

# Computer Vision

# Papers-1<sub>(2022)</sub>

## AlexNet

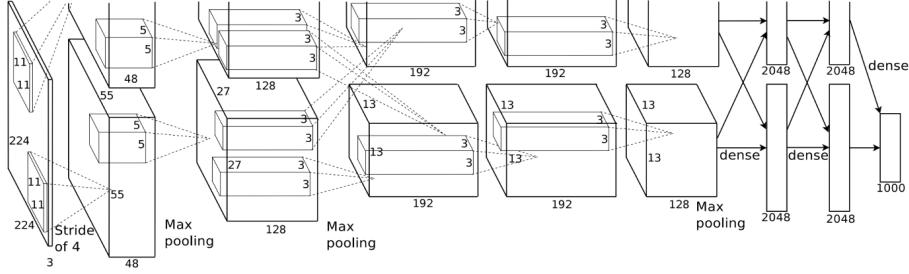


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

The architecture consists of eight layers: five convolutional layers and three fully-connected layers. Due to GPU constraints, it was trained 2 parallel stream on 2 GPUs.

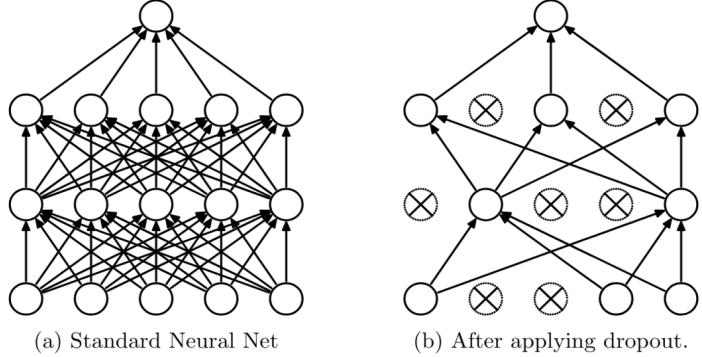
## Dropout

The key idea is to randomly drop units (along with their connections) from the neural network during training. This prevents units from co-adapting too much. Each hidden unit with dropout must learn to work with randomly chosen sample of other unit.

1. Prevent Over-fitting.

2. Approximately combine exponentially many different neural network efficiently.

$$\begin{aligned} r_j^{(l)} &\sim \text{Bernoulli}(p) \\ y^{(l)} &= \mathbf{r}^{(l)} * \mathbf{y}^{(l)} \\ z_i^{(l+1)} &= \mathbf{w}_i^{(l+1)} \mathbf{y}^{(l)} + b_i^{(l+1)} \\ y_i^{(l+1)} &= f(z_i^{(l+1)}) \end{aligned}$$



## VGG-16 & VGG-19

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64	conv3-64	conv3-64	conv3-64	conv3-64
LRN		conv3-64	conv3-64	conv3-64	conv3-64
maxpool					
conv3-128	conv3-128	conv3-128	conv3-128	conv3-128	conv3-128
	conv3-128		conv3-128		conv3-128
maxpool					
conv3-256	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256
conv3-256	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256
		conv1-256	conv1-256	conv1-256	conv1-256
maxpool					
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
		conv1-512	conv1-512	conv1-512	conv1-512
maxpool					
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
		conv1-512	conv1-512	conv1-512	conv1-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

## LRN: Local Response Normalization

Create competition for big activities amongst neuron outputs computed using different kernels.

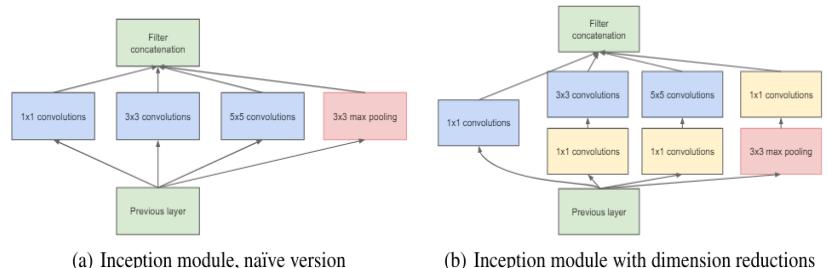
$$b_{x,y}^i = \frac{a_{x,y}^i}{\left( k + \alpha \sum_{j=\max(0, i-n/2)}^{\min(N-1, i+n/2)} (a_{x,y}^i)^2 \right)^\beta}$$

$b_{x,y}^i$  response-normalized activity

$a_{x,y}^i$  Activity of neuron computed by applying kernel  $i$  at position  $(x,y)$  and then applying the ReLU non-linearity  
 $k=2, n=2, \alpha=10^{-4}, \beta=0.75$

## GoogLeNet

In Naive module,  $5 \times 5$  convolutions can be prohibitively expensive on top of a convolutional layer with a large number of filters. The number of output filters equals to the number of filters in the previous stage. The merging of the output of the pooling layer with the outputs of convolutional layers would lead to an inevitable increase in the number of outputs from stage to stage. This leads to the second idea of the proposed architecture:  $1 \times 1$  convolutions are used to compute reductions before the expensive  $3 \times 3$  and  $5 \times 5$  convolutions.



## Batch Normalization

As the training of network continuously changes the parameter of the network, the distribution of activation changes, this is known as Internal Covariance shift.

Training

$$BN(x; \gamma, \beta, \underbrace{x_1, x_2, \dots, x_m}_{\text{mini-batch}}) = \gamma \frac{x - \mu(x_1, \dots, x_m)}{\sqrt{\sigma^2(x_1, \dots, x_m) + \epsilon}} + \beta$$

$$x \in \{x_1, x_2, \dots, x_m\} \quad x_i, \gamma, \beta, \mu, \sigma \in \mathbb{R}^d$$

Inference

$$BN(x; \gamma, \beta) = \gamma \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta$$

$\mu$  and  $\sigma$  can be taken as running averages from training data.

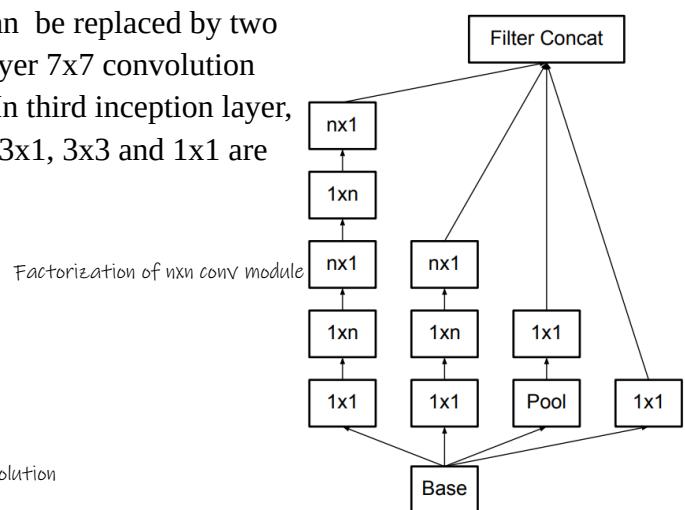
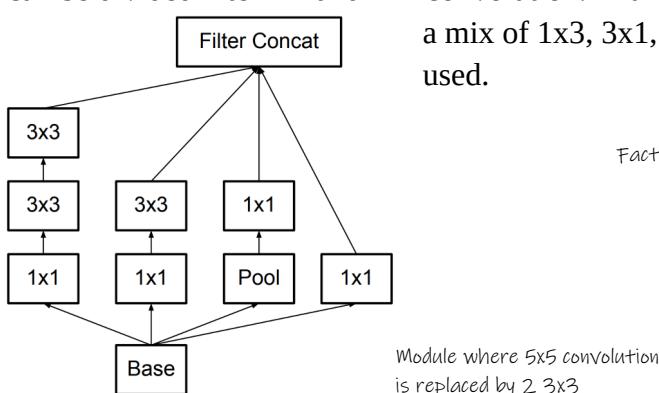
Benefits:

