

模型压缩

主讲人 王晓波

中国科学院自动化研究所生物安全中心实验室博士
京东第二届博士技术管培生





人脸识别-模型压缩



1、小网络的设计



2、低秩分解和二值化



3、知识蒸馏



4、整体总结



人脸识别-模型压缩



1、小网络的设计



2、低秩分解和二值化



3、知识蒸馏



4、整体总结



1、小网络设计

- MegaFace竞赛

Algorithm	Details	Set1	Data Size
Sogou AIGROUP	multi model + combined margin loss	99.939%	Large
SRC-Beijing	Sphereface+large training set	99.888%	Large
ST-PureFace	attention-56+A-Softmax	99.801%	Large
El Networks	resnet-150 + resnet-100	99.414%	Large
ICARE_FACE_V1	ResNet101	99.319%	Large

几乎都是大模型，大训练数据



1、小网络设计

● LFR竞赛

layer name	124-layer	output size
Input Image Crop		$112 \times 112 \times 3$
	$3 \times 3, 64, \text{stride } 1$	$112 \times 112 \times 64$
Conv2_x	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$56 \times 56 \times 64$
Conv3_x	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 13$	$28 \times 28 \times 128$
Conv4_x	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 40$	$14 \times 14 \times 256$
Conv5_x	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 5$	$7 \times 7 \times 512$
FC		$1 \times 1 \times 512$

□ 29.70 Gflops

□ 297 MB



1、小网络设计

- LFR竞赛

Lightweight Face Recognition Challenge		
DeepGlint-Large	No.1	SEResNet-152
	No.2	AttentionIRSE-156
	No.3	ResNet-100
IQIYI-Large	No.1	xxx
	No.2	DenseNet-290
	No.3	ResNetSE-152



1、小网络设计

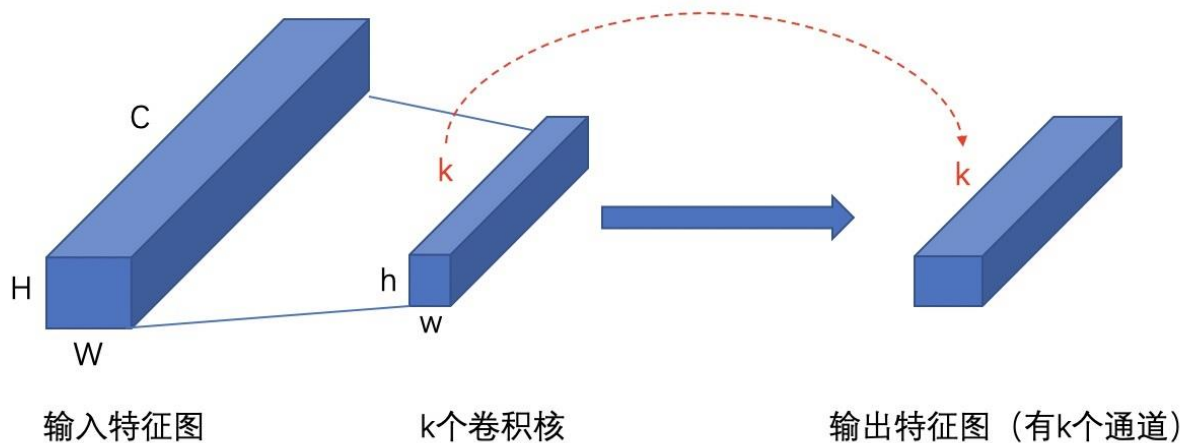
● MobileNet





1、小网络设计

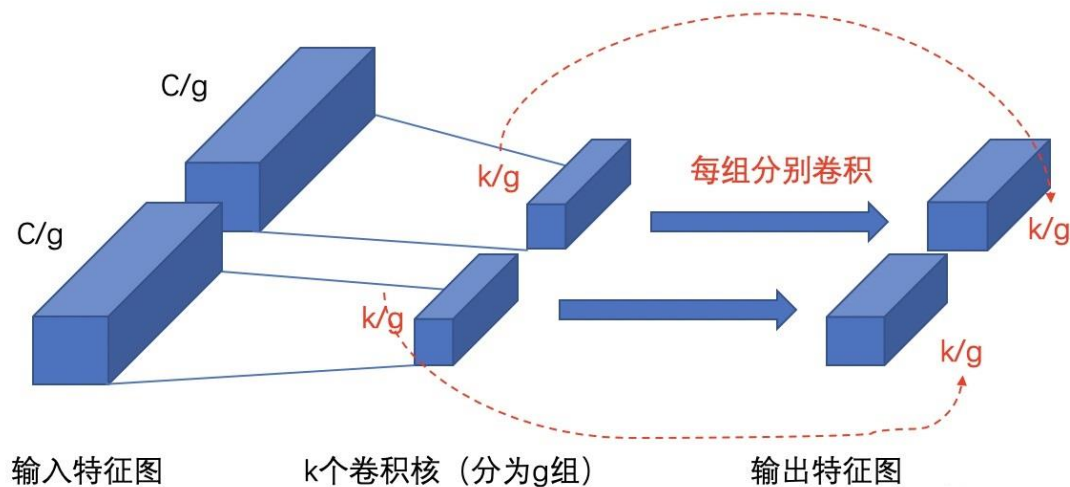
- GroupConv





1、小网络设计

- GroupConv

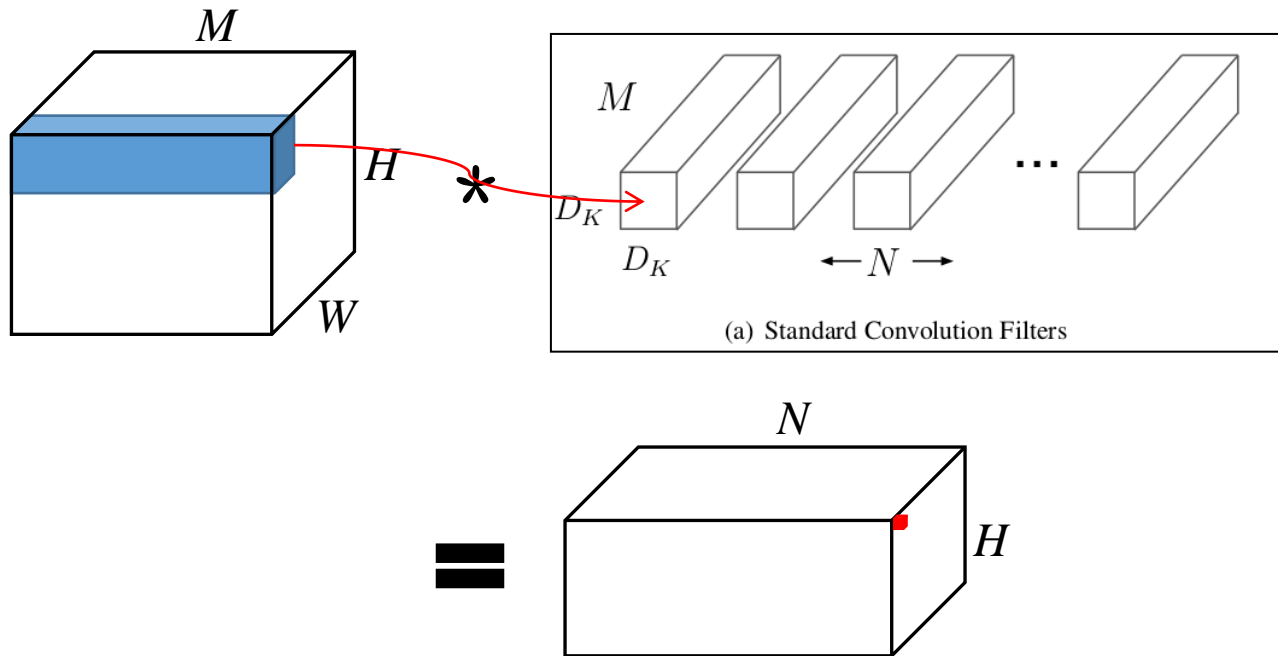


```
layer {  
  name: "conv2"  
  .....  
  convolution_param {  
    num_output: 256  
    pad: 2  
    kernel_size: 5  
    group: 2  
  }  
}
```



