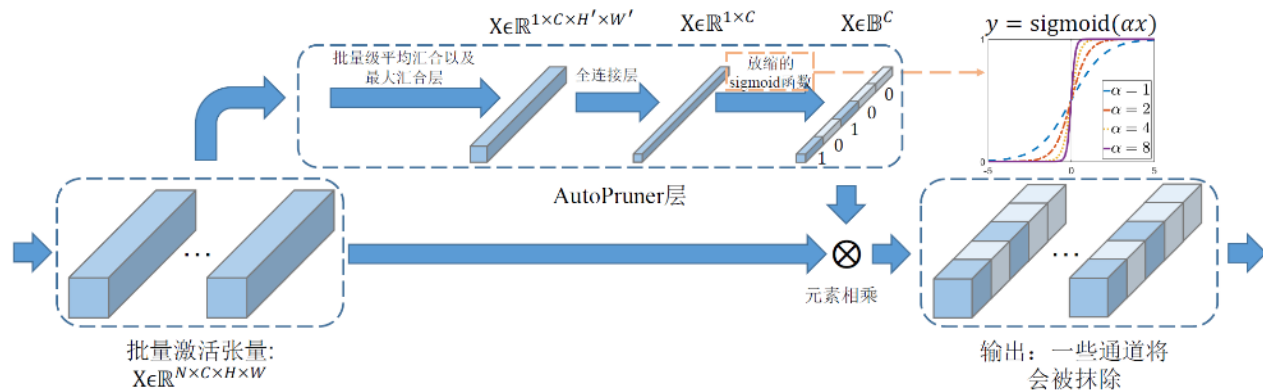
启发：能否让网络自己来决定哪些滤波器需要被剪枝？

- 新的剪枝流程



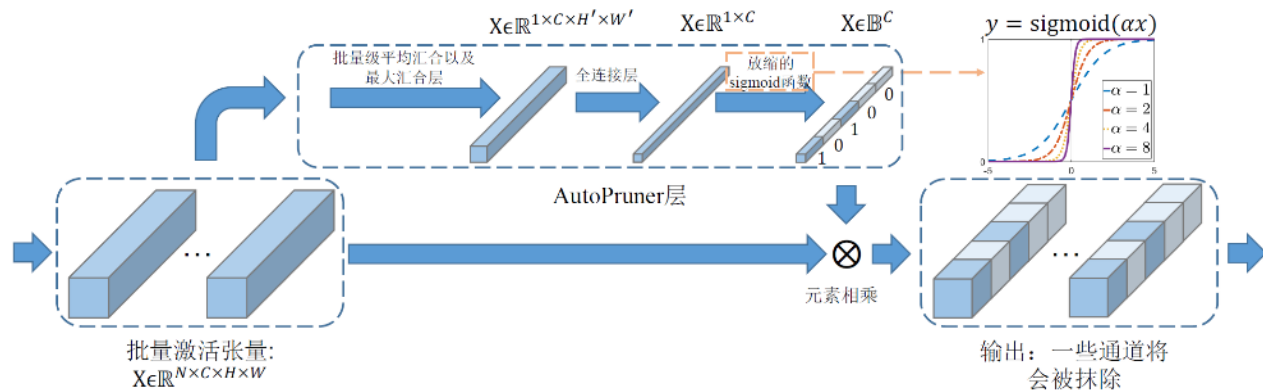
解决方案：设计新的网络层，通过**端到端**的方式**自动**完成剪枝

## • AutoPruner算法框架



- ❑ 网络结构：将AutoPruner层添加到卷积-激活函数之后，如，  
 $\text{conv1-1} - \text{ReLU} - \text{AutoPruner} - \text{conv1-2} \dots$
- ❑ 输入：批大小为N的4阶张量
- ❑ 输出：C维二值向量，0元素所对应的通道和滤波器将被移除