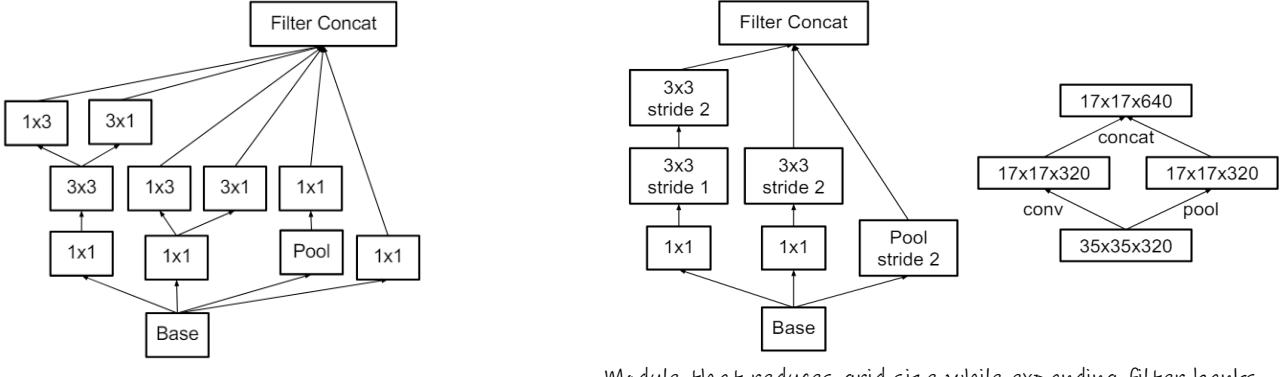
- Higher Learning Rate
- Less sensitive to initialization
- Less Sensitive to activation function
- Regularization Effect
- Preserve Gradients Magnitude (may be)

One more aspect can be look in this, which suggest the Batch Normalization helps in smoothing the Loss function landscape in the weight dimension.

## Inception-v3

In first layer of Inception ,  $5 \times 5$  convolution can be replaced by two  $3 \times 3$  convolution stacked together, in second layer  $7 \times 7$  convolution can be divided into  $1 \times 7$  and  $7 \times 1$  convolution. In third inception layer, a mix of  $1 \times 3$ ,  $3 \times 1$ ,  $3 \times 3$  and  $1 \times 1$  are used.





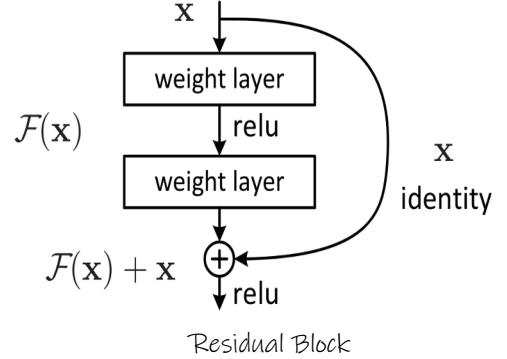
Module that reduces grid size while expanding filter banks

### Label Smoothing Normalization

Replace  $q(k|x)$  with  $\tilde{q}(k|x) = (1 - \epsilon)q(k|x) + \epsilon u(k); k = \{1, 2, \dots, K\}$

### ResNet

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112x112					
				7x7, 64, stride 2		
				3x3 max pool, stride 2		
conv2.x	56x56	$\left[ \begin{matrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{matrix} \right] \times 2$	$\left[ \begin{matrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{matrix} \right] \times 3$	$\left[ \begin{matrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{matrix} \right] \times 3$	$\left[ \begin{matrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{matrix} \right] \times 3$	
conv3.x	28x28	$\left[ \begin{matrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{matrix} \right] \times 2$	$\left[ \begin{matrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{matrix} \right] \times 4$	$\left[ \begin{matrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{matrix} \right] \times 4$	$\left[ \begin{matrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{matrix} \right] \times 4$	$\left[ \begin{matrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{matrix} \right] \times 8$
conv4.x	14x14	$\left[ \begin{matrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{matrix} \right] \times 2$	$\left[ \begin{matrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{matrix} \right] \times 6$	$\left[ \begin{matrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{matrix} \right] \times 6$	$\left[ \begin{matrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{matrix} \right] \times 23$	$\left[ \begin{matrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{matrix} \right] \times 36$
conv5.x	7x7	$\left[ \begin{matrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{matrix} \right] \times 2$	$\left[ \begin{matrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{matrix} \right] \times 3$	$\left[ \begin{matrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{matrix} \right] \times 3$	$\left[ \begin{matrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{matrix} \right] \times 3$	
	1x1			average pool, 1000-d fc, softmax		
FLOPs	$1.8 \times 10^9$	$3.6 \times 10^9$	$3.8 \times 10^9$	$7.6 \times 10^9$	$11.3 \times 10^9$	



One way to look at skip connection is, we are giving model an ability to choose a path, suppose  $F_l(\cdot) = 0$ , then input directly flow to next layer. These skip also helps in vanishing gradient problem.

### Identity Mapping

What we have is

$$y_l = h(x_l) + F(x_l, W_l)$$

$$x_{l+1} = f(y_l)$$

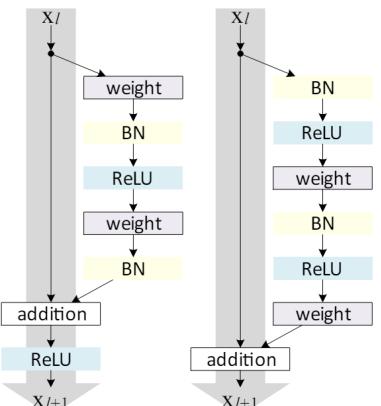
If  $h(x_l) = x_l$  and  $f$  is ReLU, we have original ResNet framework.

If  $f$  and  $h$  are identity, then  $x_{l+1} = x_l + F(x_l, W_l)$

$$x_L = x_l + \sum_{i=1}^{L-1} F(x_i, W_i)$$

One way to look at this is, information from very lower level is flowing to the final layer. (ignoring the effect of BN, MaxPool)

$$\frac{\partial \mathcal{E}}{\partial x_l} = \frac{\partial \mathcal{E}}{\partial x_L} \frac{\partial x_L}{\partial x_l} = \frac{\partial \mathcal{E}}{\partial x_L} \left( 1 + \frac{\partial}{\partial x_l} \sum_{i=l}^{L-1} F(x_i, W_i) \right)$$



(a) original

(b) proposed

During backpropagation, error are propagating from deep layers to initial layers.

If  $h(x_l) = \lambda_l x_l$ , then

$$x_{l+1} = \lambda_l x_l + F(x_l, W_l)$$

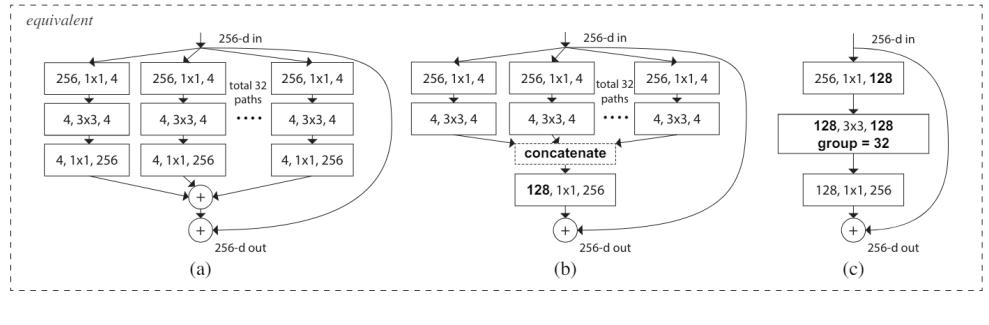
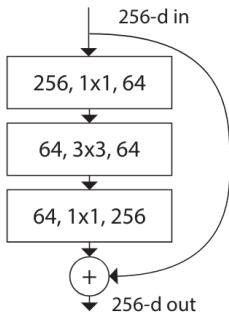
$$x_L = \left( \prod_{i=l}^{L-1} \lambda_i \right) x_l + \sum_{i=l}^{L-1} \left( \prod_{j=i+1}^{L-1} \lambda_j \right) F(x_i, W_i)$$

$$x_L = \left( \prod_{i=1}^{L-1} \lambda_i \right) x_L + \sum_{i=1}^{L-1} \hat{F}(x_i, W_i)$$

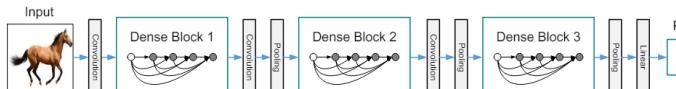
$$\frac{\partial \varepsilon}{\partial x_L} = \frac{\partial \varepsilon}{\partial x_L} \left( \left( \prod_{i=1}^{L-1} \lambda_i \right) + \frac{\partial}{\partial x_L} \sum_{i=1}^{L-1} \hat{F}(x_i, W_i) \right)$$

Now, the first term has lot of multiplication, so its either gonna explode or vanish. Same way if h is nor linear, its going to have product of h', which will produce same issue.