1、小网络设计

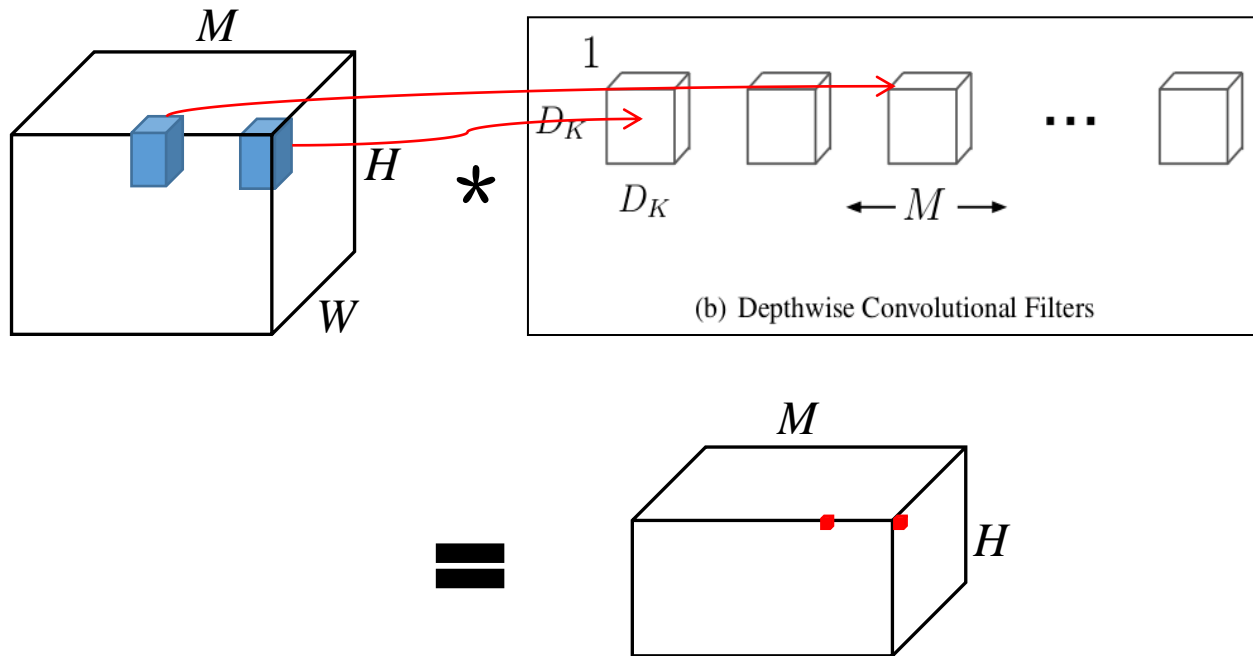
- MobileNet





1、小网络设计

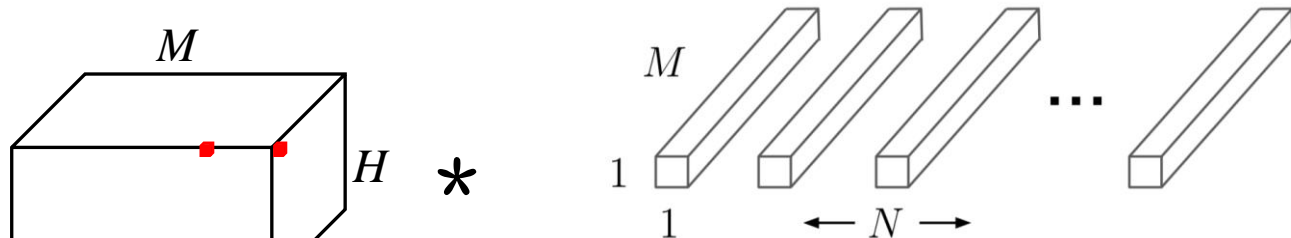
- MobileNet



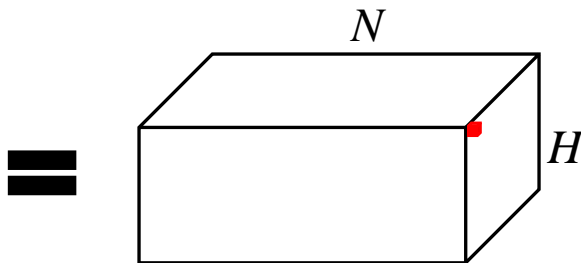


1、小网络设计

- MobileNet



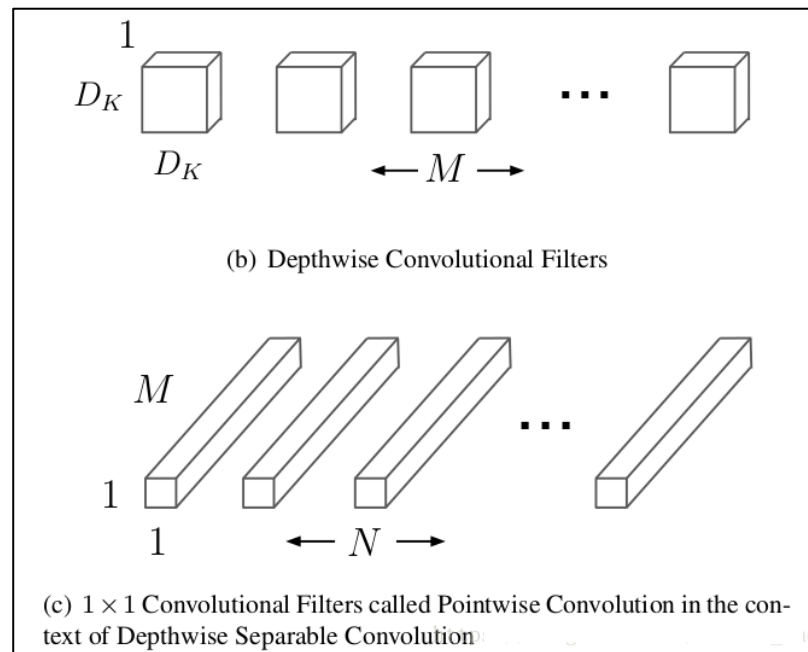
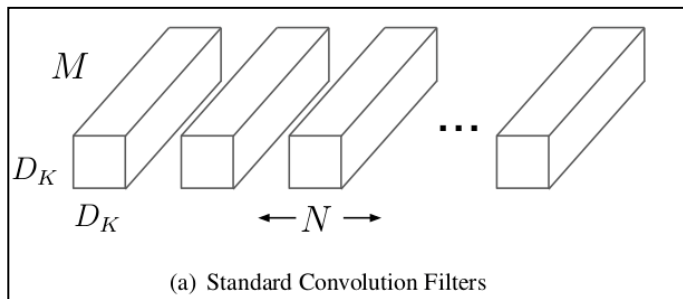
(c) 1×1 Convolutional Filters called Pointwise Convolution in the context of Depthwise Separable Convolution





