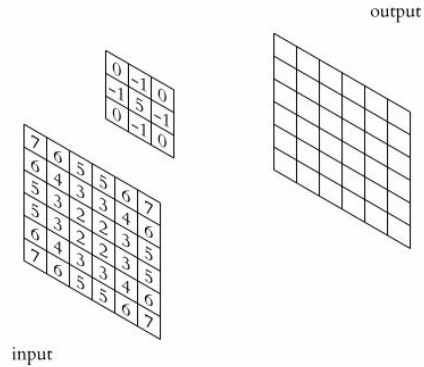
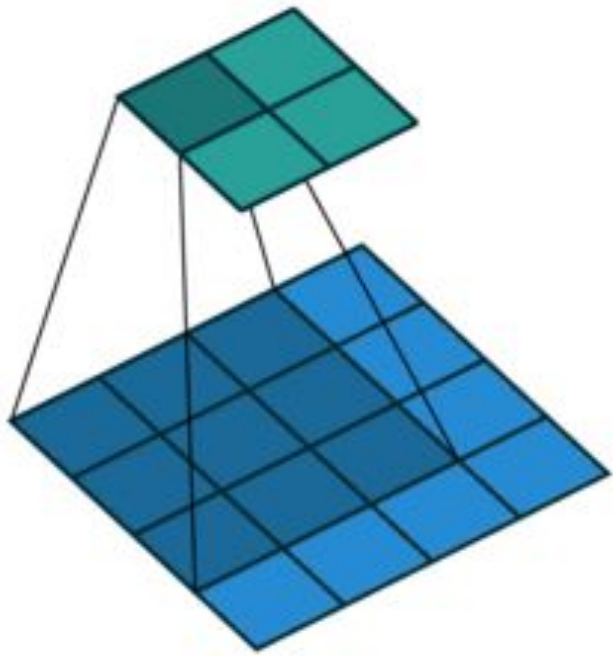


---

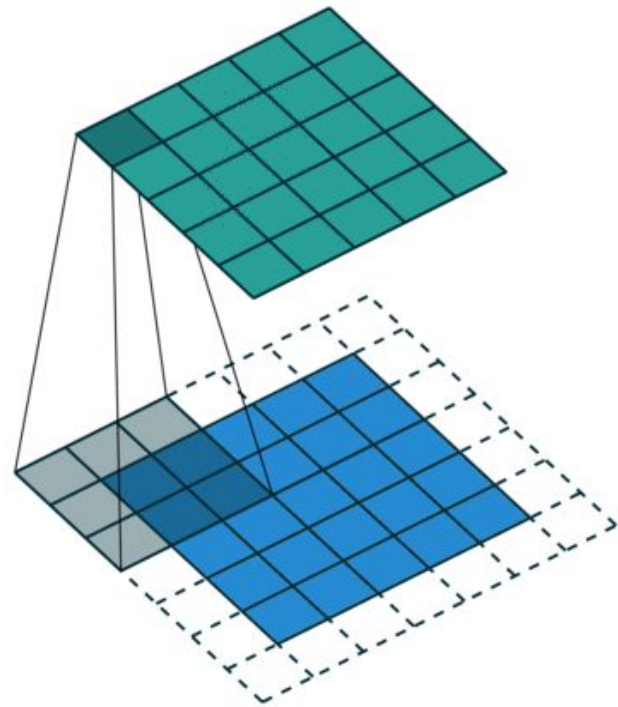
# ANN: ...trainable layers...



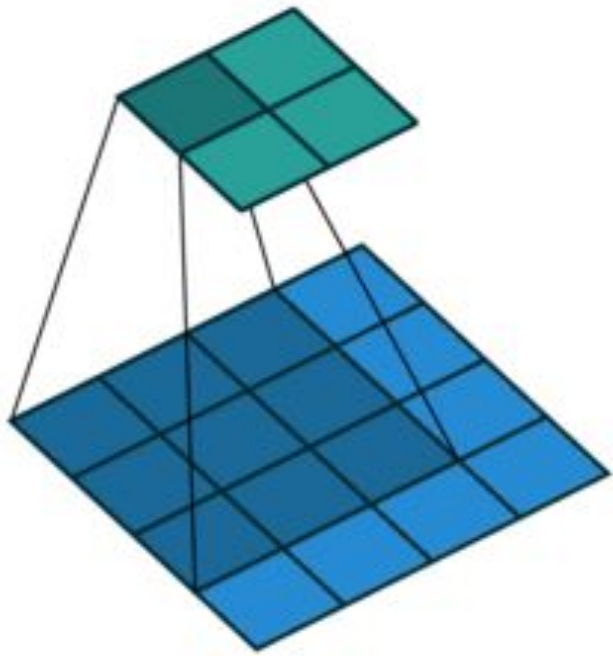
- **convolution** every input is connected to the output by a fixed kernel
  - kernel is (a)symmetrical (e.g.  $3 \times 3$ ,  $3 \times 1$ )
  - some input can be skipped (strides  $> 1$  and/or dilation  $> 1$ )
  - input can be 1D, 2D, 3D, ...  $\times$  channels
  - output channel numbers determined by number of filters
  - only one set of kernel weights (computationally efficient)
  - convolutional 1D layer with 1 kernel == dense layer
  - $\text{inputs} \cdot \text{kernel} \cdot \text{activation} = \text{output}$
  - $\text{inputs} \cdot \text{kernel} \cdot \text{activation} + \text{bias} = \text{output}$
  - output size  $> |<$  input size e.g.  $[128, 16, 16, 32] \Rightarrow [128, 8, 8, 32]$



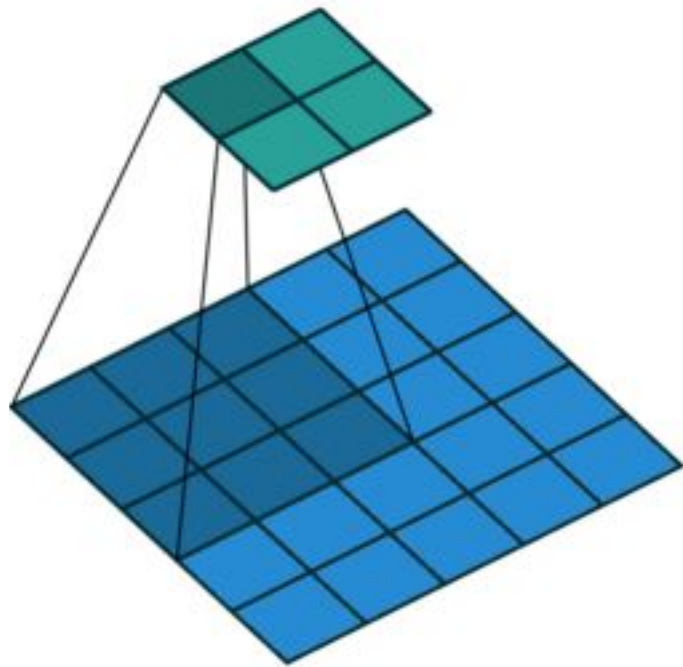
convolution  
(‘valid’ padding, strides 1)



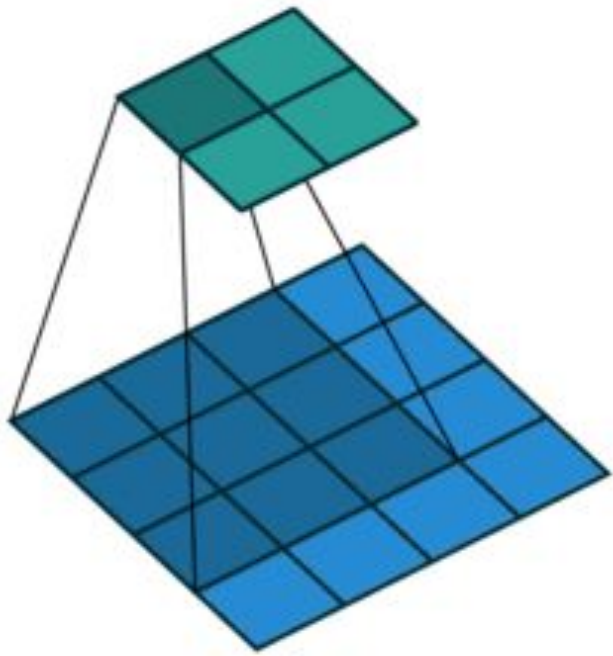
convolution  
(‘same’ padding, strides 1)