## • AutoPruner算法框架



### □ 具体实现:

1. 汇合 (pooling): 利用**批均值汇合**和**最大汇合**, 综合不同图像之间的信息
2. 编码: 使用**全连接层**得到C维向量
3. 二值化: 使用**放缩sigmoid**函数进行逐步二值化

- 稀疏控制与损失函数

□ 损失函数由分类损失与压缩率控制两部分组成：

$$\mathcal{L} = \mathcal{L}_{\text{classification}} + \lambda \left\| \frac{\|\mathbf{v}\|_1}{C} - r \right\|_2^2$$

- $\mathcal{L}_{\text{classification}}$ : 标准的分类损失项, 如cross-entropy;
- $\mathbf{v}$ :  $\mathbf{v} = \text{sigmoid}(\alpha x)$ 是AutoPruner层的输出, 由实值逐渐变为二值;
- $r$ : 预设的压缩率;

□ 如何控制 $\alpha$ 的值?

- $\alpha$ 的值过大:  $\mathbf{v}$ 的值将很快变为二值, 造成剪枝提前完成。此时, AutoPruner将退化为随机选择;
- $\alpha$ 的值过小:  $\mathbf{v}$ 有可能很难甚至根本无法收敛到二值;
- $\alpha$ 在不同网络模型、不同网络层的合适值可能不大相同;

解决方案: 提出了一种自适应的二值化控制方法

- 在CUB200数据集上的剪枝效果对比

方法	压缩率 $r = 0.2$			压缩率 $r = 0.5$		
	top-1 (%)	top-5 (%)	#FLOPs	top-1 (%)	top-5 (%)	#FLOPs
微调后的 VGG16	76.68	94.06	30.93B	76.68	94.06	30.93B
随机选择	57.28	83.52	2.62B	70.25	91.16	9.63B
ThiNet	63.12	87.54	2.62B	73.00	92.27	9.63B
AutoPruner	<b>65.06</b>	<b>87.93</b>	2.62B	<b>73.45</b>	<b>92.56</b>	9.63B

AutoPruner超过了前面介绍的ThiNet算法

- 在ImageNet上剪枝ResNet50模型

方法	Top-1	Top-5	#FLOPs	加速比 <sup>1</sup>
原始 ResNet-50 模型 <sup>2</sup>	76.15%	92.87%	8.18B	1.00×
<b>AutoPruner (<math>r = 0.3</math>)</b>	<b>73.05%</b>	<b>91.25%</b>	2.78B	2.94×
Taylor-FO-BN-56% <sup>[78]</sup>	71.69%	90.59%	2.68B	2.88×
GAL-1-joint <sup>[31]</sup>	69.31%	89.12%	2.22B	3.48×
ThiNet-30	68.42%	88.30%	2.20B	3.51×
<b>AutoPruner (<math>r = 0.5</math>)</b>	<b>74.76%</b>	<b>92.15%</b>	4.00B	2.05×
AutoPruner with block pruning ( $r = 0.5$ )	73.84%	91.75%	4.56B	1.79×
Channel Pruning (2×) <sup>[27]4</sup>	72.30%	90.80%	5.22B	1.48×
SSS (ResNet-26) <sup>[30]</sup>	71.82%	90.79%	4.00B	1.93×
ThiNet-50	71.01%	90.02%	3.41B	2.27×



- 对轻量级网络进行剪枝

算法	精度		MACs	速度	训练时间		
	top-1	top-5			剪枝阶段	微调阶段	总时间
MobileNetV2-1.0 <sup>1</sup>	72.19%	90.53%	300.79M	0.138s	-	-	-
MobileNetV2-0.75	69.95%	88.99%	209.08M	0.119s	-	-	-
AMC (70% FLOPs) <sup>[82]</sup>	70.85%	89.91%	210.63M	0.105s	22.52s×800	1687s×150	3.14d
Slimmable NN <sup>[33]</sup>	68.93%	88.34%	209.08M	0.120s	-	8618s×480	47.88d
AutoPruner	<b>71.18%</b>	<b>89.95%</b>	207.93M	0.099s	4898s×5	1658s×150	3.16d