1、小网络设计

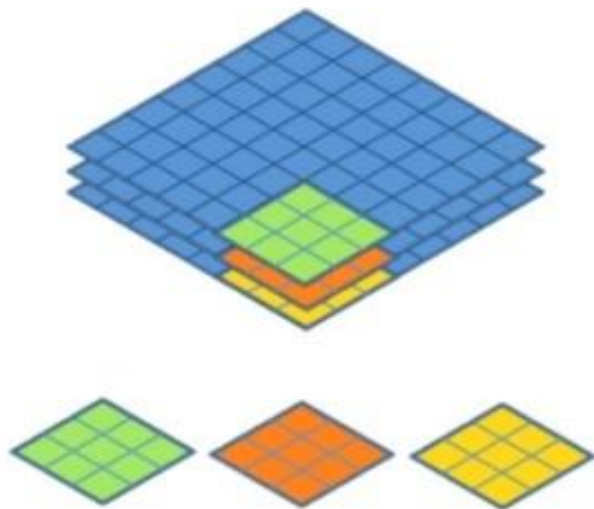
● MobileNet



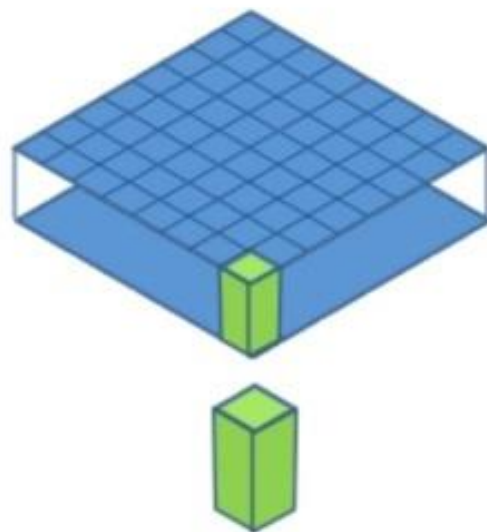


1、小网络设计

- MobileNet



Depthwise Convolution

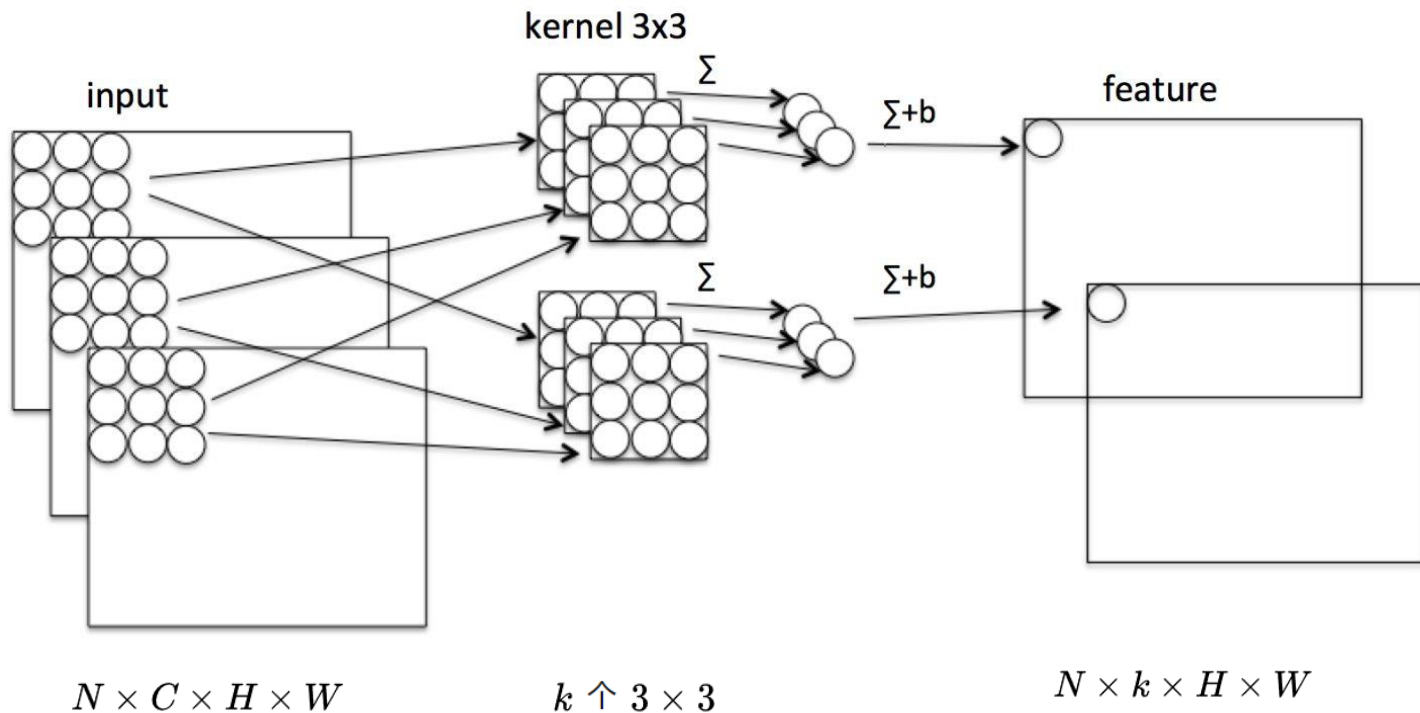


Pointwise Convolution



1、小网络设计

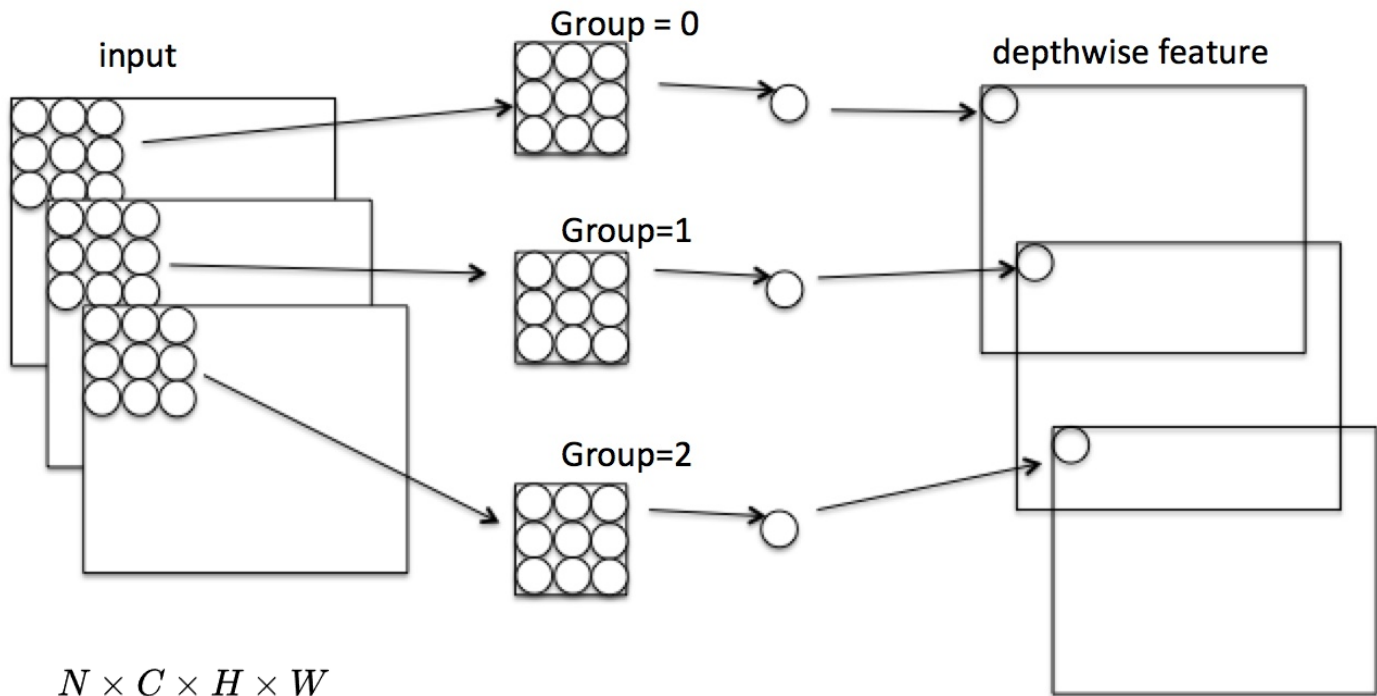
- MobileNet





1、小网络设计

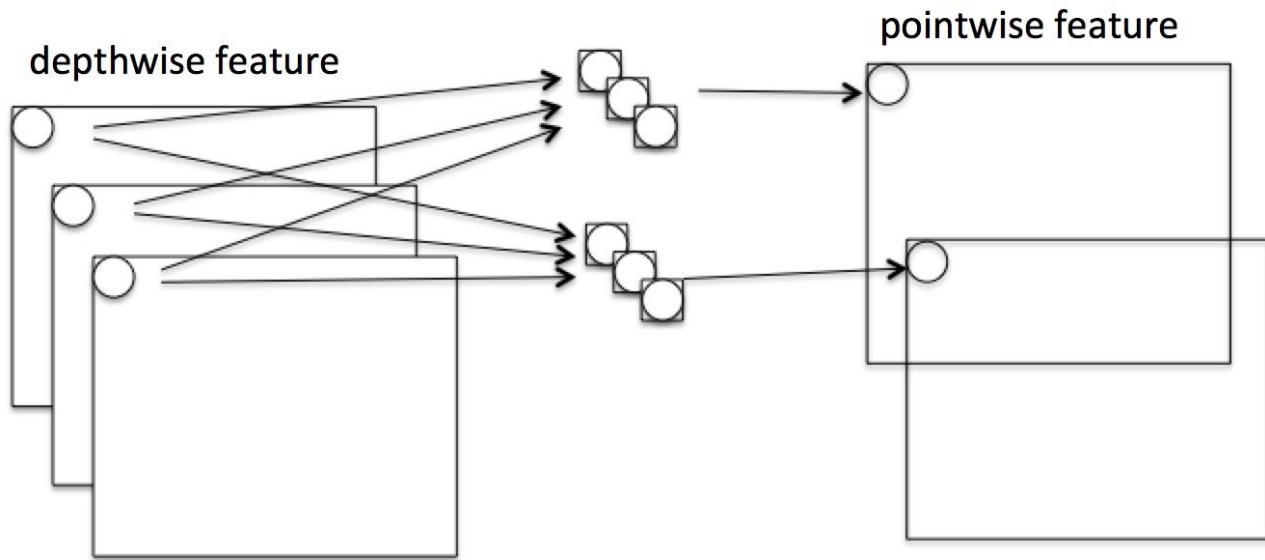
- MobileNet





1、小网络设计

- MobileNet

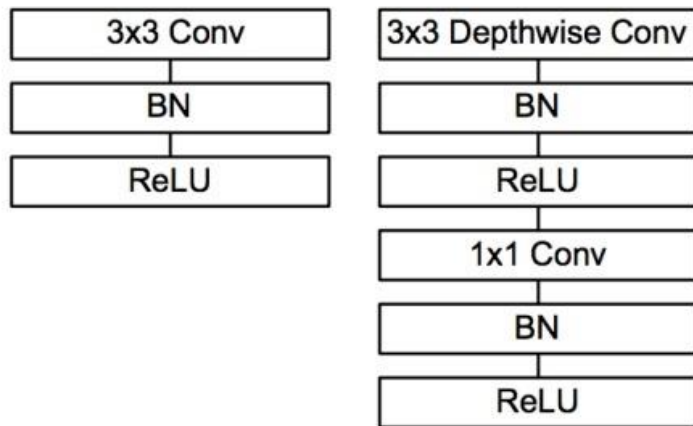


$$N \times C \times H \times W$$



1、小网络设计

- MobileNet



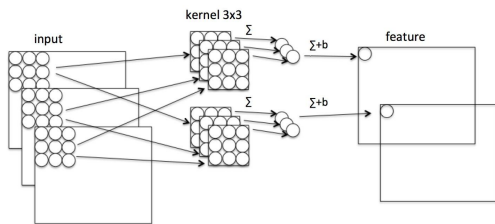
- Depthwise通道间信息隔离
- Pointwise通道间信息交换
-

Figure 3. Left: Standard convolutional layer with batchnorm and ReLU. Right: Depthwise Separable convolutions with Depthwise and Pointwise layers followed by batchnorm and ReLU.



1、小网络设计

- MobileNet



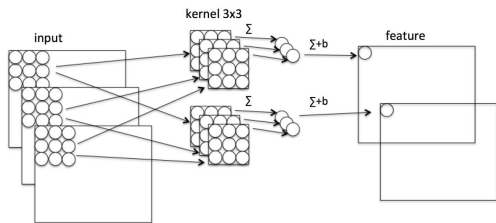
普通卷积计算量

$$H \times W \times C \times k \times 3 \times 3$$



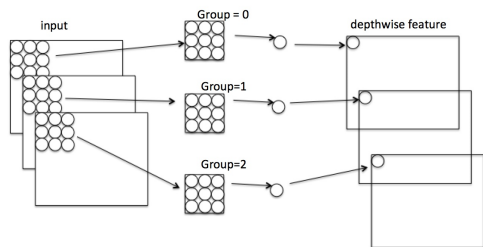
1、小网络设计

● MobileNet



普通卷积计算量

$$H \times W \times C \times k \times 3 \times 3$$



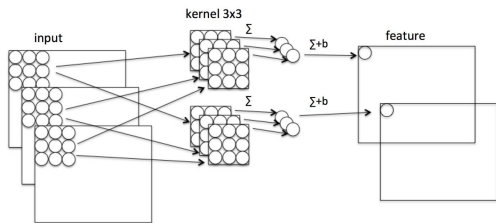
Depthwise计算量

$$H \times W \times C \times 3 \times 3$$



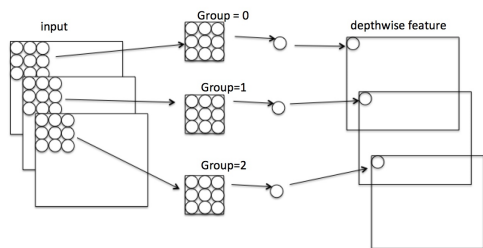
1、小网络设计

● MobileNet



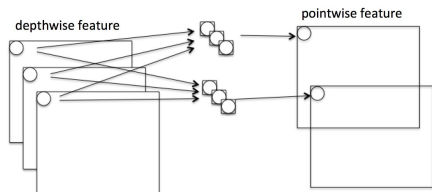
普通卷积计算量

$$H \times W \times C \times k \times 3 \times 3$$



Depthwise计算量

$$H \times W \times C \times 3 \times 3$$



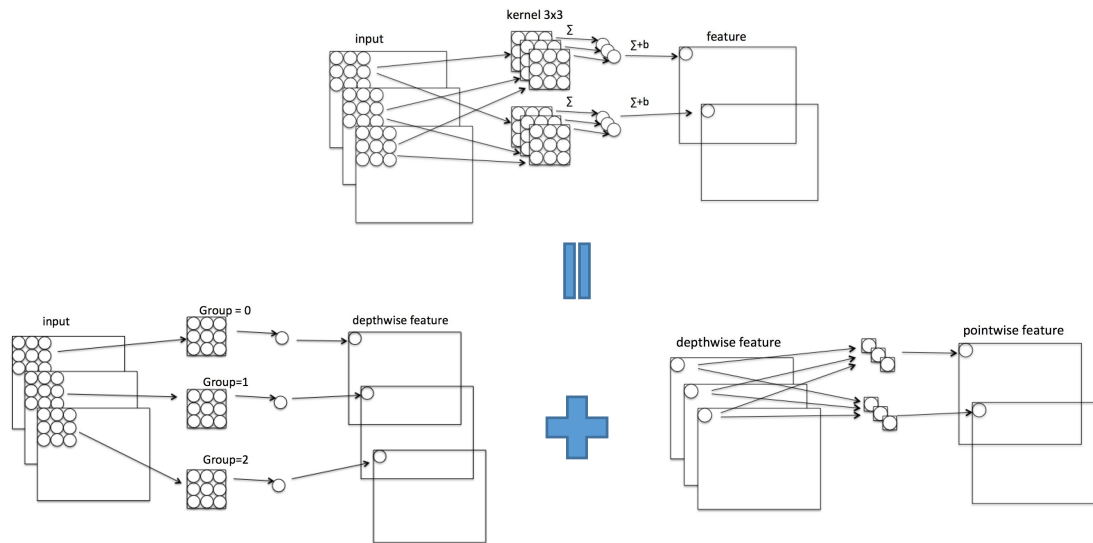
Pointwise计算量

$$H \times W \times C \times k$$



1、小网络设计

● MobileNet



$$\frac{\text{depthwise} + \text{pointwise}}{\text{conv}} = \frac{H \times W \times C \times 3 \times 3 + H \times W \times C \times k}{H \times W \times C \times k \times 3 \times 3} = \frac{1}{k} + \frac{1}{3 \times 3}$$