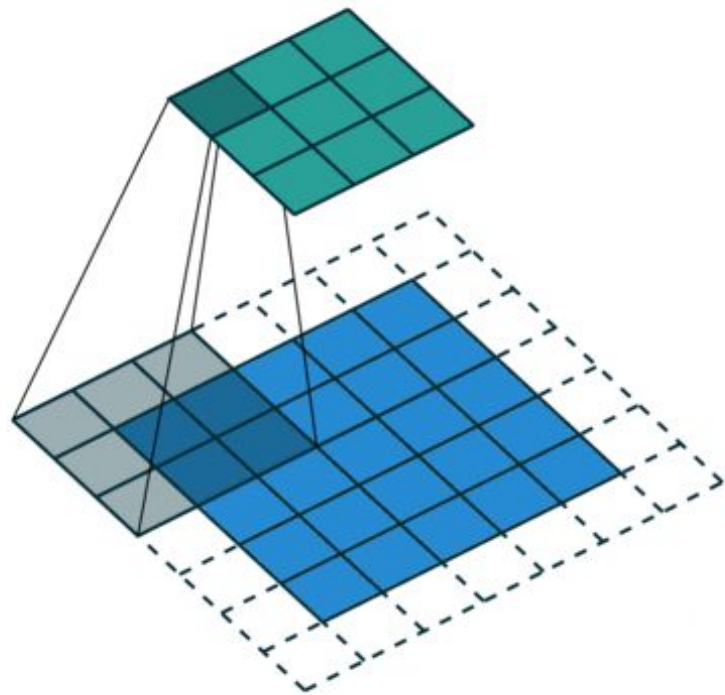
convolution  
(‘valid’ padding, strides 1)



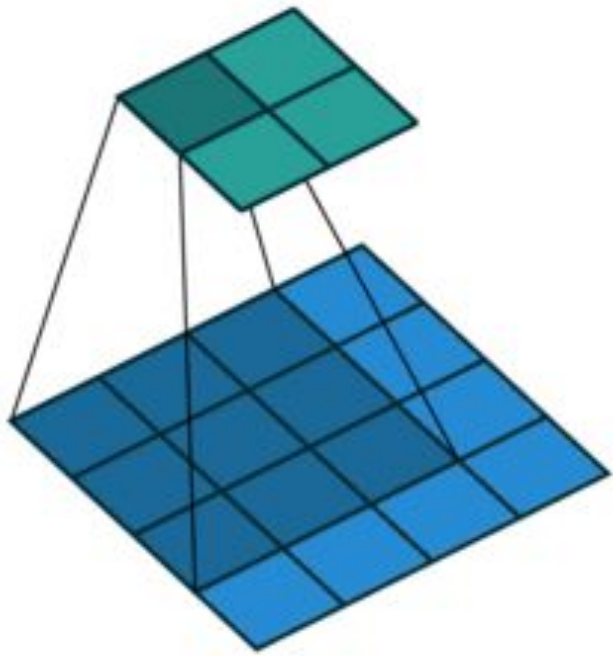
convolution  
(‘valid’ padding, strides 2)



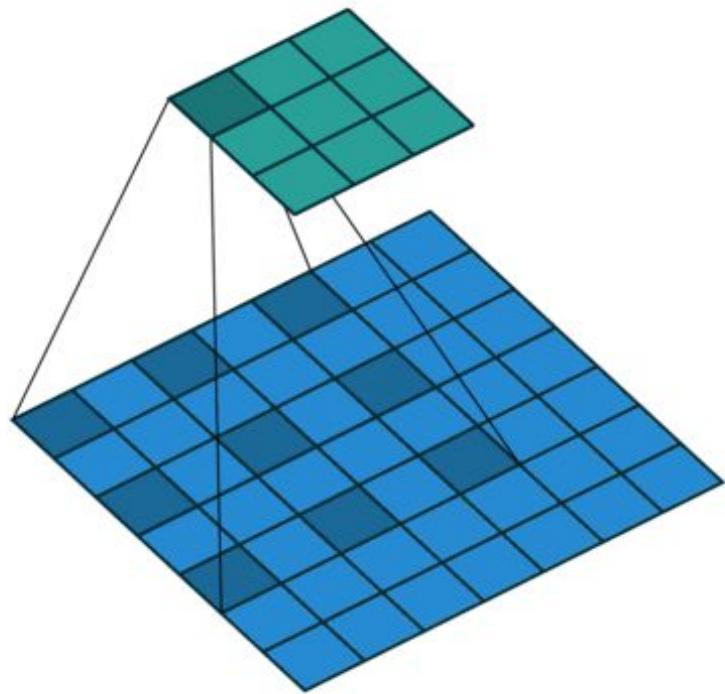
convolution  
(‘valid’ padding, strides 1)



convolution  
(‘same’ padding, strides 2)

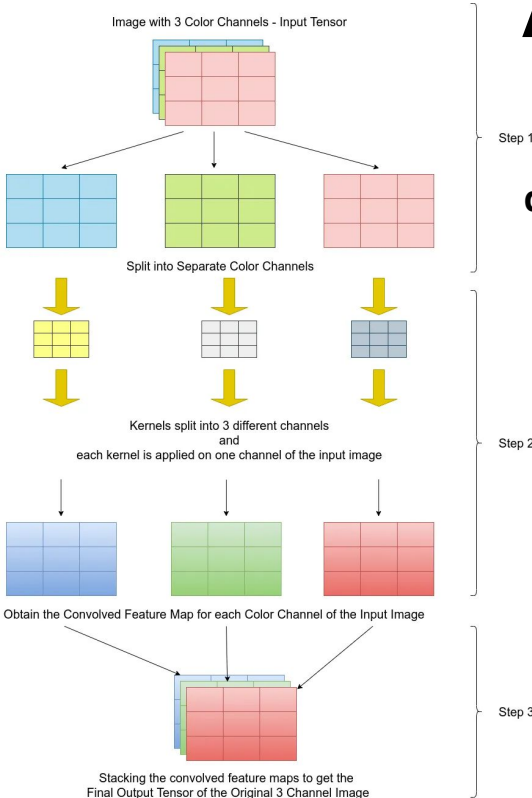


convolution  
(‘valid’ padding, strides 1, dilation 1)



convolution  
(‘valid’ padding, strides 1, dilation 2)

# ANN: ...trainable layers...



## depthwise convolution

separate convolution of channels

one kernel per channel

[0] separate channels

[1] convolute each channel

[2] concatenate channels

=> many fewer multiplications

faster on CPU, but slower on GPU

output size  $\geq$  input size

e.g.  $[128, 16, 16, 32] \Rightarrow [128, 8, 8, 32]$

# ANN: ...trainable layers...

## separable convolution

depthwise + standard convolution

one kernel per channel

[0] separate channels

[1] convolute each channel

[2] concatenate channels

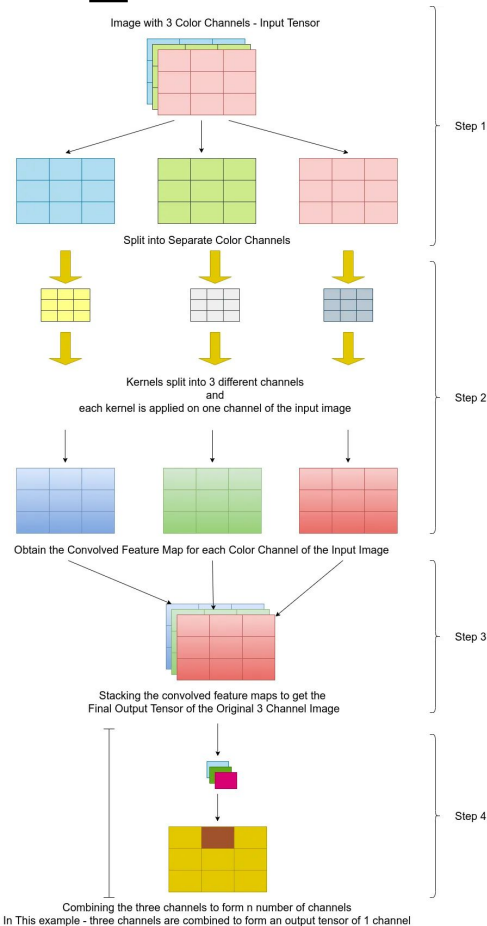
[3] 1×1 convolute

=> many fewer multiplications

faster on CPU, but slower on GPU

output size  $\geq$  input size

e.g.  $[128, 16, 16, 32] \Rightarrow [128, 8, 8, 32]$



In This example - three channels are combined to form an output tensor of 1 channel

