**ECCV 2018**

**DeepLabV3+: Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation**

**CVPR 2021**

**SETR: Rethinking Semantic Segmentation from a Sequence-to-Sequence Perspective with Transformers**
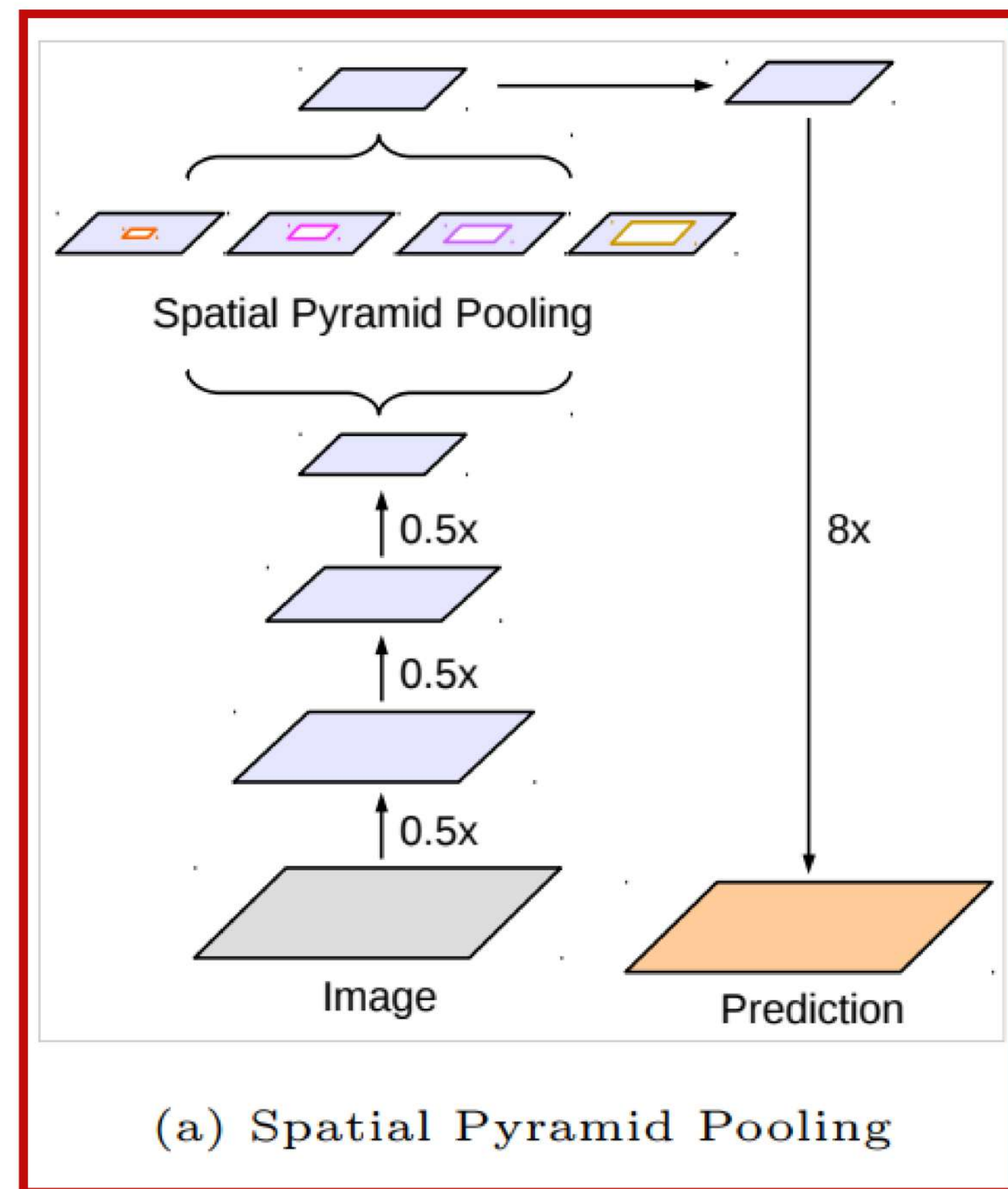
**NeurIPS 2021**