<sup>1</sup> <https://github.com/d-li14/mobilenetv2.pytorch>

AutoPruner同样适用于MobileNet的剪枝，精度更高、速度更快、所消耗的计算代价比较适中

- 剪枝后网络的泛化性

算法	MobileNetV2-1.0	MobileNetV2-0.75	AMC <sup>[82]</sup>	Slimmable NN <sup>[33]</sup>	AutoPruner
mAP	0.6743	0.6529	0.6571	0.6602	0.6576

剪枝对模型的泛化性没有显著影响，但优势也并不明显，如何对分类以外的模型进行剪枝仍然值得进一步研究

- 针对以往剪枝算法过于依赖人工经验的缺点，我们在本章提出了一种可端到端训练的自动剪枝算法AutoPruner

- Jian-Hao Luo and Jianxin Wu. AutoPruner: An end-to-end trainable filter pruning method for efficient deep model inference. **Pattern Recognition (PR)**, 2020.

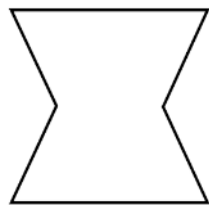
# 问题4, 5: 走向实用化

---

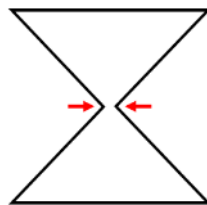
- 问题4: 小数据怎么办?
- 简单答案: 数据增强与知识蒸馏!
- 问题5: 残差结构怎么办?
- 简单答案: 残差内外一视同仁!

- 如何剪枝残差连接？

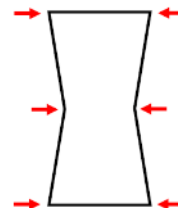
□ ThiNet、AutoPruner等算法均只剪枝残差结构的中间层



(a) 瓶颈结构



(b) 沙漏结构



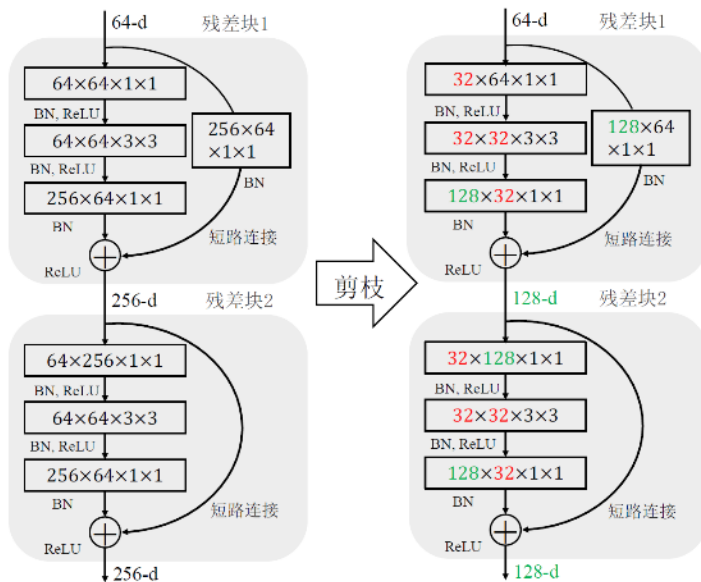
(c) 钱包结构

同时对残差结构的中间层与输出层进行剪枝是一种更好的策略

- 如何直接在小数据集上进行剪枝?
  - 直接在小数据集上进行剪枝效果不佳
    - T1：先在大数据集上进行剪枝，再迁移到小数据集上进行训练
    - T2：直接在小数据集上进行剪枝训练
    - 效果：T1 > T2
  - 直接在小数据集上进行剪枝应用更加实用
    - 由于一些实际情况（如隐私保护、法律问题等），原始大数据集可能无法获取到
    - 在大数据集上进行剪枝，代价昂贵，不利于迭代开发