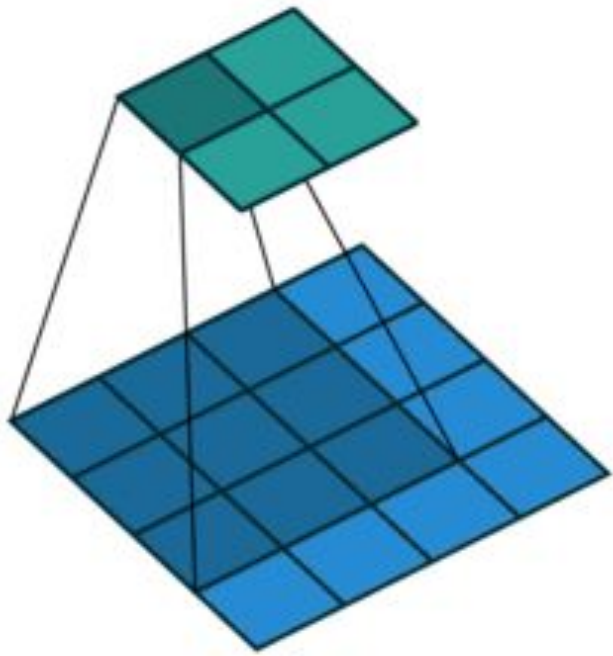
---

# ANN: ...trainable layers...

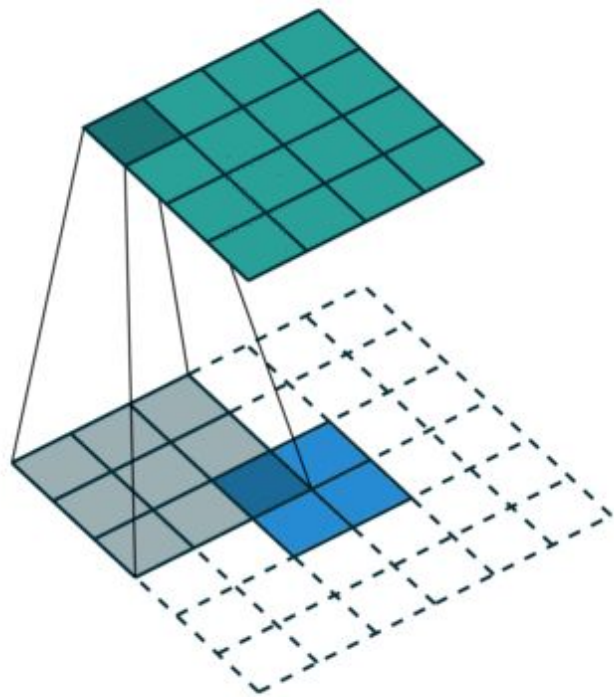
**deconvolution**     a.k.a. transposed convolution  
increases tensor time/height/width  
a learnable upsampling  
output size  $\geq$  input size  
e.g. [128, 16, 16, 32]  
=> [128, 32, 32, 32]

---





convolution  
(‘valid’ padding, strides 1)



deconvolution  
(‘valid’ padding, strides 1)

---

# ANN: ...trainable layers...

## **LSTM**

Long Short-Term Memory

stores data from the last  $n$  inputs

hidden state and cell state act as 'memory'

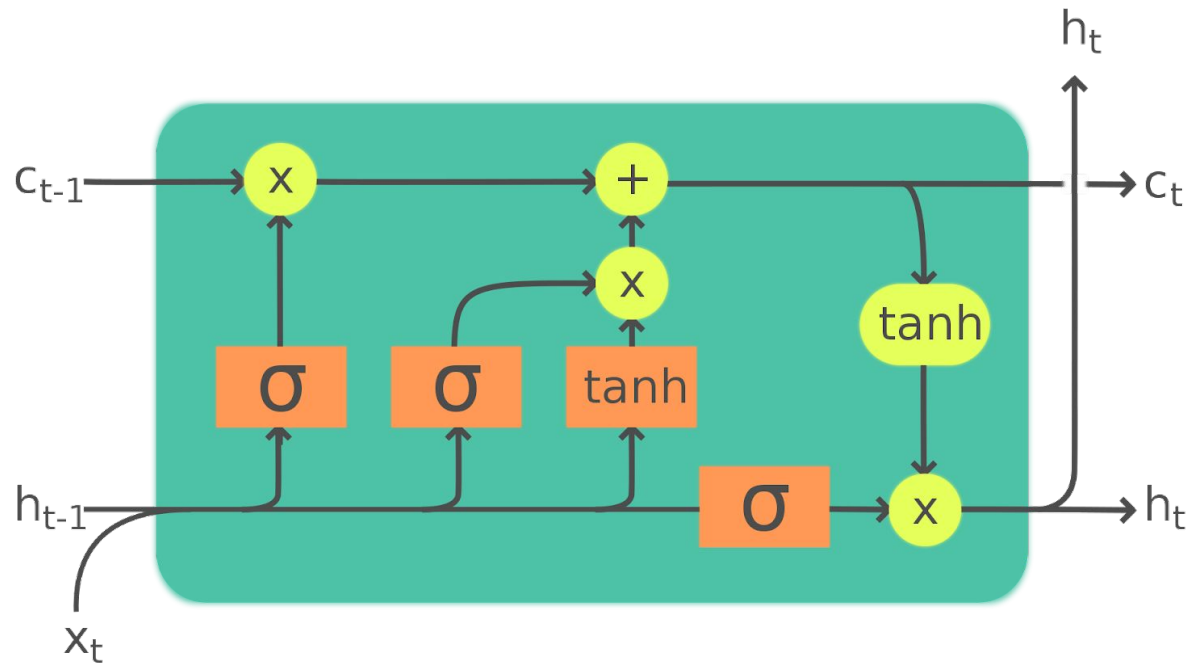
input gate, output gate, and forget gate

used singly or in a bidirectional configuration

upweight 'important' elements

downweight 'unimportant' elements

---



---

# ANN: ...trainable layers...

## GRU

Gated Recurrent Unit

a simplified LSTM

stores data from the last  $n$  inputs

hidden state acts as 'memory'

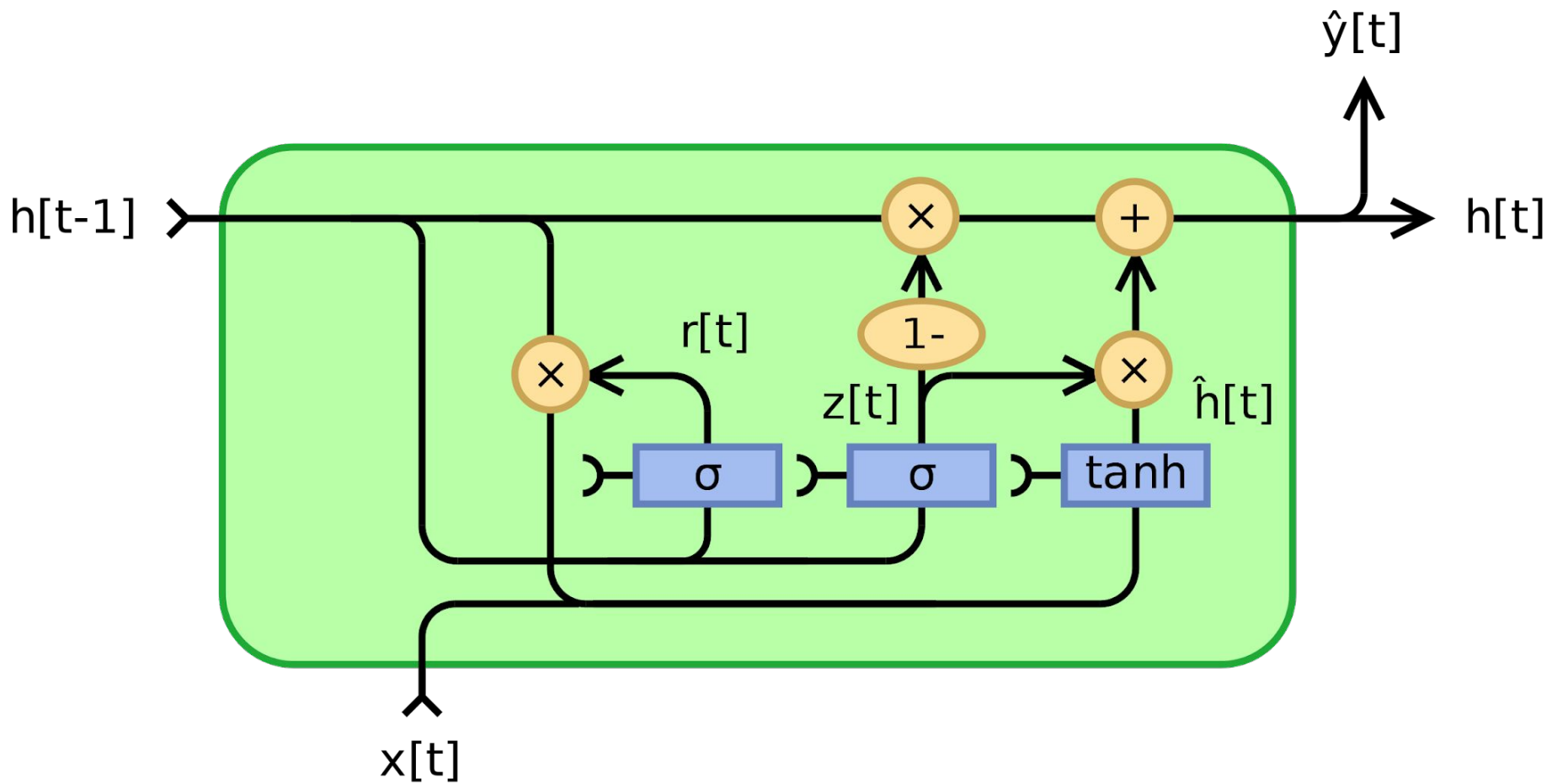
update gate and reset gate

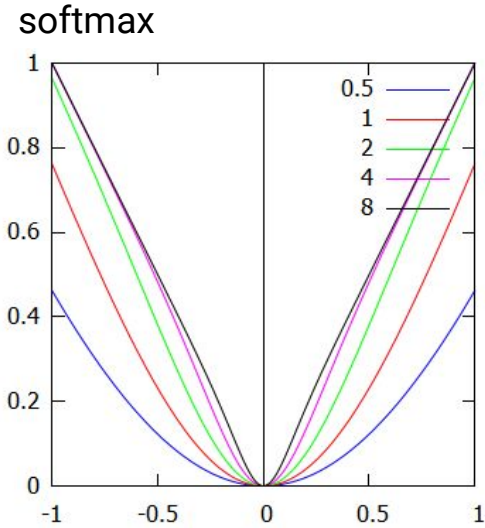
used singly or in a bidirectional configuration