## ResNeXt



## DenseNet



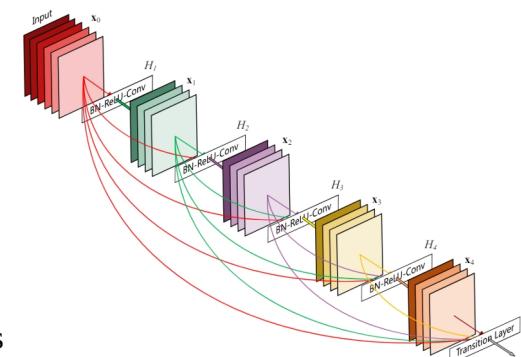
$x_l \rightarrow$  output of  $l^{\text{th}}$  layer

$$x_l = H \left[ \underbrace{x_0, x_1, \dots, x_{l-1}}_{\text{concatenate}} \right]$$

$H(\cdot) \rightarrow$  NonLinear Transformation

$k \rightarrow$  Growth rate of network

if each layer  $H_l$  produces  $k$  feature maps then  $l^{\text{th}}$  layer has  $k_0 + k(l-1)$  input feature maps



## mixup

### Empirical Risk Minimization (ERM)

$f \in F \rightarrow$  function

$X \rightarrow$  Random feature vector

$Y \rightarrow$  random target vector

$P(X, Y) \rightarrow$  Joint Distribution

$l \rightarrow$  loss function

$R(f) = \int l(f(x), y) dP(x, y) \rightarrow$  expected risk

$D = \{(x_i, y_i)\}_{i=1}^N \rightarrow$  Training Data

$R_\delta(f) = \int l(f(x), y) \delta P(x, y) = \frac{1}{n} \sum_{i=1}^N l(f(x_i), y_i) \rightarrow$  Empirical Risk

### Vicinal Risk Minimization (VRM)

$$P_\nu(\tilde{x}, \tilde{y}) = \frac{1}{n} \sum_{i=1}^N \nu(\tilde{x}, \tilde{y} | x_i, y_i)$$

$\nu \rightarrow$  vicinity distribution (Measure the probability of finding virtual feature target ( $\tilde{x}, \tilde{y}$ ) in vicinity of training feature ( $x_i, y_i$ )

$\nu(\tilde{x}, \tilde{y}|x_i, y_i) = N(\tilde{x} - x_i, \sigma^2) \delta(\tilde{y} - y_i) \rightarrow$  Gaussian vicinity

$$R_v(f) = \frac{1}{m} \sum_{i=1}^m l(f(\tilde{x}_i, \tilde{y}_i)) \rightarrow \text{Empirical Vicinal Risk}$$

# mixup

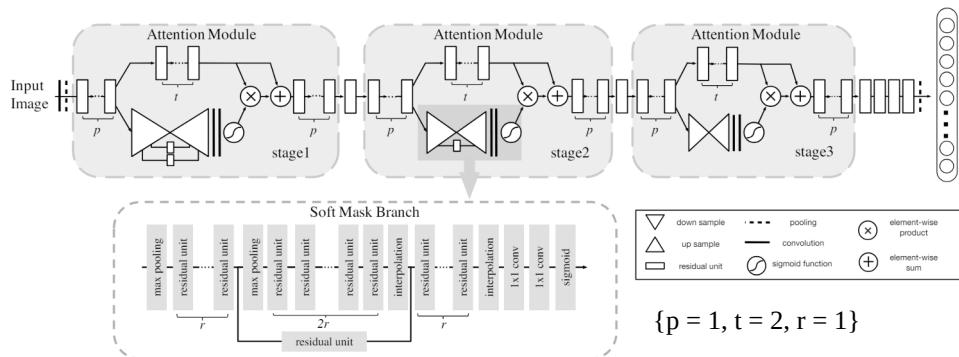
$$\mu(\tilde{x}, \tilde{y} | x_i, y_i) = \frac{1}{n} \sum_{j=1}^n E[\delta(\tilde{x} = \lambda x_i + (1-\lambda)x_j, \tilde{y} = \lambda y_i + (1-\lambda)y_j)]$$

$\lambda \sim Beta(\alpha, \alpha)$  for  $\alpha \in (0, \infty)$

mixup GAN

$$\max_g \min_d \mathbb{E}_{x,z,\lambda} l(d(\lambda x + (1-\lambda)g(z)), \lambda)$$

## Residual Attention Network



## Attention Module

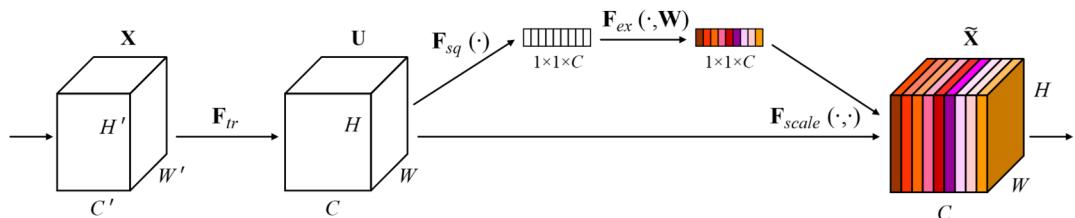
- Mask branch (bottom-up top-down structure)
  - Trunk branch (perform feature processing)

$$H_{i,c}(x) = \underbrace{(1 + M_{i,c}(x))}_{mask} * \underbrace{F_{i,c}(x)}_{trunk}$$

$i \rightarrow$  ranges over all spatial positions

$c \rightarrow$  index of the channel

## Squeeze-and-Excitation Networks



$$X \in \mathbb{R}^{H' \times W' \times C'} \quad U \in \mathbb{R}^{H \times W \times C}$$

$$F_{tr}: X \rightarrow U$$

$V \equiv [v_1, v_2, \dots, v_C]$  learned set of filter kernels

### $\gamma_c$ parameters of the $c^{th}$ filter

$$U = [u_1 \ u_2 \ \dots \ u_n] \quad u_i \in \Re^{H \times W}$$

$$u_c = v_c * X = \sum_{c'} v_c^s * x^s$$

Squeeze

$$z \in \Re^C \quad z_c = F_{sq}(u_c) = \frac{1}{HW} \sum_{i,j} u_c(i,j)$$

Excitation

$$s = F_{ex}(z, W) = \sigma(g(z, W)) = \sigma(W_2 \delta(W_1 z))$$

$$W_1 \in \Re^{r \times C} \quad W_2 \in \Re^{C \times \frac{C}{r}}$$

Scale

$$\tilde{x}_c = F_{scale}(u_c, s_c) = s_c u_c$$

$$\tilde{X} = [\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_C]$$

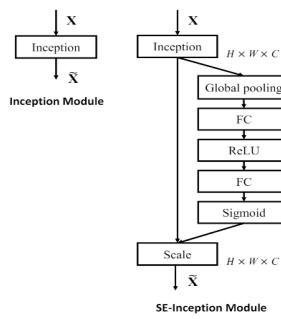


Fig. 2. The schema of the original Inception module (left) and the SE-Inception module (right).

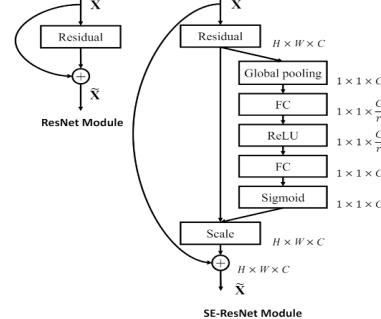


Fig. 3. The schema of the original ResNet module (left) and the SE-ResNet module (right).