1、小网络设计

● MobileNet

Table 1. MobileNet Body Architecture

Type / Stride	Filter Shape	Input Size
Conv / s2	$3 \times 3 \times 3 \times 32$	$224 \times 224 \times 3$
Conv dw / s1	$3 \times 3 \times 32$ dw	$112 \times 112 \times 32$
Conv / s1	$1 \times 1 \times 32 \times 64$	$112 \times 112 \times 32$
Conv dw / s2	$3 \times 3 \times 64$ dw	$112 \times 112 \times 64$
Conv / s1	$1 \times 1 \times 64 \times 128$	$56 \times 56 \times 64$
Conv dw / s1	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 128$	$56 \times 56 \times 128$
Conv dw / s2	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 256$	$28 \times 28 \times 128$
Conv dw / s1	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 256$	$28 \times 28 \times 256$
Conv dw / s2	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 512$	$14 \times 14 \times 256$
5×	Conv dw / s1 $3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
	Conv / s1 $1 \times 1 \times 512 \times 512$	$14 \times 14 \times 512$
Conv dw / s2	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
Conv / s1	$1 \times 1 \times 512 \times 1024$	$7 \times 7 \times 512$
Conv dw / s2	$3 \times 3 \times 1024$ dw	$7 \times 7 \times 1024$
Conv / s1	$1 \times 1 \times 1024 \times 1024$	$7 \times 7 \times 1024$
Avg Pool / s1	Pool 7×7	$7 \times 7 \times 1024$
FC / s1	1024×1000	$1 \times 1 \times 1024$
Softmax / s1	Classifier	$1 \times 1 \times 1000$

Width Multiplier	ImageNet Accuracy	Million Mult-Adds	Million Parameters
1.0 MobileNet-224	70.6%	569	4.2
0.75 MobileNet-224	68.4%	325	2.6
0.5 MobileNet-224	63.7%	149	1.3
0.25 MobileNet-224	50.6%	41	0.5

Model	ImageNet Accuracy	Million Mult-Adds	Million Parameters
1.0 MobileNet-224	70.6%	569	4.2
GoogleNet	69.8%	1550	6.8
VGG 16	71.5%	15300	138



1、小网络设计

● MobileFaceNet

Table 1. MobileNet Body Architecture

Type / Stride	Filter Shape	Input Size
Conv / s2	$3 \times 3 \times 3 \times 32$	$224 \times 224 \times 3$
Conv dw / s1	$3 \times 3 \times 32$ dw	$112 \times 112 \times 32$
Conv / s1	$1 \times 1 \times 32 \times 64$	$112 \times 112 \times 32$
Conv dw / s2	$3 \times 3 \times 64$ dw	$112 \times 112 \times 64$
Conv / s1	$1 \times 1 \times 64 \times 128$	$56 \times 56 \times 64$
Conv dw / s1	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 128$	$56 \times 56 \times 128$
Conv dw / s2	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 256$	$28 \times 28 \times 128$
Conv dw / s1	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 256$	$28 \times 28 \times 256$
Conv dw / s2	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 512$	$14 \times 14 \times 256$
5×	Conv dw / s1 $3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
	Conv / s1 $1 \times 1 \times 512 \times 512$	$14 \times 14 \times 512$
Conv dw / s2	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
Conv / s1	$1 \times 1 \times 512 \times 1024$	$7 \times 7 \times 512$
Conv dw / s2	$3 \times 3 \times 1024$ dw	$7 \times 7 \times 1024$
Conv / s1	$1 \times 1 \times 1024 \times 1024$	$7 \times 7 \times 1024$
Avg Pool / s1	Pool 7×7	$7 \times 7 \times 1024$
FC / s1	1024×1000	$1 \times 1 \times 1024$
Softmax / s1	Classifier	$1 \times 1 \times 1000$

- MobileV2, V3
- ShuffleNet etc



Global Depthwise Convolution



人脸识别-模型压缩



1、小网络的设计



2、低秩分解和二值化



3、知识蒸馏

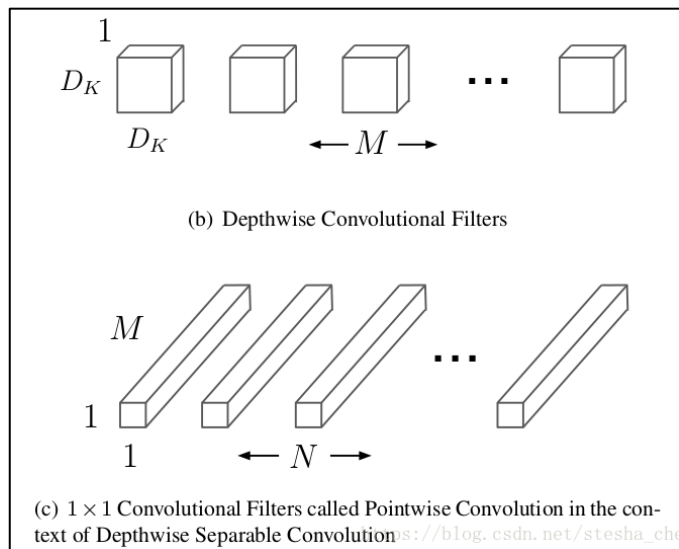
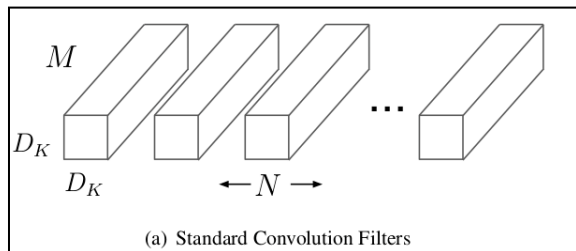