upweight 'important' elements

downweight 'unimportant' elements

---





## ANN: ...trainable layers

**output** typically a dense layer with activation

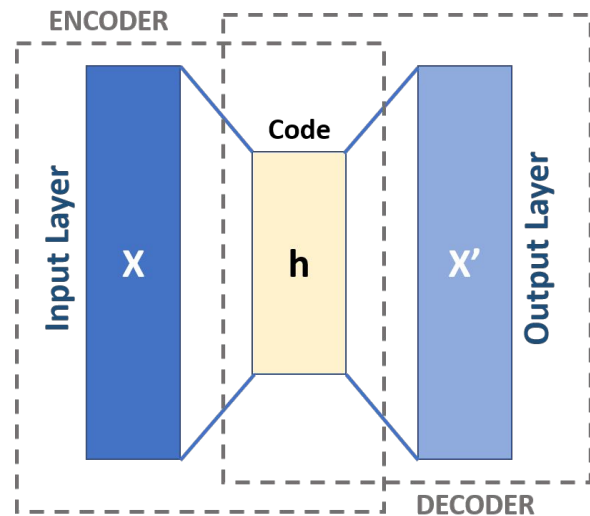
- regression =  $\tanh [-1, +1]$  or sigmoid  $[0, 1]$
- classification = softmax  $[0, 1]$  (sum to 1)
- multiclass classification = sigmoid  $[0, 1]$

often preceded by a global average pooling layer

---

# ANN: block patterns...

auto encoders/bottlenecks



minimum of three layers: large => small => large

size reduction forces compression

important information is retained

less important information and noise is lost

---

# ANN: ...block patterns...

residual connection (a.k.a. skip connection)

allows for deeper networks

signal is lost in the Back Propagation (BP) algorithm

not an issue with Direct Feedback Alignment (DFA)

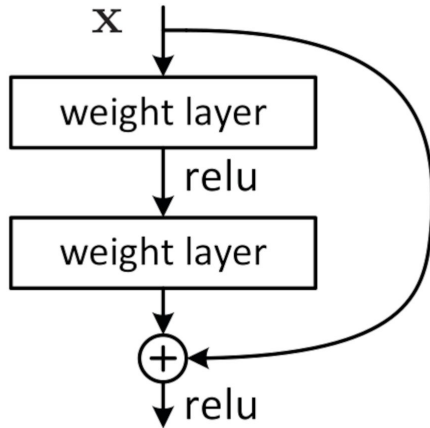
can differentially weight 'important' features

structure:

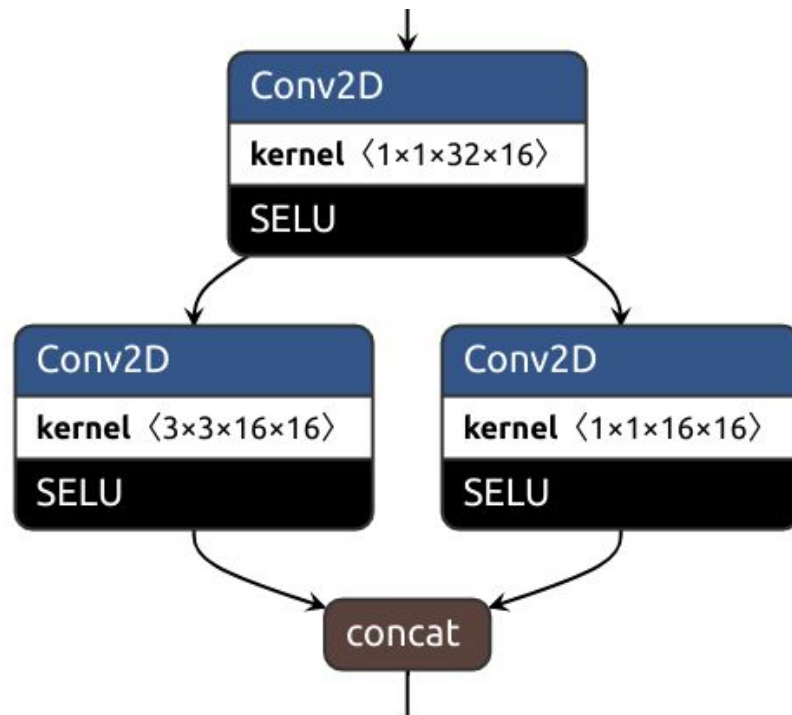
[0] copy data into two (or more) paths

[1] process each path separately (e.g. convolutions)

[2] recombine paths (usually addition or concatenation)







---

## ANN: ...block patterns...

inverted residual blocks

minimum of two layers: expand => contract

size change forces information recoding and filtering

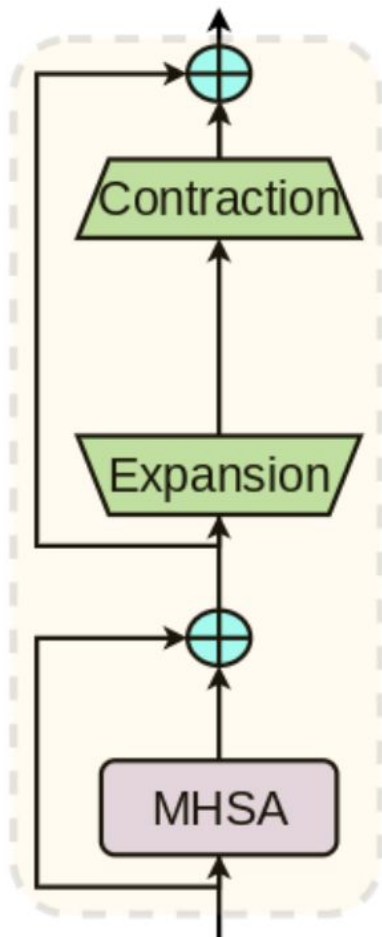
structure:

[0] split or copy data into two (or more) paths

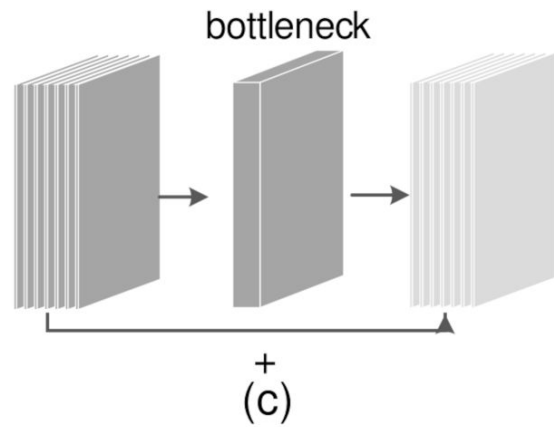
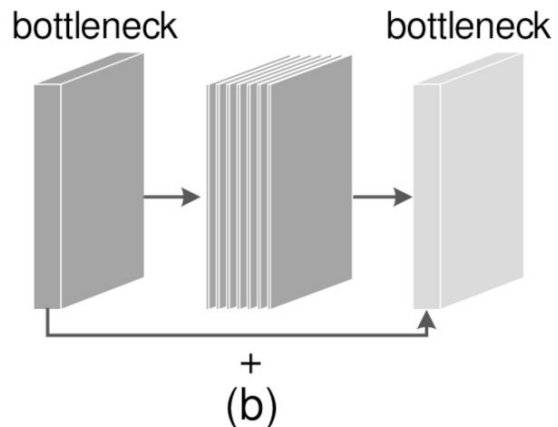
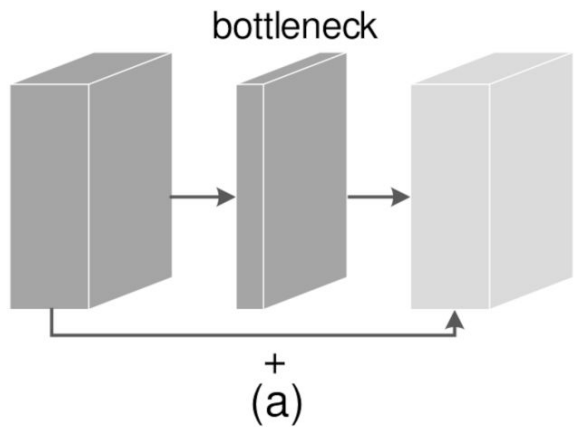
[1] process (or not) each path separately