## CBAM: Convolutional Block Attention Module

$$F \in \Re^{C \times H \times W} \text{ Intermediate feature map}$$

$$M_c \in \Re^{C \times 1 \times 1} \text{ 1D channel attention output}$$

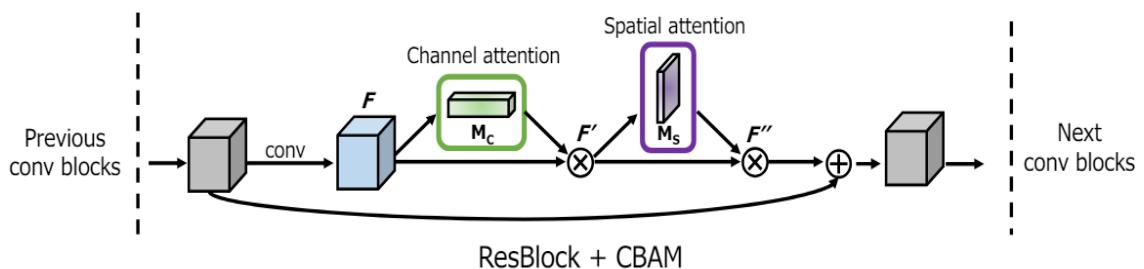
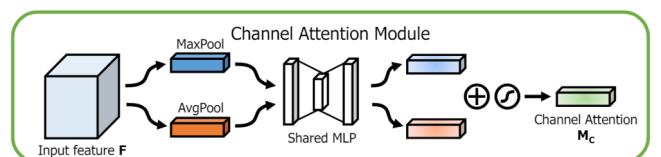
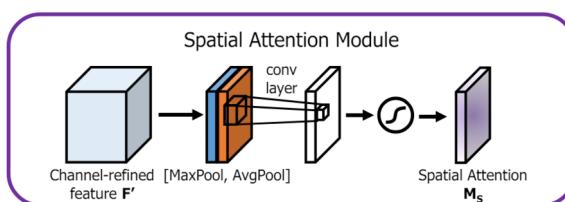
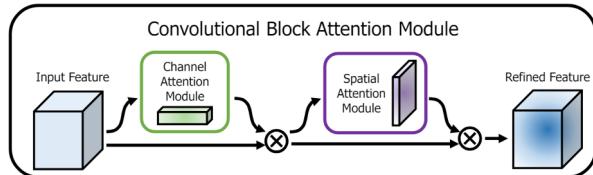
$$M_s \in \Re^{1 \times H \times W} \text{ 2d spatial attention map}$$

$$F' = M_c(F) \odot F$$

$$F'' = M_s(F') \odot F'$$

$$M_s(F) = \sigma(f^{7 \times 7}([AvgPool(F); MaxPool(F)]))$$

$$= \sigma(f^{7 \times 7}([F_{avg}^s; F_{max}^s]))$$



## Random Erasing



### Algorithm 1: Random Erasing Procedure

```

Input : Input image  $I$ ;  

        Image size  $W$  and  $H$ ;  

        Area of image  $S$ ;  

        Erasing probability  $p$ ;  

        Erasing area ratio range  $s_l$  and  $s_h$ ;  

        Erasing aspect ratio range  $r_1$  and  $r_2$ .  

Output: Erased image  $I^*$ .  

Initialization:  $p_1 \leftarrow \text{Rand}(0, 1)$ .  

1 if  $p_1 \geq p$  then  

2    $I^* \leftarrow I$ ;  

3   return  $I^*$ .  

4 else  

5   while True do  

6      $S_e \leftarrow \text{Rand}(s_l, s_h) \times S$ ;  

7      $r_e \leftarrow \text{Rand}(r_1, r_2)$ ;  

8      $H_e \leftarrow \sqrt{S_e \times r_e}$ ,  $W_e \leftarrow \sqrt{\frac{S_e}{r_e}}$ ;  

9      $x_e \leftarrow \text{Rand}(0, W)$ ,  $y_e \leftarrow \text{Rand}(0, H)$ ;  

10    if  $x_e + W_e \leq W$  and  $y_e + H_e \leq H$  then  

11       $I_e \leftarrow (x_e, y_e, x_e + W_e, y_e + H_e)$ ;  

12       $I(I_e) \leftarrow \text{Rand}(0, 255)$ ;  

13       $I^* \leftarrow I$ ;  

14    end  

15  end  

16 end  

17 end

```

## Spatial Transformers Network

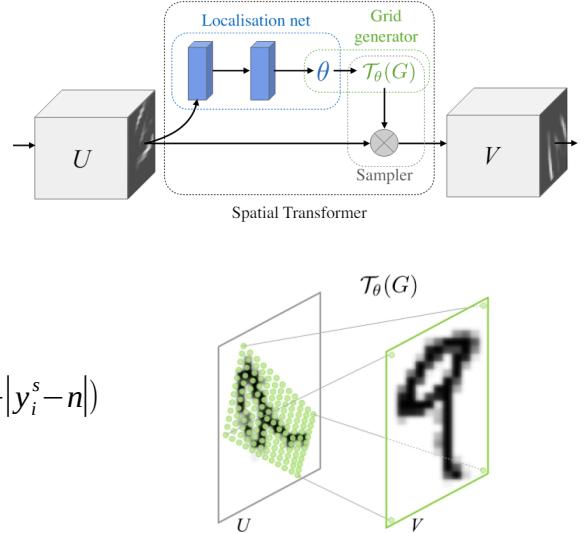
Three main differential blocks:

1. Localization Network
2. Grid Generator
3. Sampler

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = T(G_i) = A_\theta \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$

$$V_i^c = \sum_n^H \sum_m^W U_{nm}^c \max(0, 1 - |x_i^s - m|) \max(0, 1 - |y_i^s - n|)$$

Bilinear Sampling



## Capsule Networks

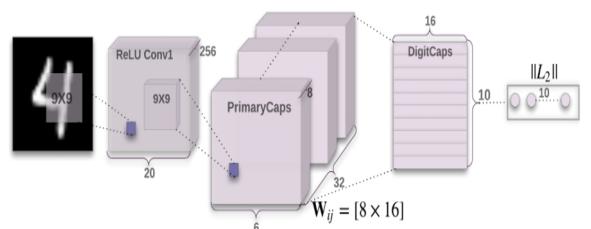
Neuron  $\rightarrow$   $\underbrace{\text{capsules}}_{\text{group of neurons}} \rightarrow$  layers

$s_j \rightarrow$  total input to capsule j

$v_j \rightarrow$  vector output of capsule j

$\|v_j\| \rightarrow$  probability that entity represented by the capsule

$\frac{v_j}{\|v_j\|} \rightarrow$  properties of entity



## Procedure 1 Routing algorithm.

```

1: procedure ROUTING( $\hat{\mathbf{u}}_{j|i}$ ,  $r$ ,  $l$ )
2:   for all capsule  $i$  in layer  $l$  and capsule  $j$  in layer  $(l + 1)$ :  $b_{ij} \leftarrow 0$ .
3:   for  $r$  iterations do
4:     for all capsule  $i$  in layer  $l$ :  $\mathbf{c}_i \leftarrow \text{softmax}(\mathbf{b}_i)$             $\triangleright \text{softmax}$  computes Eq. 3
5:     for all capsule  $j$  in layer  $(l + 1)$ :  $\mathbf{s}_j \leftarrow \sum_i c_{ij} \hat{\mathbf{u}}_{j|i}$ 
6:     for all capsule  $j$  in layer  $(l + 1)$ :  $\mathbf{v}_j \leftarrow \text{squash}(\mathbf{s}_j)$             $\triangleright \text{squash}$  computes Eq. 1
7:     for all capsule  $i$  in layer  $l$  and capsule  $j$  in layer  $(l + 1)$ :  $b_{ij} \leftarrow b_{ij} + \hat{\mathbf{u}}_{j|i} \cdot \mathbf{v}_j$ 
return  $\mathbf{v}_j$ 