4、整体总结



2、低秩分解

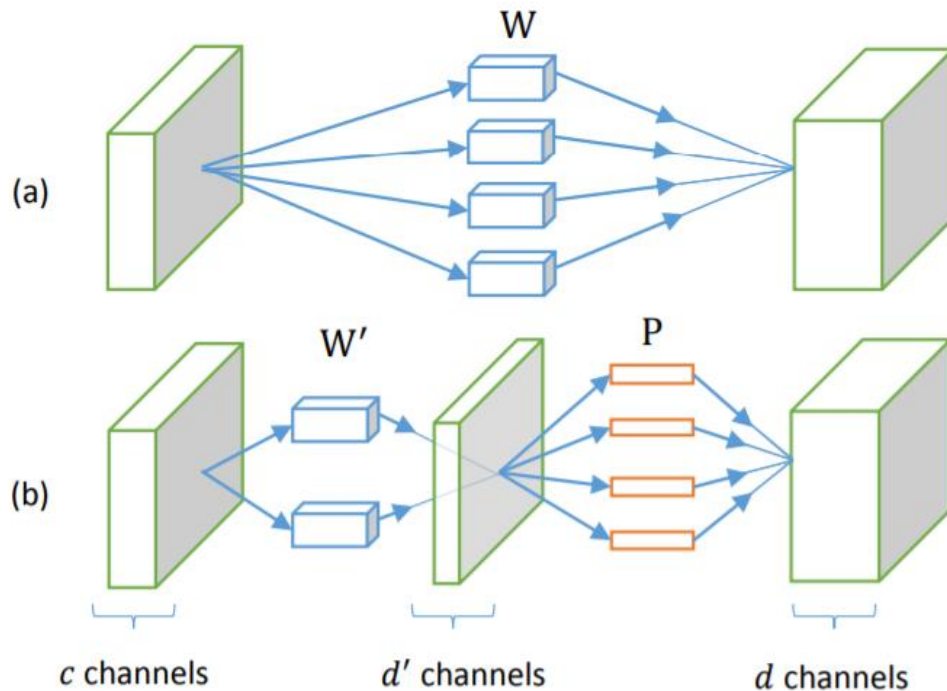
- MobileNet





2、低秩分解

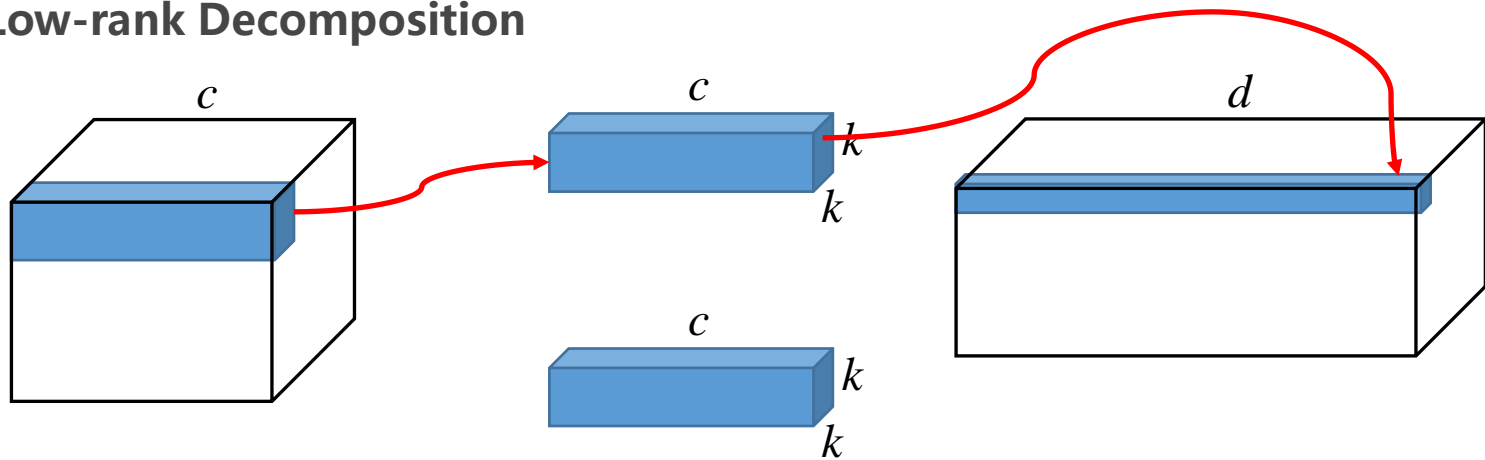
- Low-rank Decomposition





2、低秩分解

- Low-rank Decomposition



$$x \in \mathbb{R}^{k^2 c + 1}$$

$$W \in \mathbb{R}^{d \times (k^2 c + 1)}$$

$$y \in \mathbb{R}^d$$

$$y = Wx$$



2、低秩分解

- Low-rank Decomposition

$$y = M(y - \bar{y}) + \bar{y}, \quad M \in \mathbb{R}^{d \times d}$$



$$y = MWx + b, \quad b = \bar{y} - M\bar{y}$$



$$M = PQ^T, \quad W' = Q^TW$$

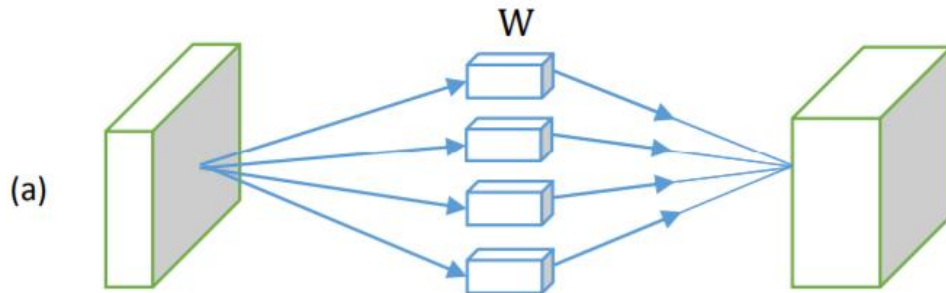
$$y = PW'x + b$$



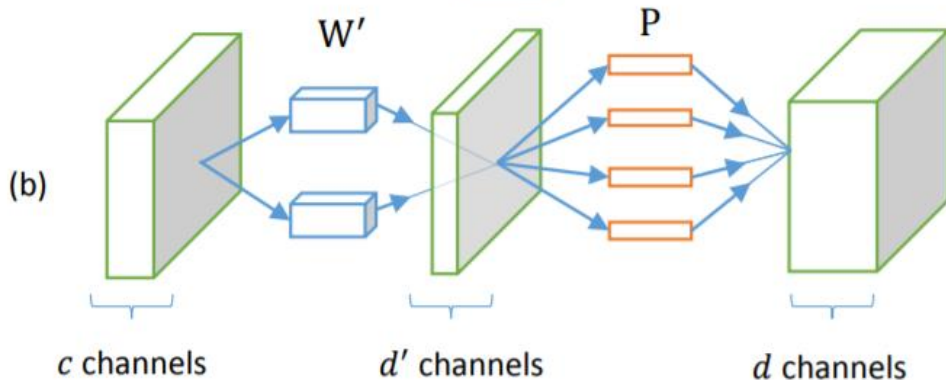
2、低秩分解

- Low-rank Decomposition

$$y = Wx$$



$$y = PW'x + b$$





2、低秩分解

- Low-rank Decomposition

$$\min_{\mathbf{M}} \sum_i \|(\mathbf{y}_i - \bar{\mathbf{y}}) - \mathbf{M}(\mathbf{y}_i - \bar{\mathbf{y}})\|_2^2,$$
$$s.t. \quad rank(\mathbf{M}) \leq d'.$$

$$\mathbf{y} = \mathbf{M}(\mathbf{y} - \bar{\mathbf{y}}) + \bar{\mathbf{y}}, \quad \mathbf{M} \in \mathbb{R}^{d \times d}$$

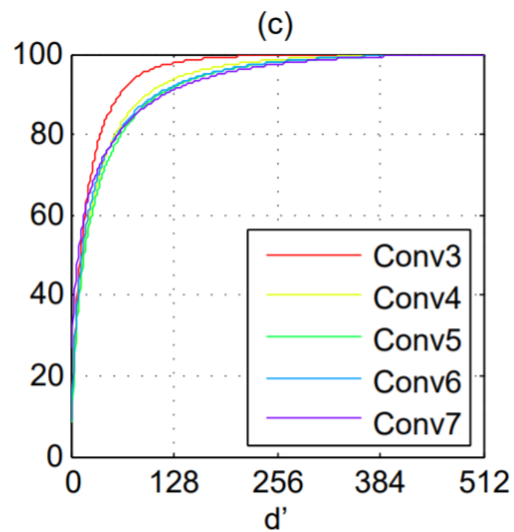
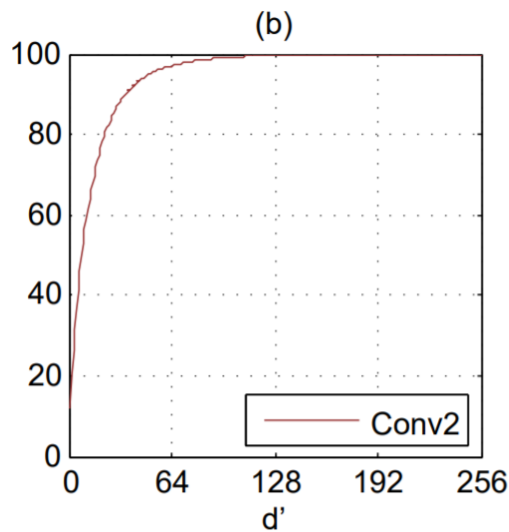
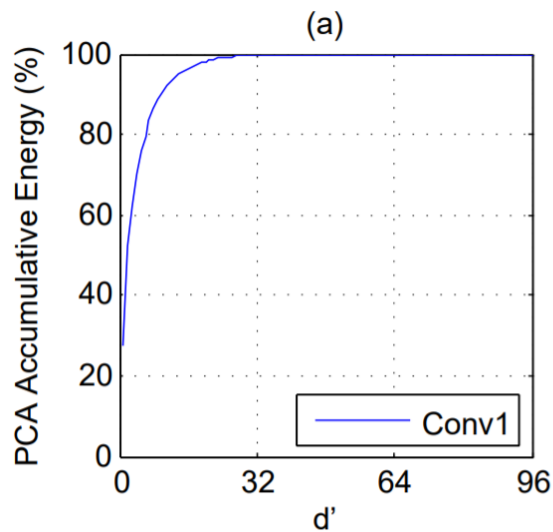
$\mathbf{Y}^T \mathbf{Y}$ 进行特征值分解, $\mathbf{Y}^T \mathbf{Y} = \mathbf{U} \mathbf{S} \mathbf{U}^T$, 其中 \mathbf{U} 是正交矩阵, \mathbf{S} 是对角线矩阵。
 $\mathbf{M} = \mathbf{U}_{d'} \mathbf{U}_{d'}^T$, 其中 $\mathbf{U}_{d'}$ 就是 \mathbf{U} 的前 d' 个特征向量。可设定 $\mathbf{P} = \mathbf{Q} = \mathbf{U}_{d'}$

$$\mathbf{W}' = \mathbf{Q}^T \mathbf{W} \quad \mathbf{y} = \mathbf{P} \mathbf{W}' \mathbf{x} + \mathbf{b}$$



2、低秩分解

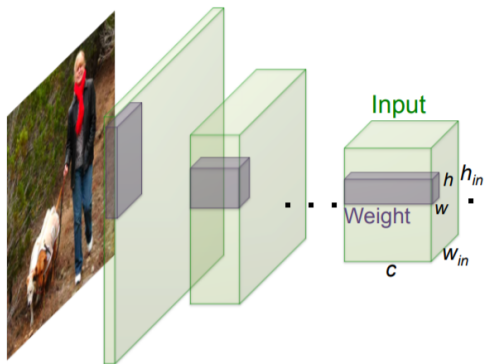
- Low-rank Decomposition





2、二值化

- XNOR-Net

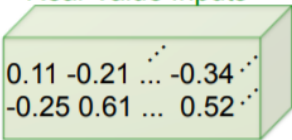
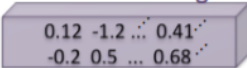
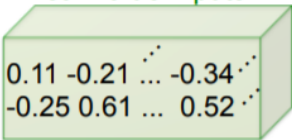
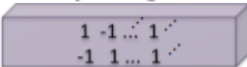


	Network Variations	Operations used in Convolution	Memory Saving (Inference)	Computation Saving (Inference)	Accuracy on ImageNet (AlexNet)
Standard Convolution	Real-Value Inputs Real-Value Weights 	$+, -, \times$	1x	1x	%56.7
Binary Weight	Real-Value Inputs Binary Weights 	$+, -$	$\sim 32x$	$\sim 2x$	%56.8
BinaryWeight Binary Input (XNOR-Net)	Binary Inputs Binary Weights 	XNOR , bitcount	$\sim 32x$	$\sim 58x$	%44.2