[2] recombine paths (usually addition)

---



MLP inverted bottleneck (Vaswani et al. 2017 fide Srinivas et al. 2101.11605; <https://arxiv.org/abs/2101.11605>)



---

## ANN: ...block patterns...

distributed convolutions

designed to increase the receptive field

more efficient than a single large-kernel convolution

structure:

[0] split or copy data into two (or more) paths

[1] convolute each path separately

[2] recombine paths (usually concatenation)

---

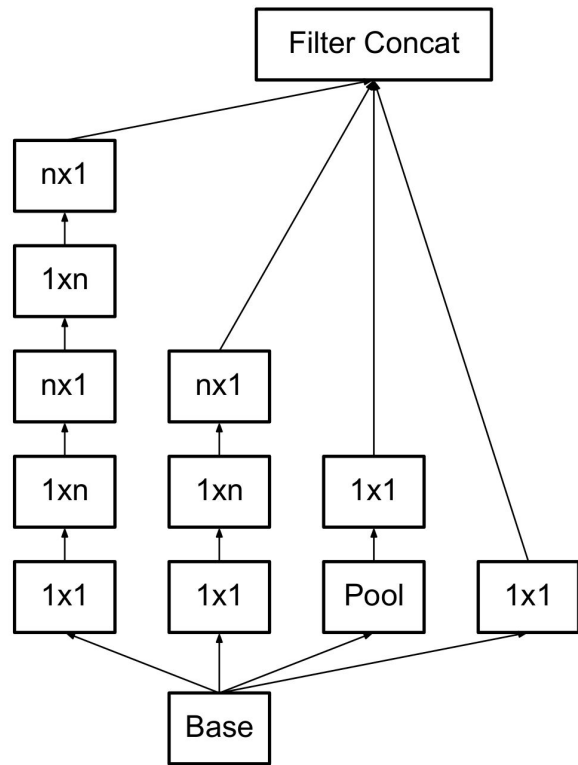


Figure 6. Inception modules after the factorization of the  $n \times n$  convolutions. In our proposed architecture, we chose  $n = 7$  for the  $17 \times 17$  grid. (The filter sizes are picked using principle 3)

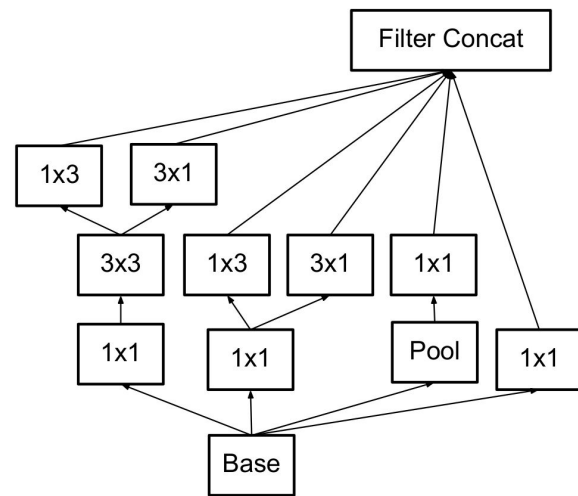
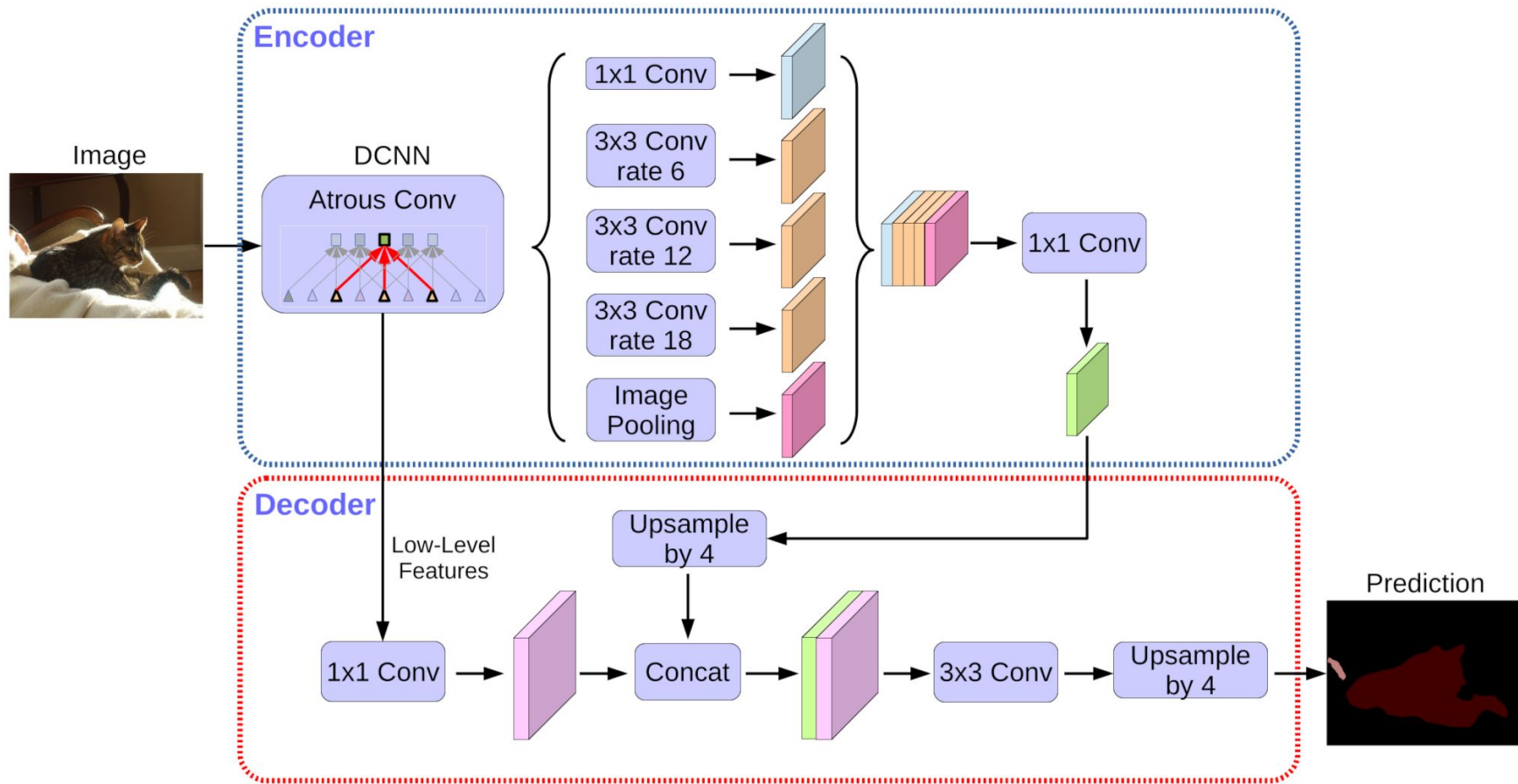


Figure 7. Inception modules with expanded the filter bank outputs. This architecture is used on the coarsest ( $8 \times 8$ ) grids to promote high dimensional representations, as suggested by principle 2 of Section 2. We are using this solution only on the coarsest grid, since that is the place where producing high dimensional sparse representation is the most critical as the ratio of local processing (by  $1 \times 1$  convolutions) is increased compared to the spatial aggregation.