```

$$v_j = \frac{\|s_j\|^2}{\underbrace{1 + \|s_j\|^2}_{\text{squashing function}}} \frac{s_j}{\|s_j\|}$$

$$s_j = \sum_i c_{ij} u_{j|i}$$

$c_{ij} \rightarrow$  Coupling coefficients

$$u_{j|i} = W_{ij} u_i$$

$u_{j|i} \rightarrow$  prediction vector from capsules in layer below

$u_i \rightarrow$  output of capsule in layer below

$$c_{ij} = \frac{\exp(b_{ij})}{\sum_k \exp(b_{ik})}$$

$b_{ij} \rightarrow$  Log prob. that capsule  $i$  should be coupled with capsule  $j$

Loss function

$$L_k = T_k \max(0, m^+ - \|v_k\|)^2 + \lambda (1 - T_k) \max(0, \|v_k\| - m^-)^2$$

$T_k = 1$  iff a digit class of class  $k$  is present

$$m^+ = 0.9 \text{ and } m^- = 0.1$$

$$\lambda = 0.5$$

## Vision Transformers

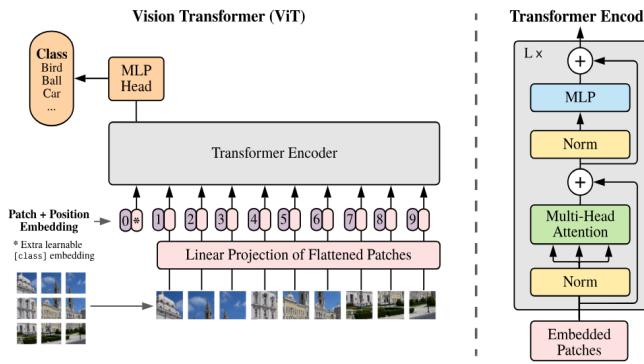


Image Patches  $\equiv$  Tokens in NLP  
 $x \in \mathbb{R}^{H \times W \times C} \rightarrow$  image  
 $x_p \in \mathbb{R}^{N \times (P^2 C)} \rightarrow$  sequence of flattened 2D patches  
 $P \times P \rightarrow$  resolution of each image patch  
 $N = \frac{HW}{P^2} \rightarrow$  resulting number of patches  
 $D \rightarrow$  latent vector size  
 $z_0 = [x_{class}; x_p^1 E; x_p^2 E; \dots; x_p^N E] + E_{pos}$   
 $E \in \mathbb{R}^{(P^2 C) \times D}, E_{pos} \in \mathbb{R}^{(N+1) \times D}$

$$z_l' = MSA(LN(z_{l-1}) + z_{l-1}), l=1, \dots, L$$

$$z_l = MLP(LN(z_l')) + z_l, l=1, \dots, L$$

$$y = LN(z_L^0)$$

LN  $\rightarrow$  Layer Normalization

MSA  $\rightarrow$  Multihead self-attention

$$\begin{aligned}
z &\in \mathbb{R}^{N \times D} \\
[q, k, v] &= zU_{qkv}, U_{qkv} \in \mathbb{R}^{D \times 3D_h} \\
A &= \text{softmax}(qk^T / \sqrt{D_h}), A \in \mathbb{R}^{N \times N} \\
SA(z) &= Av \\
MSA(z) &= [SA_1(z), SA_2(z), \dots, SA_k(z)]U_{msa}, U_{msa} \in \mathbb{R}^{kD_h \times D} \\
D_h &= D/k
\end{aligned}$$

## ML-Mixer

$X \in \mathbb{R}^{S \times C}$  → real-time input table

$S \rightarrow$  No. of image patches

$C \rightarrow$  Hidden Dimension

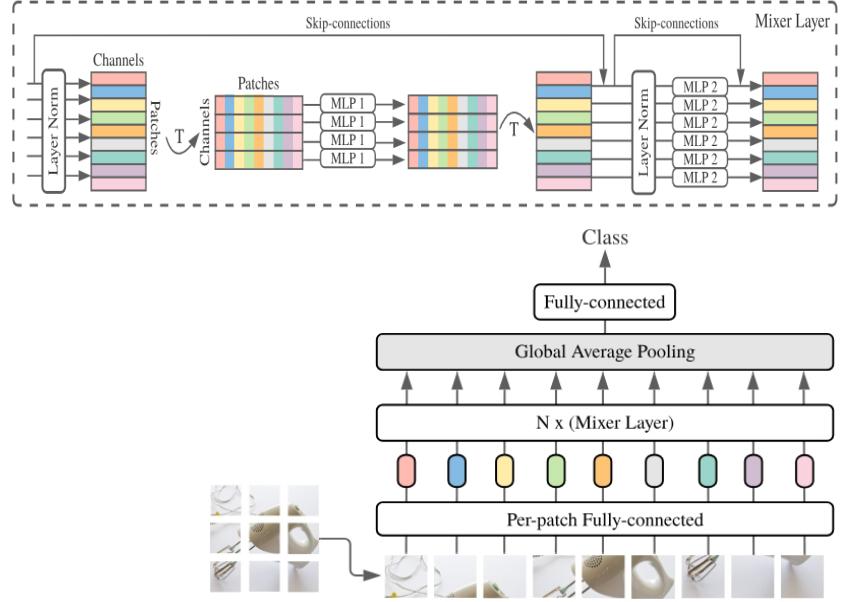
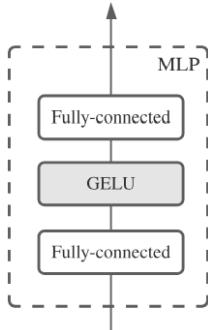
Token Mixing

$$U_{*,i} = X_{*,i} + W_2 \sigma(W_1 \text{LN}(X)_{*,i}) \quad i=1, \dots, C$$

Channel Mixing

$$Y_{j,*} = U_{j,*} + W_4 \sigma(W_3 \text{LN}(U)_{j,*}) \quad j=1, \dots, S$$

$\sigma \rightarrow$  GELU (Gaussian Error Linear unit)



## Normalizer-Free ResNets

Gradient Clipping

$L \rightarrow$  Loss

$\theta \rightarrow$  model parameters

$G = \frac{\partial L}{\partial \theta}$  → Gradient vector

$$G \rightarrow \begin{cases} \lambda \frac{G}{\|G\|} & \text{if } \|G\| > \lambda \\ G & \text{otherwise} \end{cases}$$

$\lambda \rightarrow$  Clipping threshold

Adaptive Gradient Clipping:  $\frac{\|W^l\|_F}{\|G^l\|_F}$  measure of gradient step change with respect to original weight.

$W^l \in \mathbb{R}^{N \times M}$  → weight matrix of  $l^{th}$  layer

$G^l \in \mathbb{R}^{N \times M}$  → Gradient wrt  $W^l$

$\Delta W^l = -h G^l$ ,  $h \rightarrow$  learning rate

$$\frac{\|\Delta W^l\|_F}{\|W^l\|_F} = h \frac{\|G^l\|_F}{\|W^l\|_F}$$

If  $\frac{\|\Delta W^l\|_F}{\|G^l\|_F}$  is large, we expect training to be unstable! Update i<sup>th</sup> row of Matrix G,  $G_i^l$

$$G_i^l \rightarrow \begin{cases} \lambda \frac{\|W_i^l\|_F^*}{\|G_i^l\|_F^*} G_i^l & \text{if } \frac{\|G_i^l\|_F^*}{\|W_i^l\|_F^*} > \lambda, \|W_i\|_F^* = \max(\|W_i\|_F, \epsilon) \\ G_i^l & \text{otherwise} \end{cases}$$

Normalizer Free ResNet

$$h^{l+1} = h^l + \alpha f^l \left( \frac{h^l}{\beta^l} \right) \rightarrow \text{Residual Block}$$

We want  $\text{Var}(f(z)) = \text{Var}(z)$

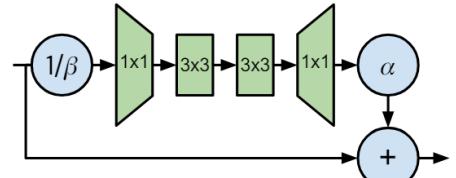
$$\beta^l = \sqrt{(\text{Var}(h^l))} ; \alpha = 0.2$$

Scale Weight Standardization

$$\hat{W}_{ij} = \frac{W_{ij} - \mu_i}{\sqrt{N} \sigma_i}$$

$$\mu_i = \frac{1}{N} \sum_j W_{ij}$$

$$\sigma_i^2 = \frac{1}{N} \sum_j (W_{ij} - \mu_i)^2$$



The activation functions are also scaled by non-linearity specific scalar gain  $\gamma$ ,  $\gamma = \sqrt{2/1 - 1/\pi}$  for ReLU. Combination of  $\gamma$ -scaled activation and scaled weight Standardization is variance preserving.