2、二值化

- XNOR-Net

	Network Variations	
Standard Convolution	<p>Real-Value Inputs</p> 	<p>Real-Value Weights</p> 
▪ Binary Weight	<p>Real-Value Inputs</p> 	<p>Binary Weights</p> 

real-value weight filter $\mathbf{W} \in \mathcal{W}$

binary filter $\mathbf{B} \in \{+1, -1\}^{c \times w \times h}$

$$\mathbf{W} \approx \alpha \mathbf{B}$$

a scaling factor $\alpha \in \mathbb{R}^+$



2、二值化

- XNOR-Net

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

1	0	1
0	1	0
1	0	1

$$\mathbf{I} * \mathbf{W} \approx (\mathbf{I} \oplus \mathbf{B}) \alpha$$

\oplus indicates a convolution without any multiplication



2、二值化

- XNOR-Net

$$J(\mathbf{B}, \alpha) = \|\mathbf{W} - \alpha\mathbf{B}\|^2$$

$$\alpha^*, \mathbf{B}^* = \underset{\alpha, \mathbf{B}}{\operatorname{argmin}} J(\mathbf{B}, \alpha)$$

$$J(\mathbf{B}, \alpha) = \alpha^2 \mathbf{B}^\top \mathbf{B} - 2\alpha \mathbf{W}^\top \mathbf{B} + \mathbf{W}^\top \mathbf{W}$$

since $\mathbf{B} \in \{+1, -1\}^n$, $\mathbf{B}^\top \mathbf{B} = n$ is a constant $n = c \times w \times h$

$$\mathbf{c} = \mathbf{W}^\top \mathbf{W}$$

$$\mathbf{B}^* = \underset{\mathbf{B}}{\operatorname{argmax}} \{\mathbf{W}^\top \mathbf{B}\} \quad s.t. \quad \mathbf{B} \in \{+1, -1\}^n$$



2、二值化

- XNOR-Net

$$J(\mathbf{B}, \alpha) = \alpha^2 \mathbf{B}^\top \mathbf{B} - 2\alpha \mathbf{W}^\top \mathbf{B} + \mathbf{W}^\top \mathbf{W}$$



$$\mathbf{B}^* = \underset{\mathbf{B}}{\operatorname{argmax}} \{ \mathbf{W}^\top \mathbf{B} \} \quad s.t. \quad \mathbf{B} \in \{+1, -1\}^n$$



最优的解：当 $\mathbf{W}_i \geq 0$ 时， $\mathbf{B}_i = 1$ ；当 $\mathbf{W}_i < 0$ 时， $\mathbf{B}_i = -1$

$$\mathbf{B}^* = \operatorname{sign}(\mathbf{W})$$



2、二值化

- XNOR-Net

$$J(\mathbf{B}, \alpha) = \alpha^2 \mathbf{B}^T \mathbf{B} - 2\alpha \mathbf{W}^T \mathbf{B} + \mathbf{W}^T \mathbf{W}$$



$$J(B, \alpha) = \alpha^2 n - 2\alpha \mathbf{W}^T \mathbf{B} + c$$



$$\alpha = \frac{\mathbf{W}^T \mathbf{B}^*}{n} = \frac{\mathbf{W}^T \text{sign}(\mathbf{W})}{n} = \frac{\sum |\mathbf{W}_i|}{n} = \frac{1}{n} \|\mathbf{W}\|_1$$



2、二值化

- XNOR-Net

$\mathbf{X}^\top \mathbf{W} \approx \beta \mathbf{H}^\top \alpha \mathbf{B}$, where $\mathbf{H}, \mathbf{B} \in \{+1, -1\}^n$ and $\beta, \alpha \in \mathbb{R}^+$

$$\alpha^*, \mathbf{B}^*, \beta^*, \mathbf{H}^* = \underset{\alpha, \mathbf{B}, \beta, \mathbf{H}}{\operatorname{argmin}} \|\mathbf{X} \odot \mathbf{W} - \beta \alpha \mathbf{H} \odot \mathbf{B}\|$$



$$\gamma^*, \mathbf{C}^* = \underset{\gamma, \mathbf{C}}{\operatorname{argmin}} \|\mathbf{Y} - \gamma \mathbf{C}\|$$

$$\mathbf{Y}_i = \mathbf{X}_i \mathbf{W}_i$$

$$\mathbf{C}_i = \mathbf{H}_i \mathbf{B}_i$$

$$\gamma = \beta \alpha$$



2、二值化

- XNOR-Net

$$\gamma^*, \mathbf{C}^* = \underset{\gamma, \mathbf{C}}{\operatorname{argmin}} \|\mathbf{Y} - \gamma \mathbf{C}\|$$

$$\mathbf{C}^* = \operatorname{sign}(\mathbf{Y}) = \operatorname{sign}(\mathbf{X}) \odot \operatorname{sign}(\mathbf{W}) = \mathbf{H}^* \odot \mathbf{B}^*$$

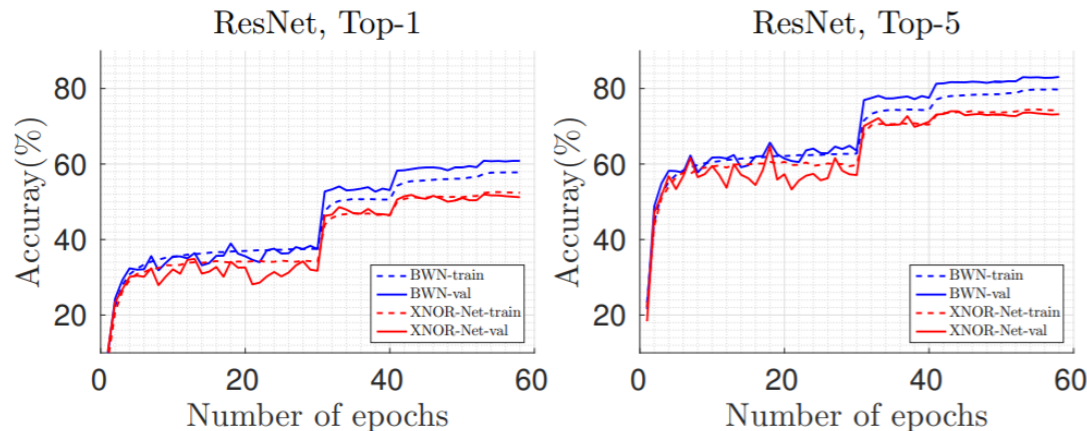
$$\gamma^* = \frac{\sum |\mathbf{Y}_i|}{n} = \frac{\sum |\mathbf{X}_i| |\mathbf{W}_i|}{n} \approx \left(\frac{1}{n} \|\mathbf{X}\|_{\ell_1} \right) \left(\frac{1}{n} \|\mathbf{W}\|_{\ell_1} \right) = \beta^* \alpha^*$$



2、二值化

● XNOR-Net

Classification Accuracy(%)									
Binary-Weight				Binary-Input-Binary-Weight				Full-Precision	
BWN		BC[11]		XNOR-Net		BNN[11]		AlexNet[1]	
Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
56.8	79.4	35.4	61.0	44.2	69.2	27.9	50.42	56.6	80.2