---

# ANN: ...block patterns...

(self) attention

designed to differentially weight 'important' features

(usually) more efficient than convolutions or residual blocks alone

structure:

[0] split or copy data into two (or more) paths (heads)

[1] process each path separately (e.g. convolutions)

[2] activate with [0, 1] output (e.g. softmax, sigmoid)

[3] recombine paths (usually with multiplication)

---



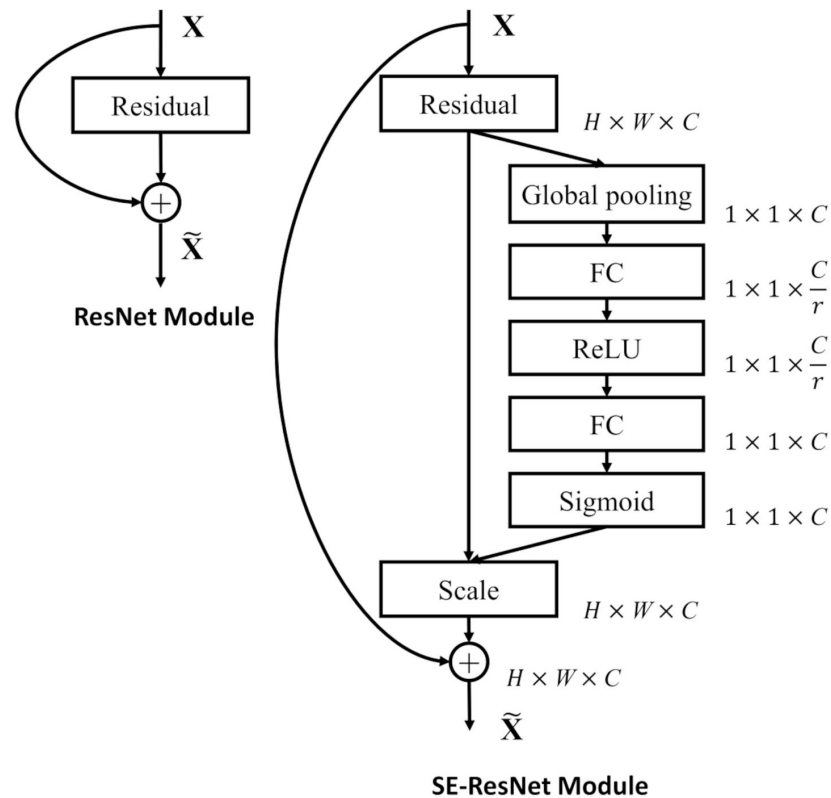
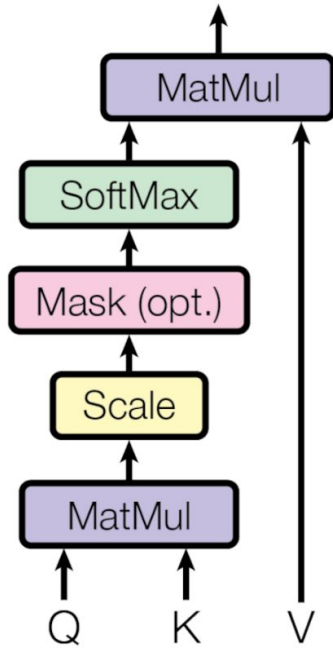
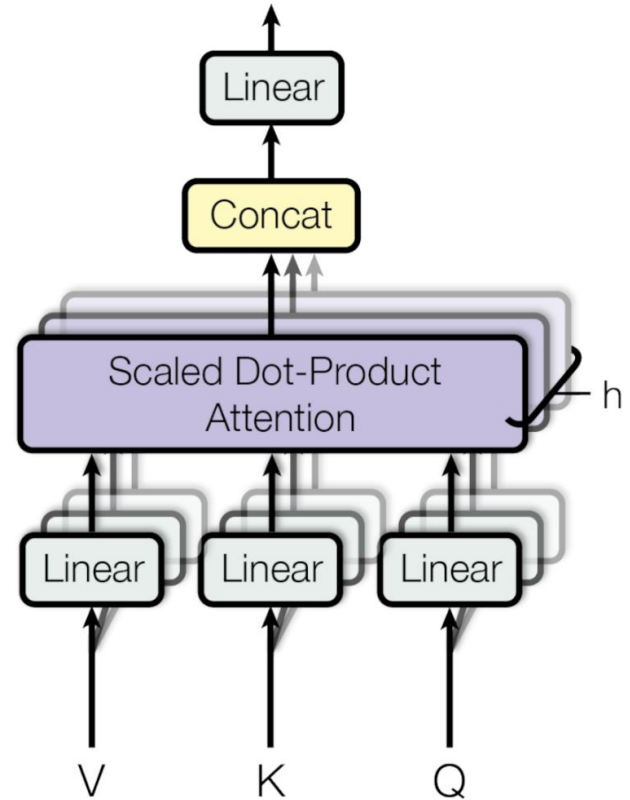


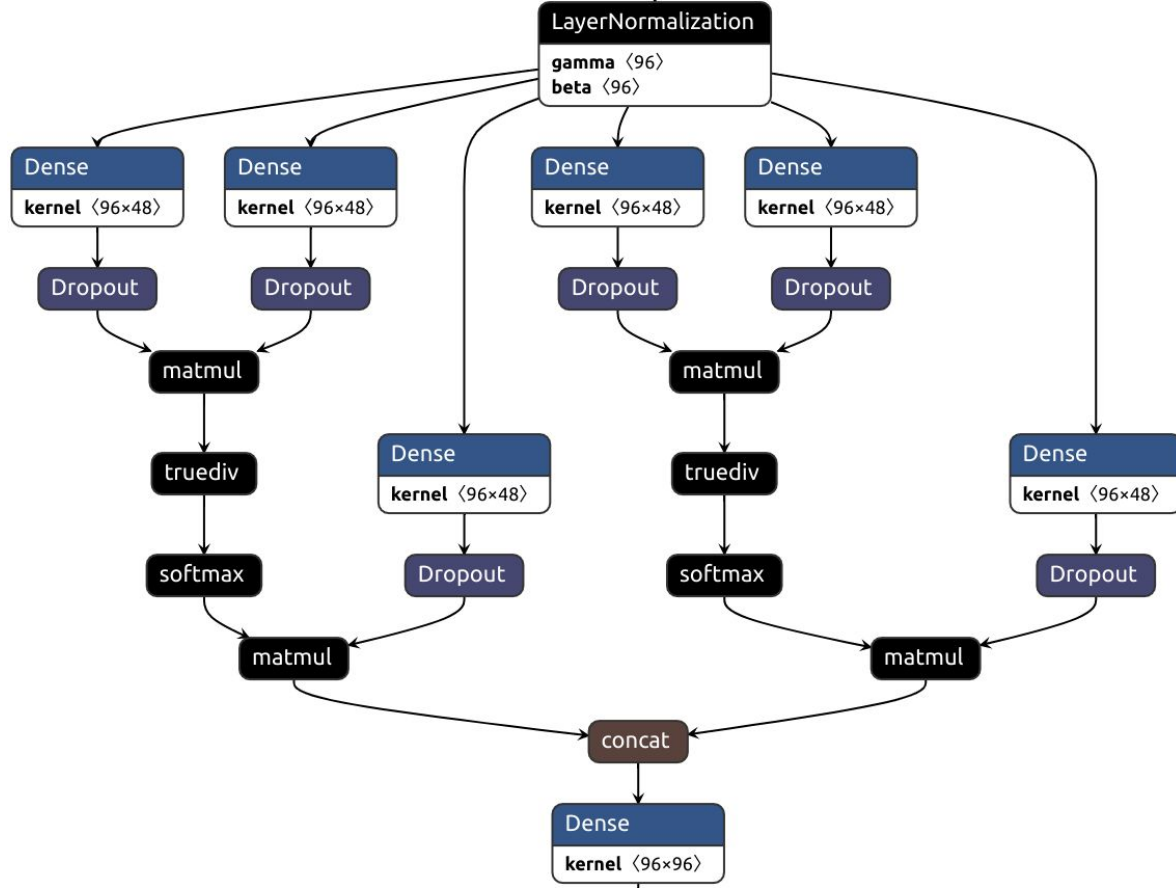
Fig. 3. The schema of the original Residual module (left) and the SE-ResNet module (right).

## Scaled Dot-Product Attention

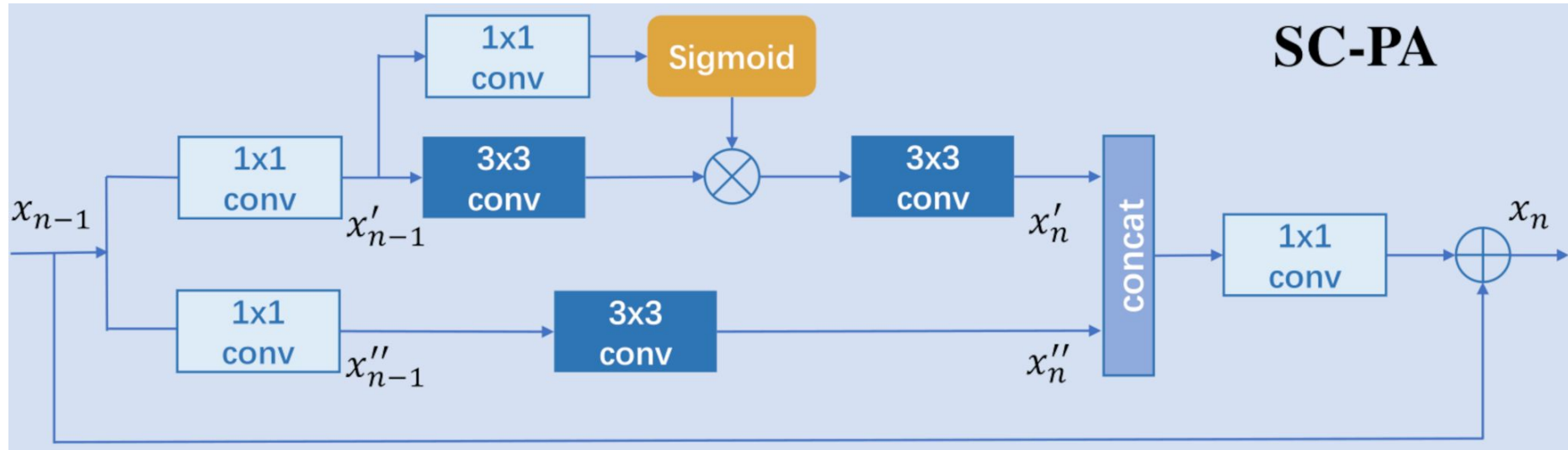


## Multi-Head Attention





Multi-Head Self-Attention (Vaswani et al. 2017; <https://arxiv.org/abs/1706.03762>)