**SegFormer: Simple and Efficient Design for Semantic Segmentation with Transformers**
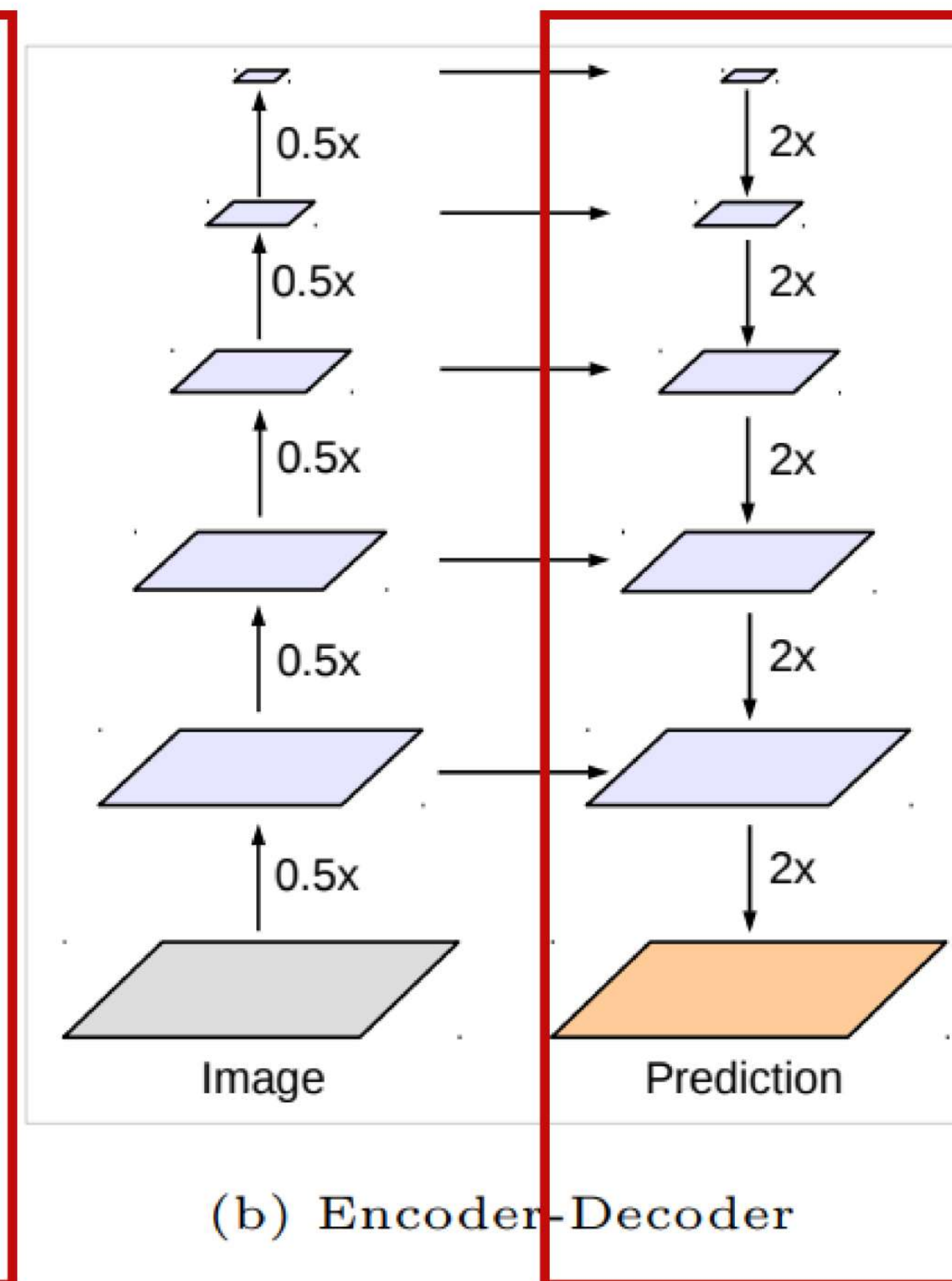
고민수