人脸识别-模型压缩



1、小网络的设计



2、低秩分解和二值化



3、知识蒸馏

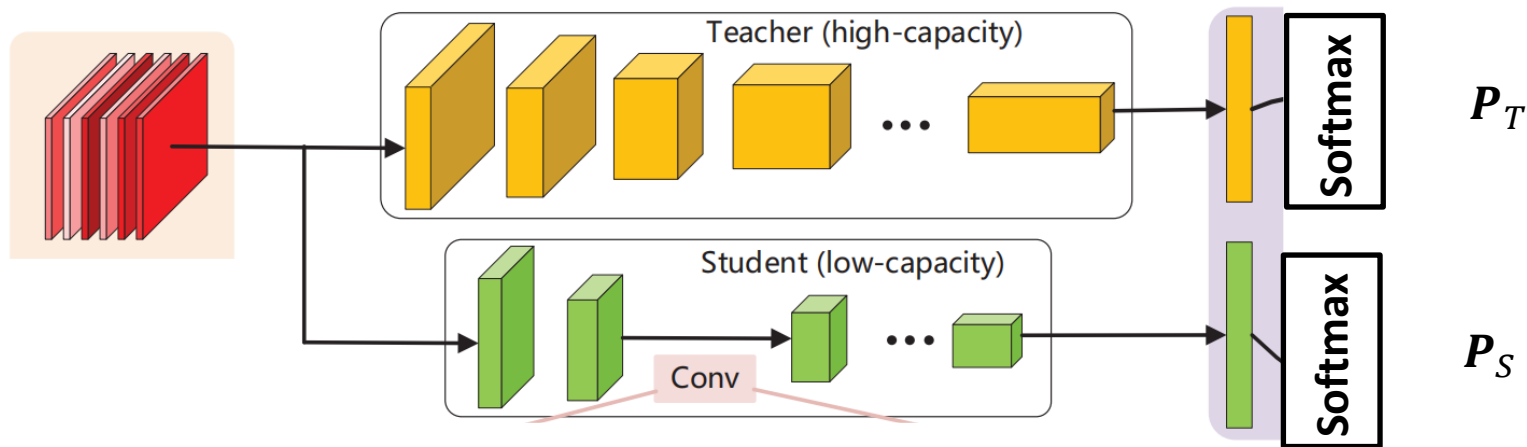


4、整体总结



3、知识蒸馏

- Knowledge Distillation

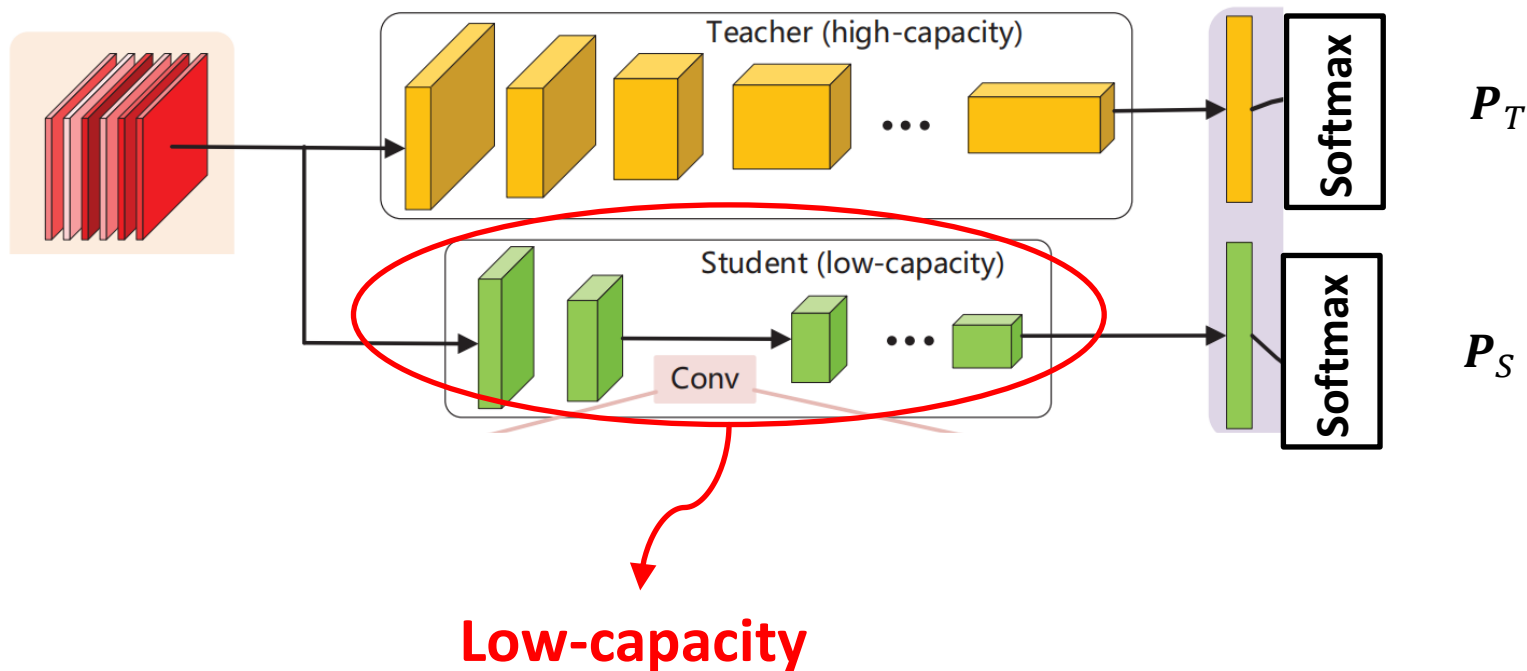


$$\mathcal{L}_{\text{PC}} := \mathcal{L}(P_T^\tau, P_S^\tau) = \mathcal{L}((z_T/\tau), (z_S/\tau))$$



3、知识蒸馏

- EC-KD





3、知识蒸馏

- EC-KD

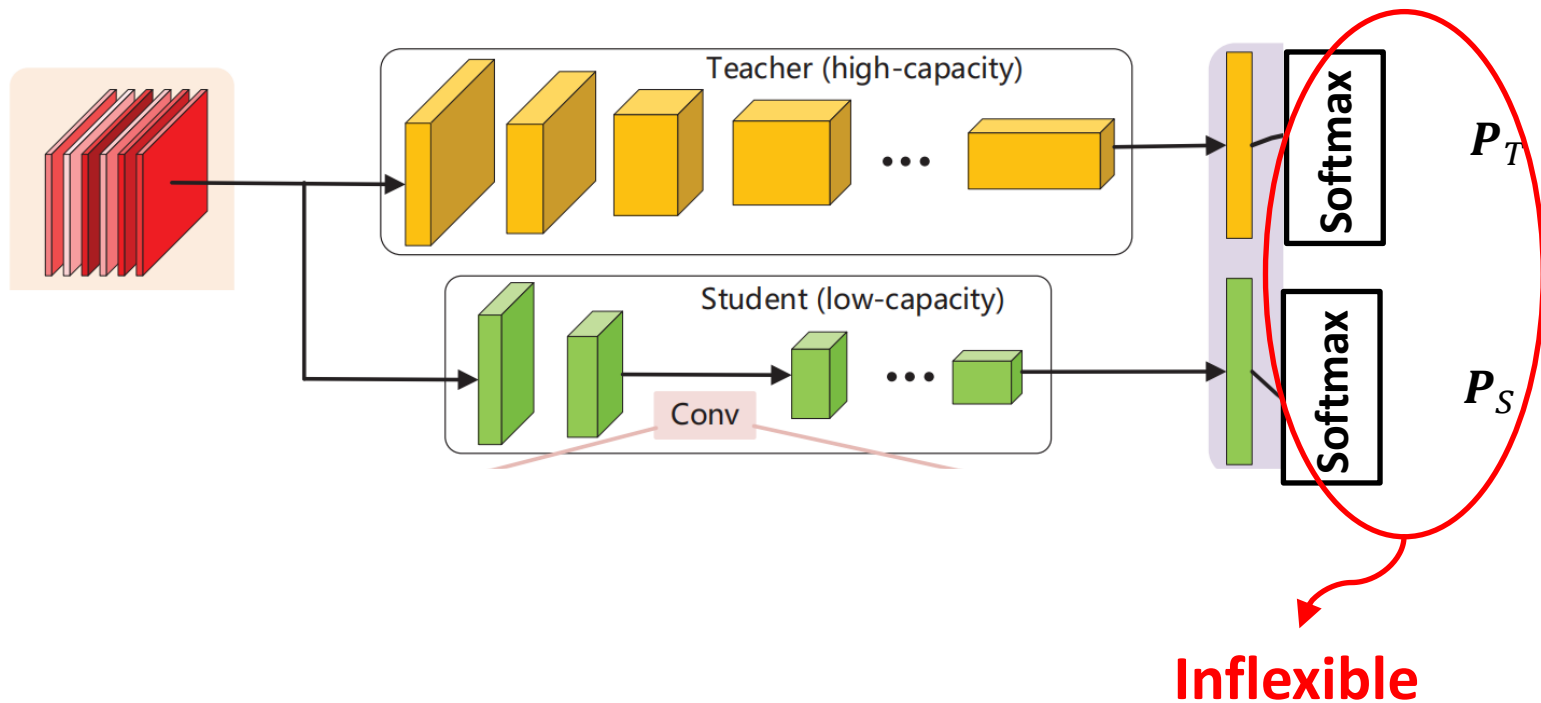
#Filters	Model Size	Flops	Infer Time	LFW	MF-Id.	MF-Veri.
(2×)128	16MB	1.44G	158ms	99.34	87.19	90.82
(Orig.)64	4.8MB	0.38G	84ms	99.11	83.96	87.57
(1/2)32	1.7MB	0.11G	49ms	98.55	74.32	78.71
(1/4)16	648KB	0.03G	34ms	97.60	52.60	58.69
(1/8)8	304KB	0.01G	28ms	94.29	25.32	27.04

Student 每一层的kernel数目越少，模型越小，性能则一般越低



3、知识蒸馏

- EC-KD





3、知识蒸馏

- EC-KD

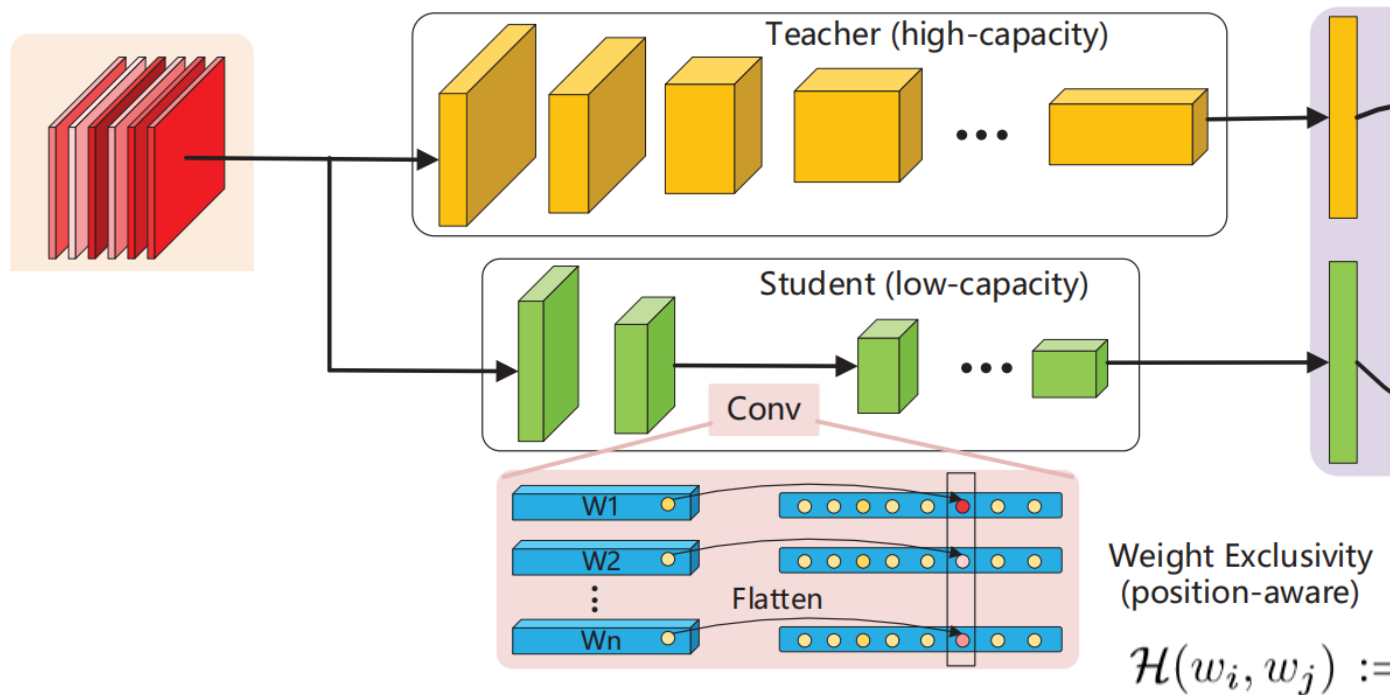
$$\mathcal{L}_{\text{PC}} := \mathcal{L}(P_{\text{T}}^{\tau}, P_{\text{S}}^{\tau}) = \mathcal{L}((z_{\text{T}}/\tau), (z_{\text{S}}/\tau))$$

- ❑ Teacher的训练类别和Student的训练类别不一致，或者Teacher模型由Contrastive Loss等训练时，上述公式的Softmax没法计算。
- ❑ Student的训练数据包含噪声标签时，性能没有办法保证。
- ❑ Student训练类别数较多时，训练时间长，收敛速度比较慢。