## Knowledge Distillation

Knowledge is transferred to distilled model by training it on a transfer set and using soft target distribution for each case in the transfer set that is produced by using the cumbersome model with a high temperature in its softmax. The same high temperature is used when training distilled model, but after it has trained it uses a temperature of 1.

$$q_i = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}; T \rightarrow \text{Temperature}$$

## Pruning

Dropout Ratio Adjustment

Do: Original Dropout rate

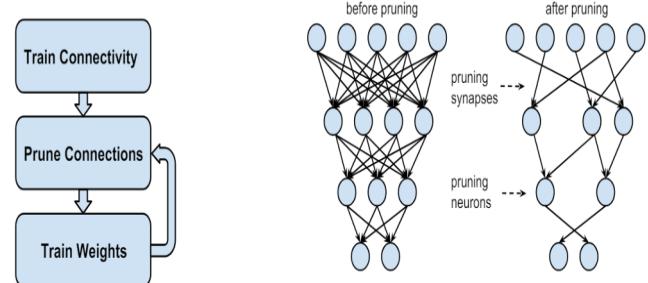
Dr: Dropout rate during training

$C_{ir}$ : #connections in layer i after retraining.

$C_{io}$ : #connections in layer i for original network.

$$D_r = D_o \sqrt{\frac{C_{ir}}{C_{io}}}$$

$C_{io}$  and  $C_{ir}$  vary quadratically with no. Of neurons!



## Deep Compression

Network Pruning

All connections with weights below a threshold are removed from the network.

- Compresses Sparse Row
- Compressed Sparse Column

No. of Elements =  $2a + n + 1$

a: No. of Non-zero elements

n: No. Of Rows

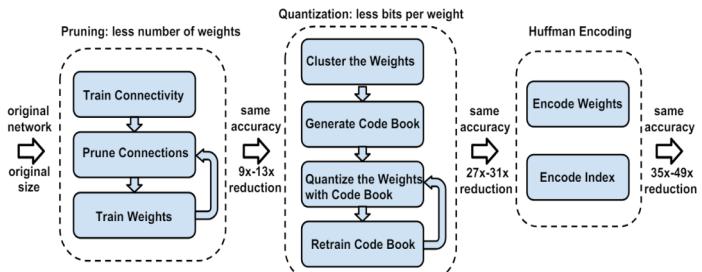
To compress further, store the index difference instead of the absolute position, and encode this difference in 8 bits for conv layer and 5 bits for fc layer.

Trained Quantization and weight sharing

Suppose 4x4 weight matrix. These weights are quantized into 4 bins (by colour). We store the index of weight into the weight matrix. During update, all gradients are grouped by colour, multiplied by learning rate and subtracted form centroid.

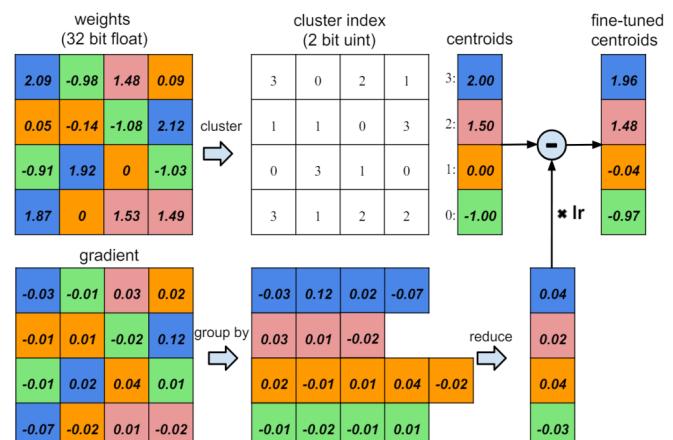
Compression rate

$$r = \frac{nb}{n \log_2(k) + kb}$$



Span Exceeds 8=2^3																
idx	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
diff		1			3								8		3	
value	3.4				0.9							0		1.7		

Filler Zero



k: #clusters;  $\log_2(k)$ : bits required to encode index; n: #weights; b: bits representing each connection  
Use k-means cluster to get he centroid. Original weights  $W = \{w_1, \dots, w_n\}$ , Clusters  $C = \{c_1, \dots, c_k\}$

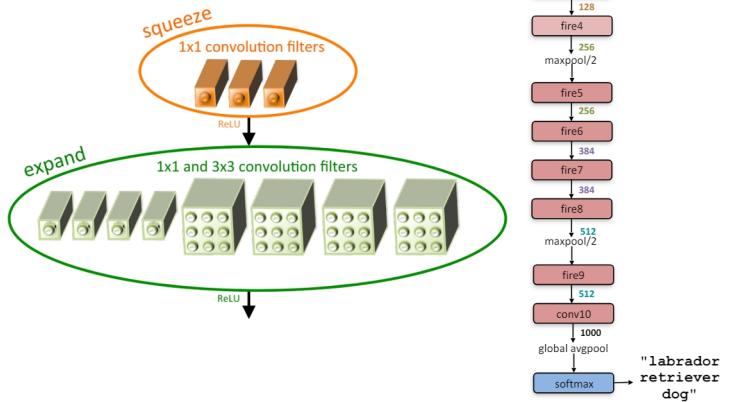
Minimize the within-cluster sum of squares (WCSS)

$$\arg \min_C \sum_{i=1}^k \sum_{w \in c_i} |w - c_i|$$

## SqueezeNet

Strategies:

- Replace 3x3 filter with 1x1
- Decrease number of input channels to 3x3 filters
- Downsample late in network so that convolution have large activation map



### Fire Module

$s_{1 \times 1}$  → Number of  $1 \times 1$  filters in squeeze layer

$e_{1 \times 1}$  → Number of  $1 \times 1$  filters in expand layer

$s_{3 \times 3}$  → Number of  $3 \times 3$  filters in expand layer

Also used skip connections over module.

## XNOR-Net

Two variations:

**Binary weights** here the weight matrix is binary tensor.

Estimate weight matrix by  $W \approx \alpha B$ , B is binary filter  $B \in \{-1, +1\}^{c \times w \times h}$ ,  $\alpha \in \mathbb{R}^+$ .

Convolution can be approximated as

$I * W \approx (I \oplus B) \alpha$ .  $\oplus$  indicates convolution without multiplication.

For Estimation, assume W, B are vector in  $\mathcal{R}^n$ , where  $n = c \times w \times h$  and solve

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

$$\alpha^*, B^* = \arg \min_{\alpha, B} J(B, \alpha)$$

$$J(B, \alpha) = \alpha^2 B^T B - 2\alpha W^T B + W^T W, \quad B^T B = n, \quad W^T W \text{ is constant.}$$

$$B^* = \arg \max_B W^T B \quad s.t. B \in \{-1, +1\}^n$$

	Network Variations	Operations used in Convolution	Memory Saving (Inference)	Computation Saving (Inference)	Accuracy on ImageNet (AlexNet)
Standard Convolution	Real-Value Inputs 0.11 - 0.21 ... - 0.34 - 0.25 - 0.61 ... 0.52  Real-Value Weights 0.12 0.2 ... 0.64 0.23 0.53 ... 0.88	+ , - , $\times$	1x	1x	55.67
Binary Weight	Real-Value Inputs 0.11 - 0.21 ... - 0.34 - 0.25 - 0.61 ... 0.52  Binary Weights 1 - 1 ... - 1 1 - 1 ... 1	+ , -	$\sim 32x$	$\sim 2x$	55.8
BinaryWeight Binary Input (XNOR-Net)	Binary Inputs 1 - 1 ... - 1 - 1 ... 1  Binary Weights 1 - 1 ... - 1 1 - 1 ... 1	XNOR , bitcount	$\sim 32x$	$\sim 58x$	44.2