如何提高小数据集上的直接剪枝效果？

## • 如何剪枝残差连接？



### □ 分析原因

由于残差结构限制，同一组内的残差块的**输出维度**需时刻保持**一致**来完成残差相加的操作

### □ 解决方法

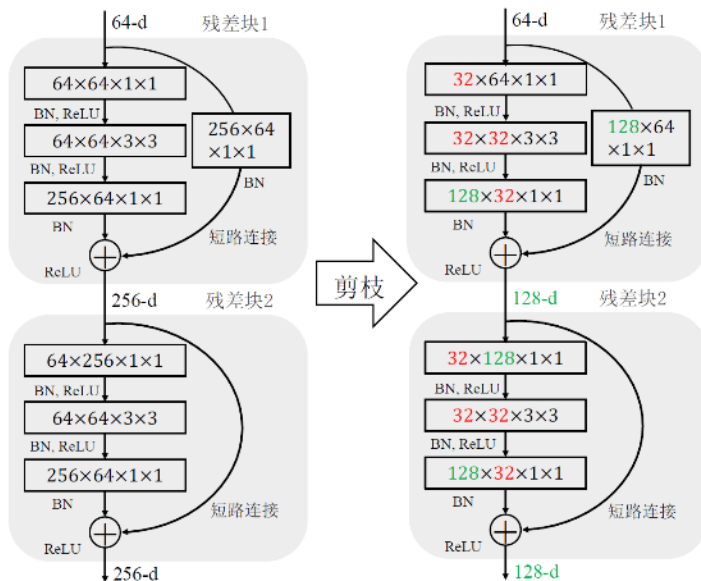
将**不同残差块**的输出维度看作一个整体，**同时**进行剪枝

### □ 新的挑战

如何衡量多个输出维度的重要性

评价标准：关注于网络输出，计算移除输出维度后的信息损失

## • 如何剪枝残差连接？



□ 剪枝输出维度 (绿色部分) :  
同时移除多个残差块同一输出维度的滤波器

□ 全局重要性衡量标准:

$$s = D_{KL}(p \parallel q) = \sum_{i=1}^n p_i \log \frac{p_i}{q_i}$$

p: 原网络的预测概率输出

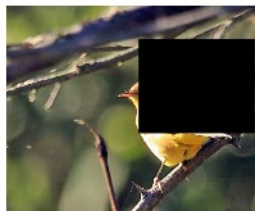
q: 移除滤波器后的概率输出



- 如何直接在小数据集上进行剪枝?
  - 分析原因: 数据量不足 + 监督信息弱
  - 解决数据量不足的问题: 图像变换以扩充数据集



(a) 原图



(b) 裁剪



(c) 旋转



(d) 2×2打散

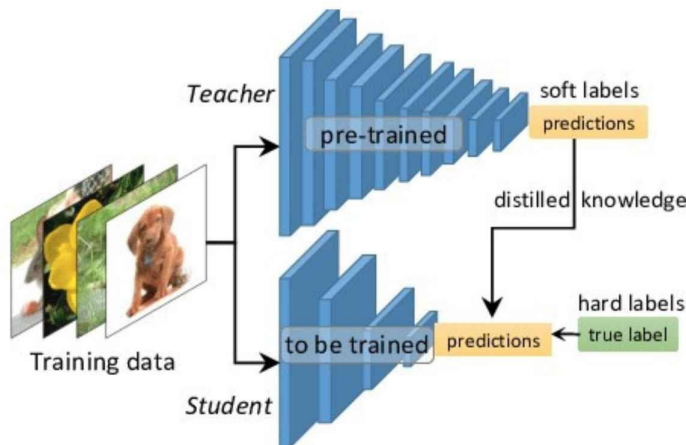


(e) 3×3打散



(f) 4×4打散

- 如何直接在小数据集上进行剪枝？
  - 分析原因：数据量不足 + 监督信息弱
  - 解决监督信息弱的问题：使用知识蒸馏技术



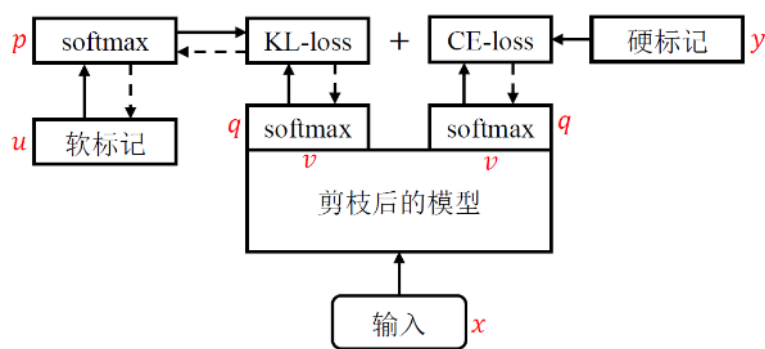
- 老师模型：原始网络
- 学生模型：剪枝后的网络

老师模型未见过扩充后的图像，其logits含噪声，不利于学生模型的学习

- 如何直接在小数据集上进行剪枝?

- 解决方法: **标记细化**

- 第一步: 在原始数据集上使用mixup和知识蒸馏技术微调剪枝后的模型, 获得一个精确的模型
- 第二步: 在扩充后的数据集上使用标记细化来移除噪声标记的干扰



- $p$ : 老师模型的输出;
- $q$ : 学生模型的输出;

$$\mathcal{L} = \alpha T^2 \mathcal{L}_{KL}(q, p) + (1 - \alpha) \mathcal{L}_{CE}(q, y)$$

- 同时对软标记进行更新以移除噪声标记的影响
- 逆转 $p, q$ 以促使网络更关注于软标记

- Kun Yi and Jianxin Wu. Probabilistic End-to-end Noise Correction for Learning with Noisy Labels. **In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)**, Long Beach, CA, USA, June 2019.

- 在ImageNet上剪枝ResNet50

方法	Top-1	Top-5	MACs	# 参数
ResNet50	76.15%	92.87%	4.09G	25.56M
ThiNet-30	68.42%	88.30%	1.10G	8.66M
GAL-1-joint <sup>[31]</sup>	69.31%	89.12%	1.11G	10.21M
GDP-0.5 <sup>[93]</sup>	69.58%	90.14%	1.57G	-
Taylor-FO <sup>[78]</sup>	71.69%	-	1.34G	7.90M
Slimmable NN <sup>[33]</sup>	72.10%	90.57%	1.05G	6.92M
AutoPruner	73.05%	91.25%	1.39G	12.60M
CURL	<b>73.39%</b>	<b>91.46%</b>	1.11G	6.67M

CURL的速度更快、精度更高、模型更小，证明了我们剪枝策略的有效性

- 小数据集上的剪枝

基模型	剪枝方法	CUB200 <sup>[60]</sup>		Pets <sup>[71]</sup>		Dogs <sup>[94]</sup>		Car <sup>[95]</sup>	
		精度	MACs	精度	MACs	精度	MACs	精度	MACs
MobileNetV2	MobileNetV2-1.0	78.77%	299.77M	89.94%	299.56M	78.94%	299.67M	87.27%	299.76M
	MobileNetV2-0.5	73.96%	96.12M	86.32%	95.91M	72.15%	96.02M	81.50%	96.11M
	Slimmable NN <sup>[33]</sup>	72.20%	96.12M	83.37%	95.91M	67.81%	96.02M	87.02%	96.11M
	CURL	<b>78.72%</b>	96.07M	<b>86.89%</b>	95.91M	<b>74.72%</b>	96.04M	<b>87.64%</b>	96.15M
ResNet50	ResNet50	84.76%	4.09G	92.59%	4.09G	83.82%	4.09G	91.89%	4.09G
	ResNet50-CURL	81.33%	1.11G	<b>90.84%</b>	1.11G	<b>81.60%</b>	1.11G	88.60%	1.11G
	Slimmable NN <sup>[33]</sup>	80.05%	1.05G	89.67%	1.05G	77.80%	1.05G	91.47%	1.05G
	CURL	<b>83.64%</b>	1.10G	90.30%	1.11G	79.79%	1.11G	<b>92.19%</b>	1.11G