---

# ANN: ...block patterns...

Multi-Layer Perceptron (MLP)

in general, any multilayer network

specifically, used as a recoding layer in metaformers

architecture:

- a layer-wise normalization layer

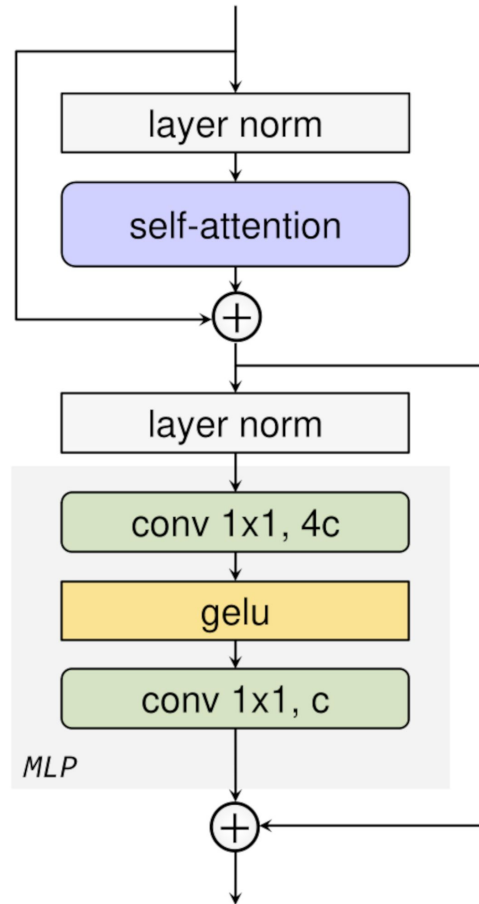
- 2 (or more) dense layers

  - same size or expand then contract channels

  - activate the first and/or the second

- an additive residual connection to the input

---



MLP in a transformer (Yang et al. 2023; <https://arxiv.org/abs/2210.01820>)

---

# ANN: ...block patterns...

metaformers (transformers and their modifications)

- initially designed as a better LSTM for NLP

  - used with tokenized text

  - used for computer vision with tokenized images

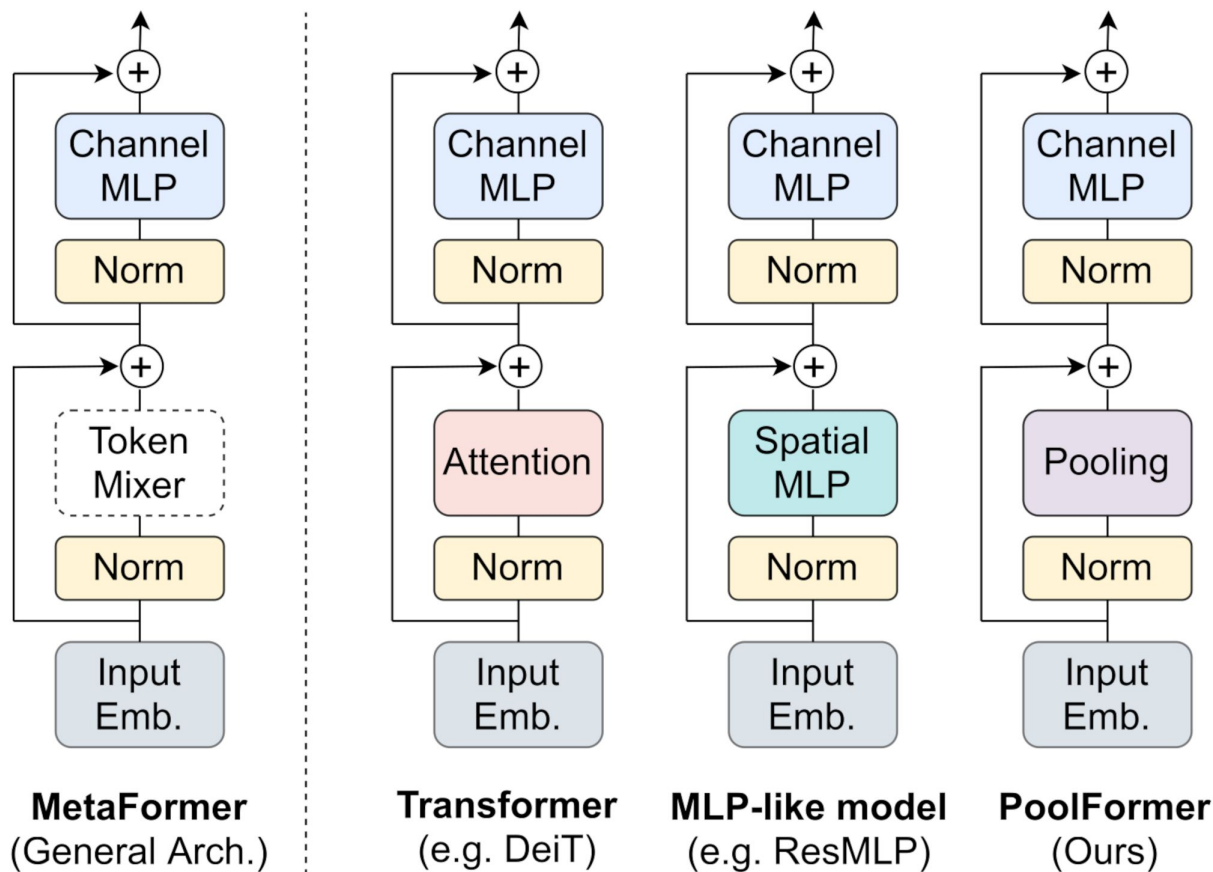
architecture:

- self-attention, smoothing, or mixing unit

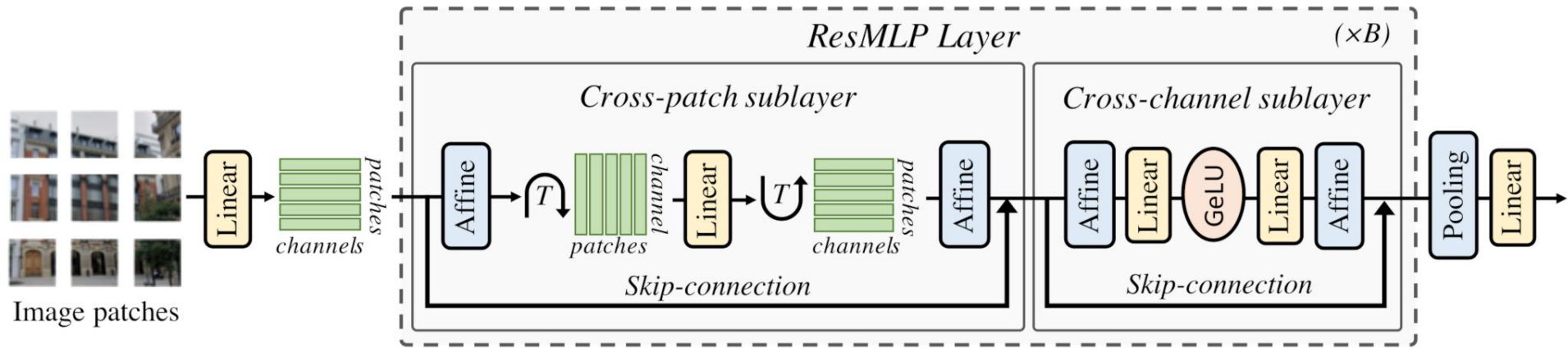
  - an additive residual connection to the input