This optimization can be solved  $B^* = \text{sign}(W)$  and  $\alpha^* = \frac{W^T B^*}{n}$ , replacing  $B^*$  with  $\text{sign}(W)$ , we get

$$\alpha^* = \frac{1}{n} \|W\|_1$$

### Training

Given,  $W_t$ , calculate  $\alpha_t = \frac{1}{n} \|W_t\|_1$  and  $B_t = \text{sign}(W_t)$

$$\tilde{W}_t = \alpha_t B_t$$

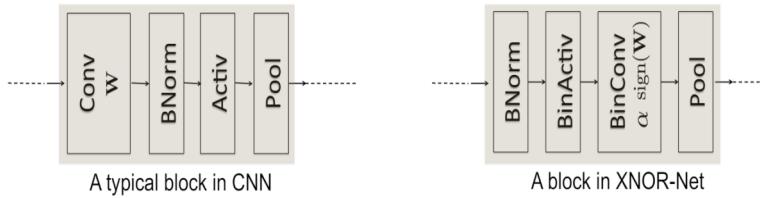
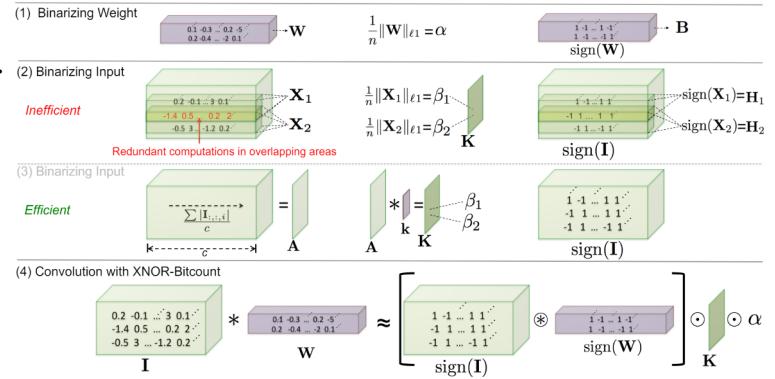
$$W_{t+1} = W_t - \eta \nabla_{\tilde{W}_t} J$$

**XNOR-Networks** here elements of both input feature map and weight matrix are binary.

For efficiency you compute A, which is average of feature map across channels.

Convolve A with k to get K, where  $k_{ij} = 1/(w \times h) \forall i, j$ . Hence we can approximate convolution of Image I with Weight W as

$$I * W \approx (\text{sign}(I) * \text{sign}(W)) \odot K \alpha$$



BinActiv computes k and I, BinConv does the convolution shown in above figure

## MobileNets

During Normal convolution, we have

$$F \in \mathcal{R}^{D_F \times D_F \times M} \rightarrow \text{Input feature maps}$$

$$G \in \mathcal{R}^{D_F \times D_F \times N} \rightarrow \text{Output feature maps}$$

$$K \in \mathcal{R}^{D_K \times D_K \times M \times N}$$

$$G_{k,l,n} = \sum_{i,j,m} K_{i,j,m,n} F_{k+i-1, l_j-1, m}$$

Computational cost:  $D_K^2 MN D_F^2$

Depthwise separable convolution: For every input channel, produces same number of output channels, by using 2D filter of size  $D_K \times D_K$

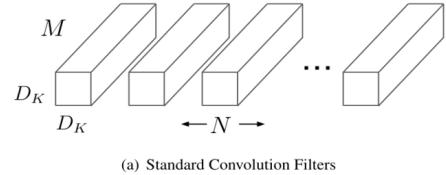
$$\hat{G}_{k,l,m} = \sum_{i,j} K_{i,j,m} F_{k+i-1, l_j-1, m}$$

$$\hat{K} \in \mathcal{R}^{D_K \times D_K \times M}$$

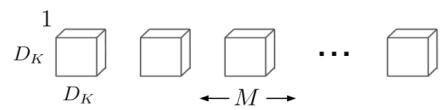
Computational cost:  $D_K^2 M D_F^2$

Pointwise convolution: Uses  $1 \times 1$  convolution to produce linear combination of input feature maps

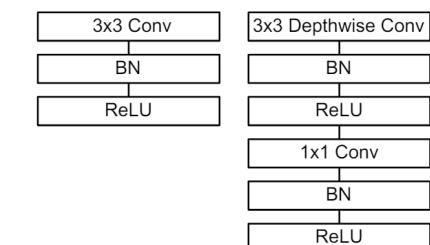
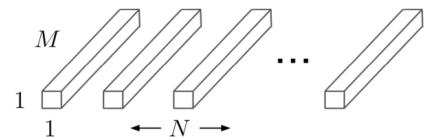
Total computation Cost  $D_K^2 M D_F^2 + MND_F^2$



(a) Standard Convolution Filters



(b) Depthwise Convolutional Filters



## ShuffleNet

channel shuffle

$g \rightarrow$  Number of groups

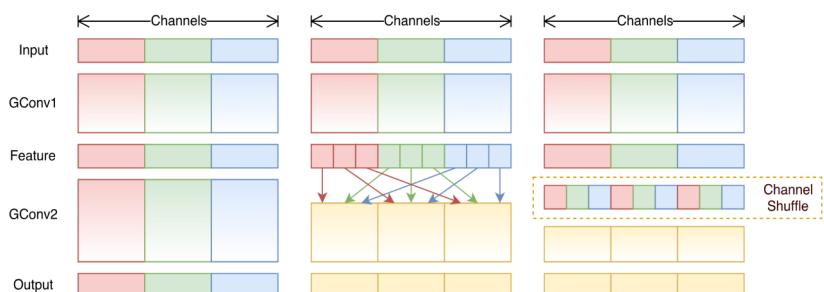
$n \rightarrow$  Number of channels per group

$gn \rightarrow$  number of channels

reshape to  $(g, n)$

transpose

flatten



## Adversarial Examples

$x \in \mathcal{R}^m \rightarrow$  an input image

$f: \mathcal{R}^m \rightarrow \{1, 2, \dots, k\}$

For a given  $x \in \mathcal{R}^m$ , we try to solve

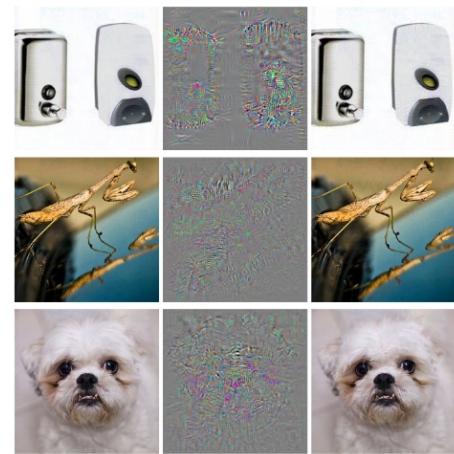
$$\begin{aligned} & \text{minimize } \|r\|_2 \\ \text{subject to } & \begin{cases} f(x+r) = l \\ x+r \in [0,1]^m \end{cases} \end{aligned}$$

label  $l \in \{1, \dots, k\}$

Informally,  $x+r$  is closest image to  $x$  classified to  $l$  by  $f$ . The minimizer is not unique and problem is non-trivial. So we try to solve below problem

$$\begin{aligned} & \text{minimize } c|r| + \text{loss}_f(x+r, l) \\ \text{subject to } & x+r \in [0,1]^m \end{aligned}$$

Box constraint L-BFGS for optimization



Left: correctly predicted sample  
Middle: Noise magnified to 10x  
Right: Left Image + Noise, predicted as "ostrich, Struthio, camelus"

## Fast Gradient Sign Method

$$\underbrace{\tilde{x}}_{\text{Adversarial Input}} = \underbrace{x}_{\text{Input}} + \eta$$

The classifier should assign same class to  $\tilde{x}$  and  $x$ , as long as  $\|\eta\|_\infty < \epsilon$ .