直接在小数据集上剪枝效果不佳

□ CURL能够在小数据集上实现良好的压缩效果

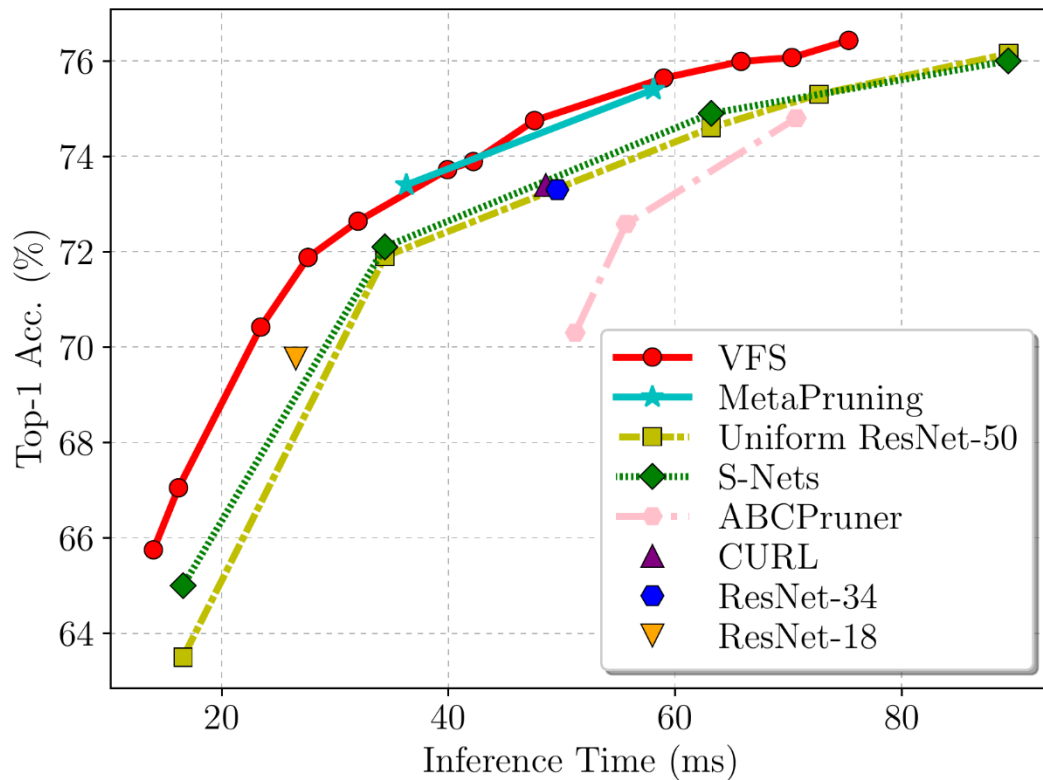
- 针对如何更好地剪枝残差结构以及如何在小数据集上进行网络剪枝的问题，我们提出了一种新的模型剪枝方案CURL
  - 我们会同时剪枝残差结构的中间层与输出维度，实现更快的运行速度、更高的模型精度，并有效降低了网络的内存消耗
  - 为了解决小数据集上训练困难的问题，我们使用图像变换来扩充训练数据集，并使用标记细化的技术来避免标记噪声的干扰
- 
- Jian-Hao Luo and Jianxin Wu. Neural Network Pruning with Residual-Connections and Limited-Data. **In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)**, Seattle, WA, USA, June 2020, accepted. (CCF-A类会议)

## 多设备模型生成

- 问题6：同时压缩滤波器、层、块？多应用压缩？(Versatile)
- 问题7：全光谱快速压缩？（Full-spectrum）
- 问题8：生成时间短、推断速度快？（Swift）

根据实际应用环境的需求  
确定需要解决的问题！

# VFS!





# Thank you!