3、知识蒸馏

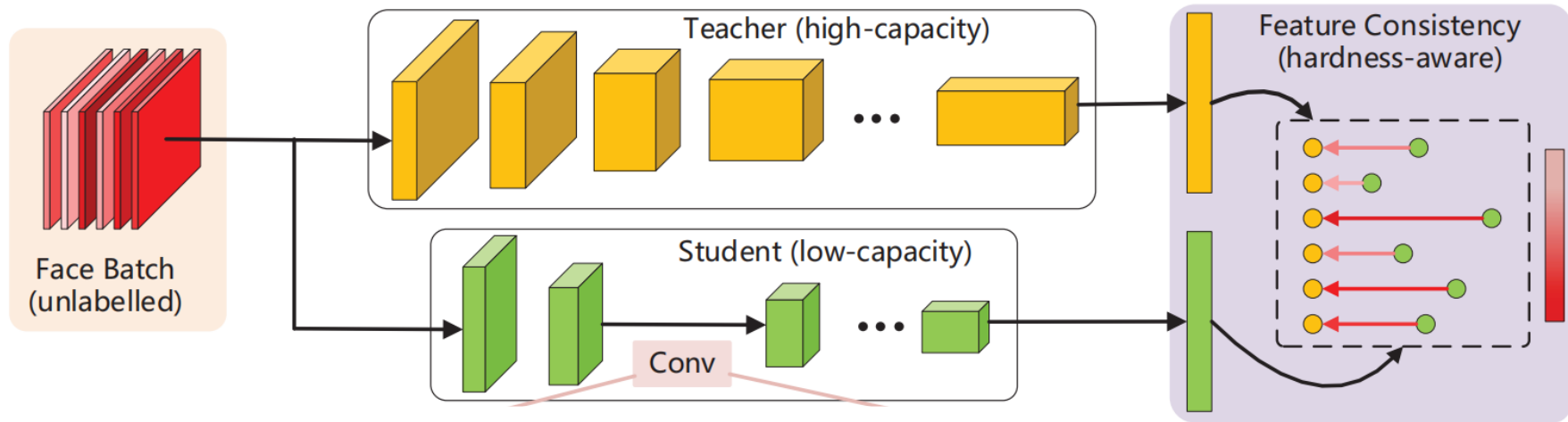
- EC-KD





3、知识蒸馏

- EC-KD

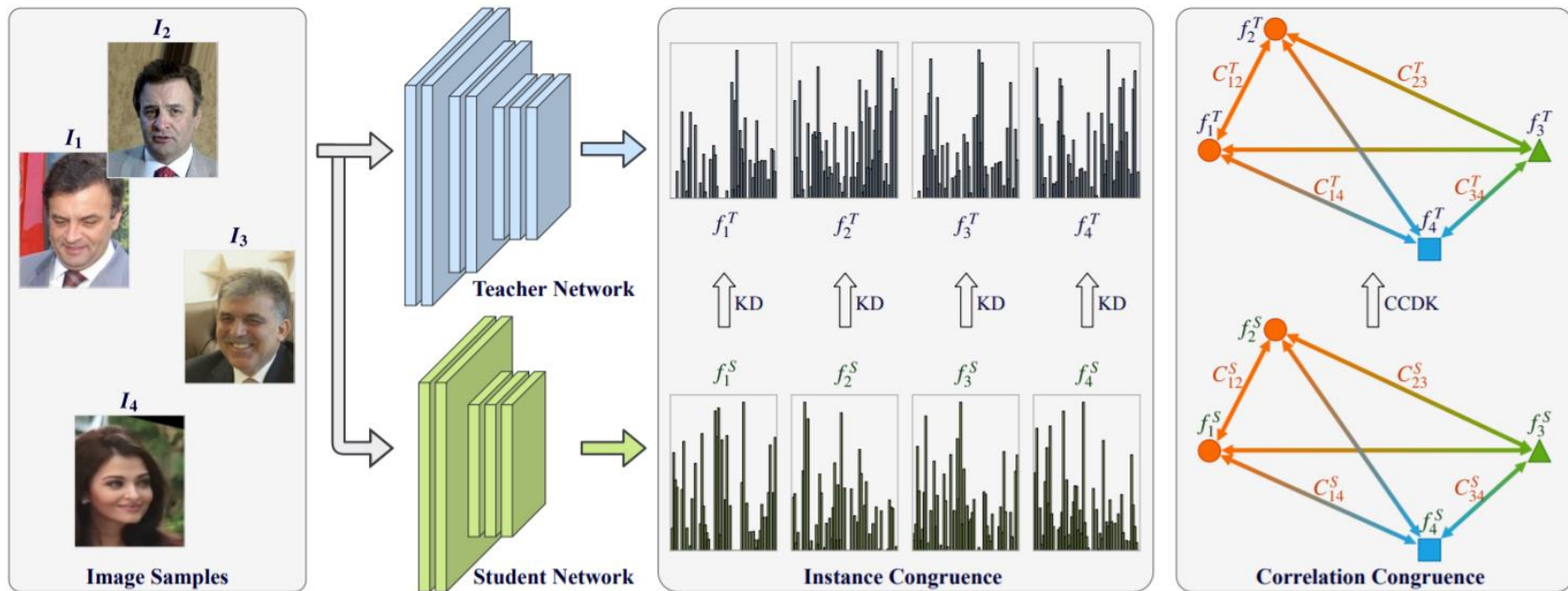


$$\mathcal{L}_{\text{FC}} := \mathcal{H}(F_S, F_T) = ||F_S - F_T||.$$



3、知识蒸馏

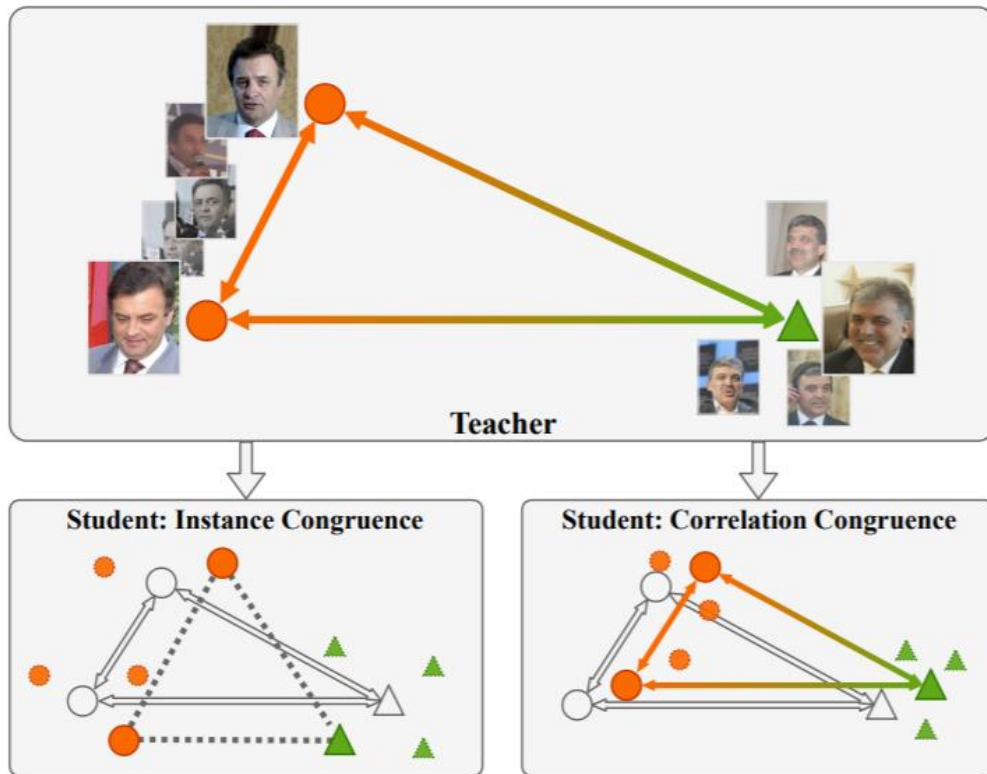
- CC-KD





3、知识蒸馏

- CC-KD



Correlation Congruence for Knowledge Distillation



3、知识蒸馏

- CC-KD

$$\mathbf{F}_t = \text{matrix}(\mathbf{f}_1^t, \mathbf{f}_2^t, \dots, \mathbf{f}_n^t),$$

$$\mathbf{F}_s = \text{matrix}(\mathbf{f}_1^s, \mathbf{f}_2^s, \dots, \mathbf{f}_n^s).$$

$$\psi : \mathbf{F} \rightarrow \mathbf{C} \in \mathbb{R}^{n \times n} \quad C_{ij} = \varphi(\mathbf{f}_i, \mathbf{f}_j), \quad C_{ij} \in \mathbb{R}$$

$$\begin{aligned} L_{CC} &= \frac{1}{n^2} \|\psi(\mathbf{F}_t) - \psi(\mathbf{F}_s)\|_2^2 \\ &= \frac{1}{n^2} \sum_{i,j} (\varphi(\mathbf{f}_i^s, \mathbf{f}_j^s) - \varphi(\mathbf{f}_i^t, \mathbf{f}_j^t))^2. \end{aligned}$$



人脸识别-模型压缩



1、小网络的设计



2、低秩分解和二值化



3、知识蒸馏



4、整体总结



4、整体总结

数据

第一章：常见训练测试数据库

第六章：常见数据分布

特征提取

第二章：传统手工特征

第三和第四章：深度学习模型

第七章：模型压缩

分类损失

第二章：传统分类器

第四章：深度学习损失函数



4、整体总结

数据

第一章：常见训练测试数据库，比如CASIA-WebFace和LFW，各种测试协议的理解；

第六章：常见数据分布，比如噪声、长尾、无标签数据、监控人脸等。很多特定场景的数据，现有算法处理不是很好，发论文相对比较好。



4、整体总结

特征提取

第二章 传统手工特征: LBP, HOG; PCA, LDA; 字典学习等;

第三和第四章 深度学习模型: DeepID系列定义了深度学习人脸识别的基本流程,
分类网络架构和人脸特征的网络架构;

第七章 模型压缩: 小网络模型的设计, 网络结构搜索。



4、整体总结

分类损失

第二章 传统分类器：卡阈值, Joint Bayesian, SVM等;

第四章 深度学习损失函数: Softmax, Margin-based Softmax (SphereFace, CosFace, ArcFace, ? ?), Mining-Softmax (Hard Example Mining, FocalLoss), Metric Learning (Contrastive Loss, Triplet Loss)。



课程作业

作业：复现知识蒸馏算法

1. 从零训练MobileFaceNet;
2. 从零训练(1/2)MobileFaceNet (所有卷积层kernel数目减半)
3. 知识蒸馏方式训(1/2)MobileFaceNet, 拟合概率和拟合特征两种方式做LFW结果对比



结语

感谢聆听！

Thanks for Listening