$$w^T \tilde{x} = w^T x + w^T \eta \quad \text{Activation grows by } w^T \eta$$

$$\eta = \text{sign}(w)\epsilon \Rightarrow \|\eta\|_\infty = \epsilon$$

$$w^T \eta = \sum_{i=1}^n w_i \eta_i = \sum_{i=1}^n |w_i| \epsilon = \epsilon n m$$

'm' is average magnitude of element of weight vector

"Fast gradient sign" method of generating adversarial example

$\theta \rightarrow$  parameters of model

$x \rightarrow$  input to the model

$y \rightarrow$  target associated with  $x$   $J(\theta, x, y) \rightarrow$  cost

$$\eta = \text{sign}(\nabla_x J(\theta, x, y))$$

$$\begin{array}{ccc} x & + .007 \times & \text{sign}(\nabla_x J(\theta, x, y)) \\ \text{"panda"} & & \text{"nematode"} \\ 57.7\% \text{ confidence} & & 8.2\% \text{ confidence} \\ & & \frac{x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))}{\text{"gibbon"} \\ & & 99.3 \% \text{ confidence}} \end{array}$$

## Adversarial Training

$$\tilde{J}(\theta, x, y) = \alpha J(\theta, x, y) + (1-\alpha) J(\theta, x + \text{sign}(\nabla_x J(\theta, x, y))\epsilon, y)$$

$$\tilde{J}(\theta, x, y) \approx \alpha J(\theta, x, y) + (1-\alpha) \epsilon \text{sign}(\nabla_x J(\theta, x, y))^T \nabla_x J(\theta, x, y)$$

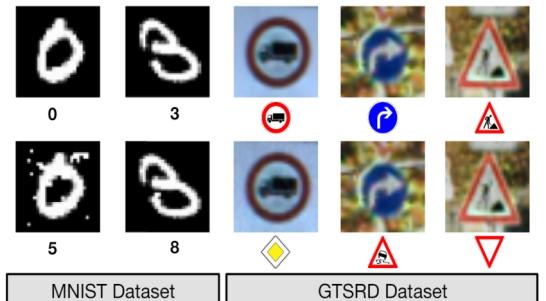
$$= \alpha J(\theta, x, y) + (1-\alpha) \epsilon \|\nabla_x J(\theta, x, y)\|_1$$

## Black-Box Attacks

The adversary has no information about the structure or parameters of DNN, and does not have access to any large training data. The adversary is only capable of assigning labels assigned by the DNN for chosen input.

Attack Strategy: Train a local substitute DNN with "synthetic" dataset

$O \rightarrow$  Oracle (Targeted DNN)



$x \rightarrow \text{input}$   
 $\tilde{O}(x) \rightarrow \text{label} = \arg \max_j O_j(x)$   
 $O(x) \rightarrow \text{Probability vector}$   
 $x^* = x + \delta_x \rightarrow \text{Adversarial Example}$   
 $\delta_x = \arg \min \{z : \tilde{O}(x_z) \neq \tilde{O}(x)\}$   
 $F \rightarrow \text{Substitute Network (learn similar decision boundary to } O)$   
 $J_F(x) = \left[ \frac{\partial F_j(x)}{\partial x_i} \right]_{i,j} \rightarrow \text{model's Jacobian Matrix}$   
 $x + \lambda \text{ sign}(J_F(x)[\tilde{O}(x)])$

---

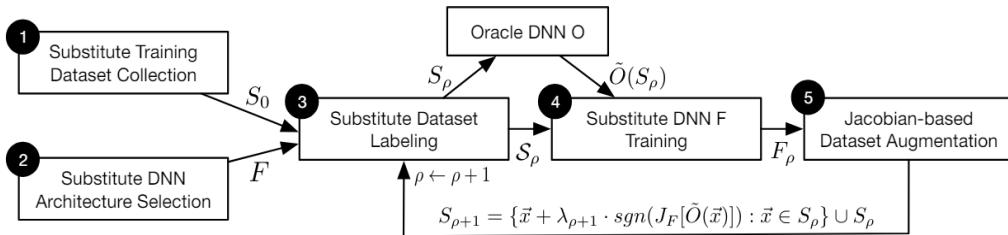
**Algorithm 1 - Substitute DNN Training:** for oracle  $\tilde{O}$ , a maximum number  $\max_\rho$  of substitute training epochs, a substitute architecture  $F$ , and an initial training set  $S_0$ .

---

**Input:**  $\tilde{O}, \max_\rho, S_0, \lambda$

- 1: Define architecture  $F$
- 2: **for**  $\rho \in 0 .. \max_\rho - 1$  **do**
- 3:   // Label the substitute training set
- 4:    $D \leftarrow \{(\vec{x}, \tilde{O}(\vec{x})) : \vec{x} \in S_\rho\}$
- 5:   // Train  $F$  on  $D$  to evaluate parameters  $\theta_F$
- 6:    $\theta_F \leftarrow \text{train}(F, D)$
- 7:   // Perform Jacobian-based dataset augmentation
- 8:    $S_{\rho+1} \leftarrow \{\vec{x} + \lambda \cdot \text{sgn}(J_F[\tilde{O}(\vec{x})]) : \vec{x} \in S_\rho\} \cup S_\rho$
- 9: **end for**
- 10: **return**  $\theta_F$

---



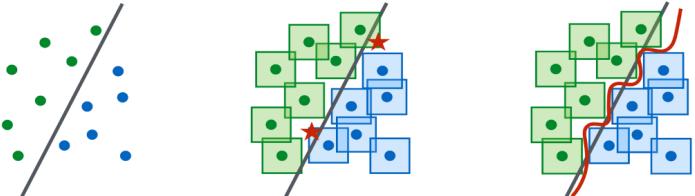
## Robust Optimization

Empirical Risk  $\min_\theta \rho(\theta) = \mathbb{E}_{(x,y) \sim D} \left[ \max_{\delta \in S} L(\theta, x + \delta, y) \right]$   $S \subset \mathbb{R}^d$  Set of allowed perturbation

Fast Gradient sign method,  $L_\infty$  bounded adversary  $x + \epsilon \text{ sign}(\nabla_x L(\theta, x, y))$

Projected Gradient Ascent  $x^{t+1} = \prod_{x \in S} (x^t + \alpha \text{ sign}(\nabla_x L(\theta, x, y)))$

A conceptual illustration of standard vs. adversarial decision boundaries. Left: A set of points that can be easily separated with a simple decision boundary. Middle: The simple decision boundary does not separate the  $\ell_\infty$ -balls. Right: Separating the  $\ell_\infty$ -balls requires a significantly more complicated decision boundary.



## One pixel Attack

$f \rightarrow$  target Image classifier

$x = (x_1, \dots, x_n) \rightarrow$  original image correctly classified as class  $t$

$f_t(x) \rightarrow$  probability of  $x$  belonging to class  $t$

$e(x) = (e_1, \dots, e_n) \rightarrow$  Additive adversarial perturbation

Black Box Attack

$$\begin{aligned} & \max_{e^*(x)} f_{adv}(x + e(x)) \\ & \text{s.t. } \|e(x)\|_0 \leq d, \text{ here } d=1 \end{aligned}$$

Candidate solution  $\{(\underbrace{x_1, y_1, r_1, g_1, b_1}_{\chi_1}), \dots, (x_d, y_d, r_d, g_d, b_d)\}$

Initially,  $x_i, y_i \sim \underbrace{U(1, 32)}_{\text{CIFAR-10}}$  or  $\underbrace{U(1, 127)}_{\text{ImageNet}}$  and  $r_i, g_i, b_i \sim \mathcal{N}(\mu=128, \sigma=127)$

$\chi_i(g+1) = \chi_j(g) + \lambda (\chi_k(g) - \chi_l(g))$   $j \neq k \neq l \rightarrow$  random number  $\lambda = 0.5 \rightarrow$  Scale Parameter

In CIFAR-10, paper did targeted attack while in ImageNet they did untargeted attack.

Fitness function = prob. Label of target class (CIFAR-10)

= prob. Label of the true class (ImageNet)

fitness function > 90%  $\Rightarrow$  stop (CIFAR-10 targeted attack)

fitness function < 5%  $\Rightarrow$  stop (ImageNet, non-targeted attack)

max  $\neq$  iter = 100