  - an ML unit

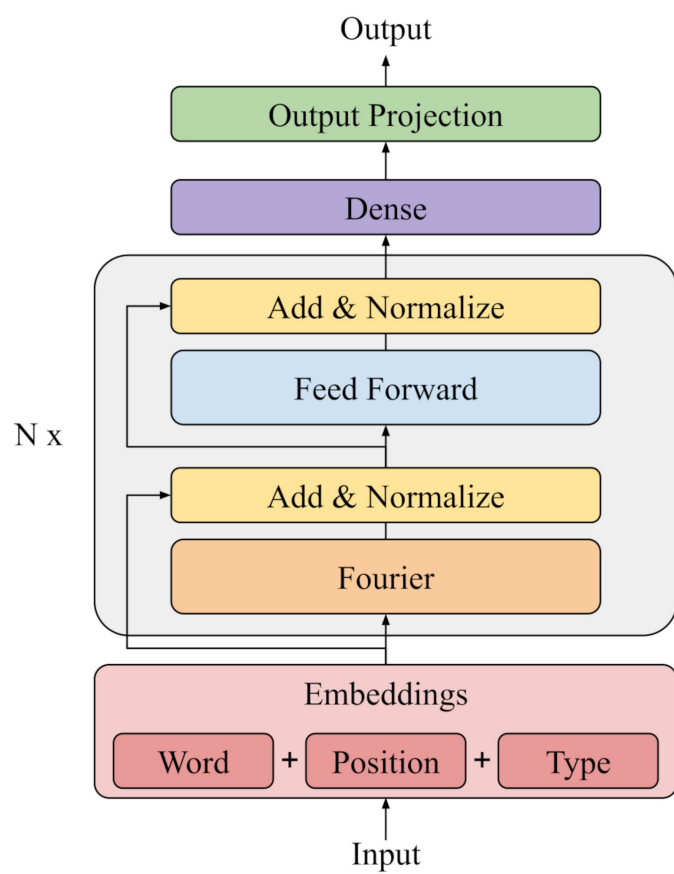
---

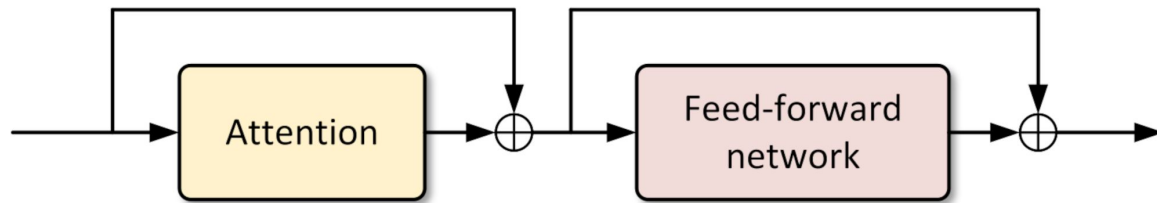




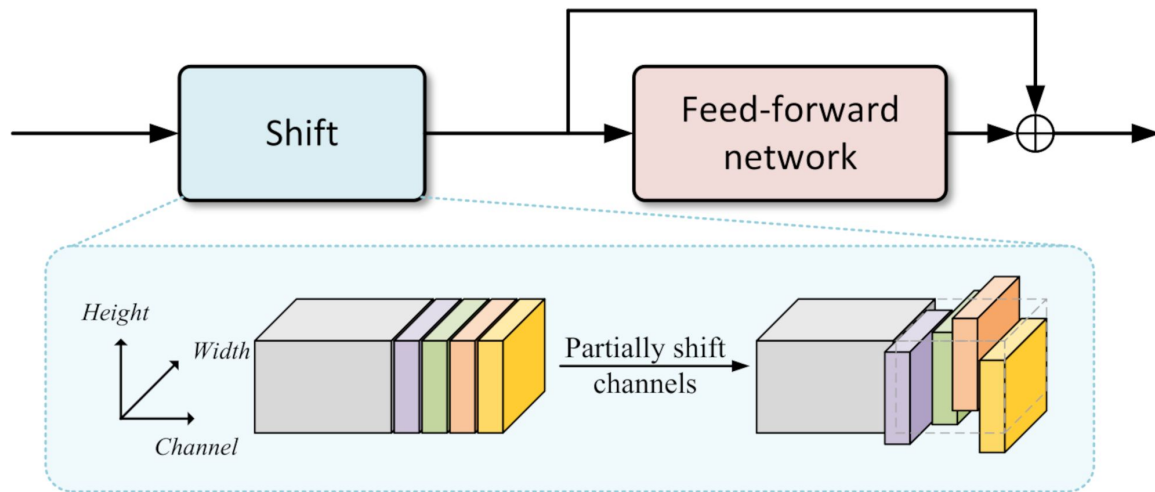


Touvron et al. 2021: Fig. 1 (ResMLP; Residual Multi-Layer Perceptrons)





(a) Standard attention building block.



(b) Our shift building block.

---

# ANN: model designs...

inversion (e.g. SqueezeNet, Inception, ResNet, EfficientNet)

used for classification or regression

$h \times w \times c \Rightarrow (1/n)h \times (1/n)w \times mc \mid w \times c \Rightarrow (1/n)w \times mc$

signal size: large  $\Rightarrow$  small

channels: few  $\Rightarrow$  many

many repeated blocks with residual connections

$h \times w$  reduction between blocks or in first block of a repeat

usually bottlenecks or inverted bottlenecks within each block

sometimes with attention mechanisms within or between blocks

---

---

# **ANN: ...model designs...**

constant size and shape (e.g. transformers)

used for classification, regression, or translation

many repeated blocks with residual connections

usually with attention, smoothing, or mixing

bottlenecks or inverted bottlenecks within each block

---

---

# ANN: ...model designs

bottleneck and unet (e.g. autoencoders)

- used for translation

- rarely with bottlenecks or inverted bottlenecks in blocks

- many repeated blocks

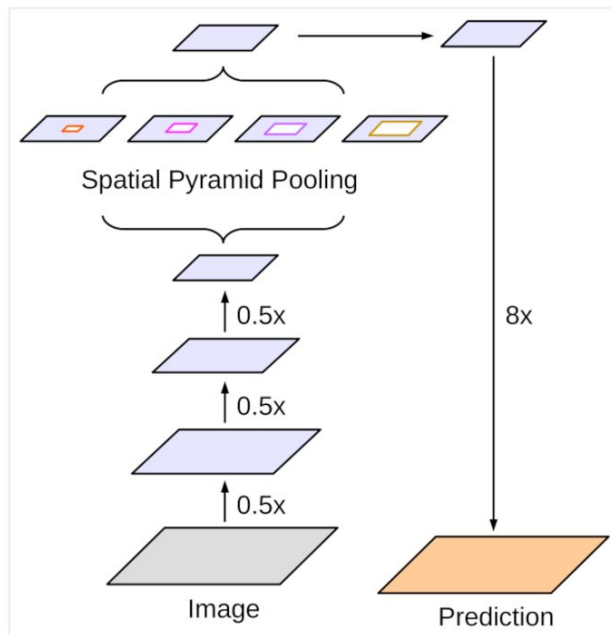
  - rarely with bottlenecks or inverted bottlenecks in blocks

    - the whole network is one giant bottleneck

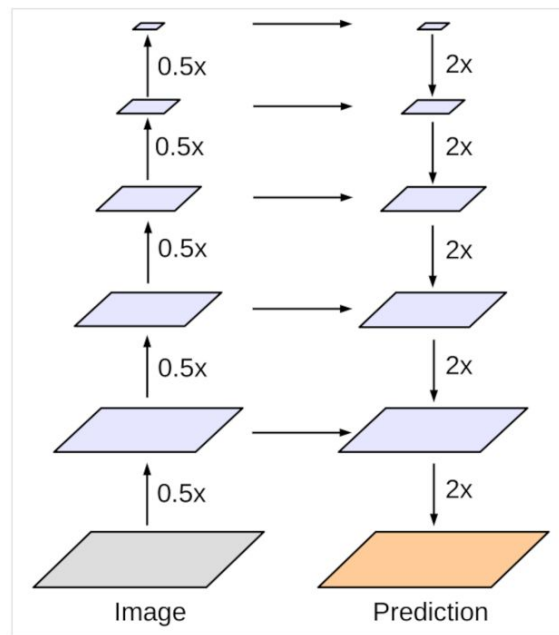
  - with no residual connections between blocks

    - or long residual connections between blocks (unet)

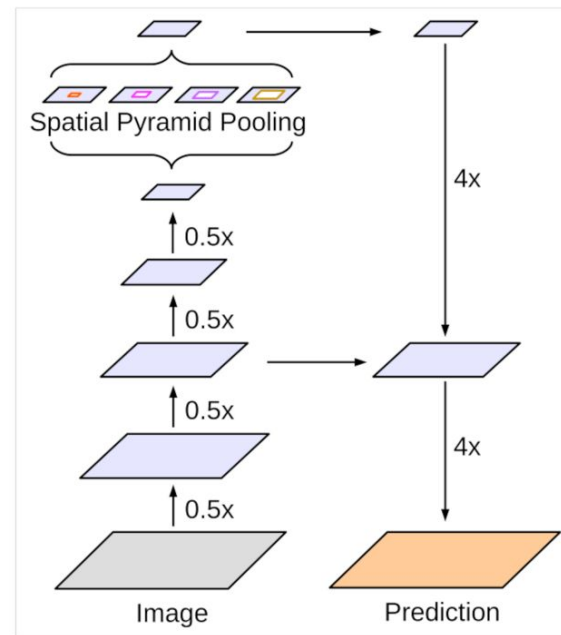
---



(a) Spatial Pyramid Pooling



(b) Encoder-Decoder



(c) Encoder-Decoder with Atrous Conv