# DeepLabV3+

Encoder-Decoder with Atrous Separable Convolution
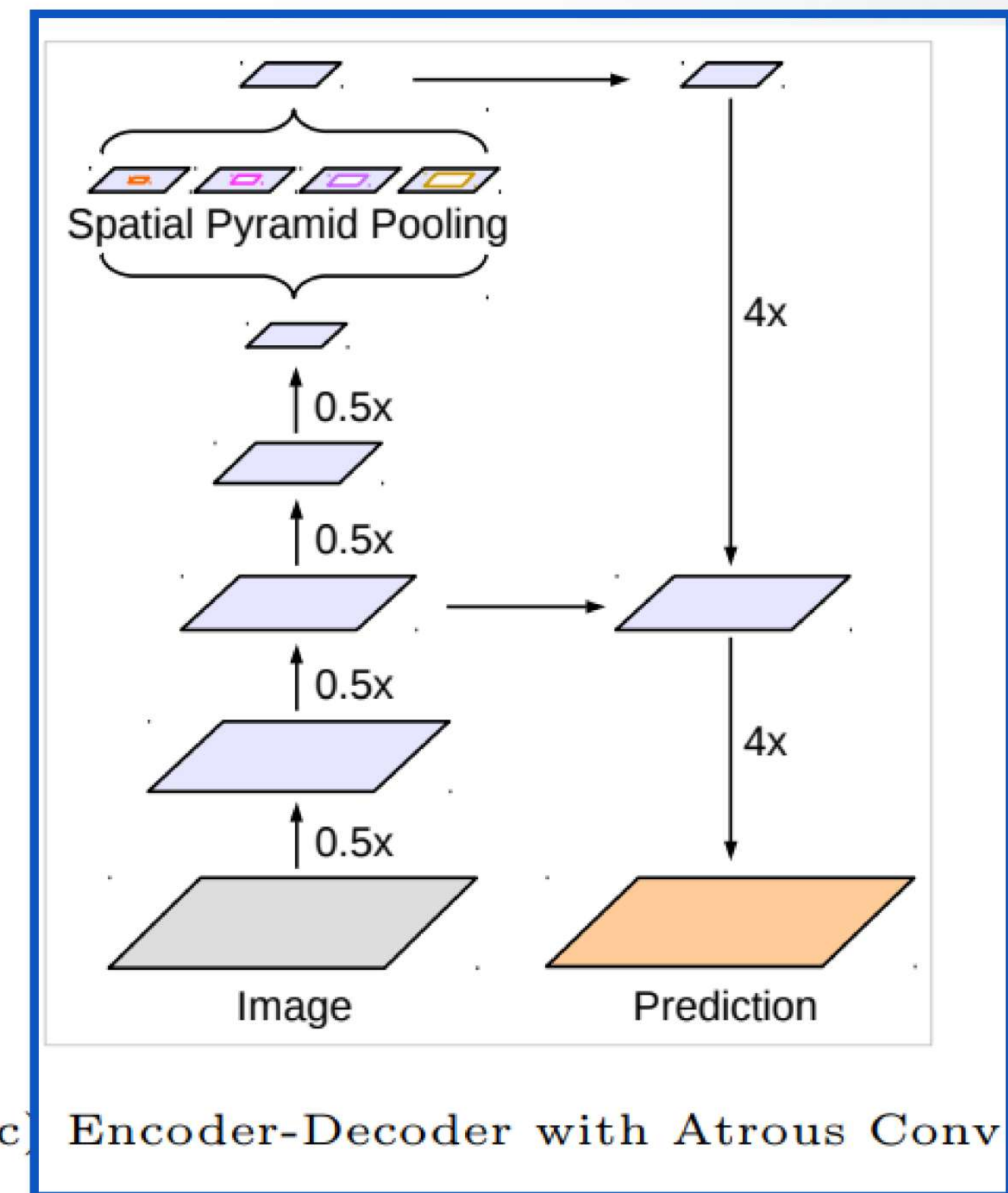for Semantic Image Segmentation

# Summary



(a) Spatial Pyramid Pooling

(b) Encoder-Decoder

(c) Encoder-Decoder with Atrous Conv

**DeepLabV3**

**Decoder module**

**DeepLabV3+**

exploiting
the multi-scale information

gradually recovers
the spatial information

# Encoder-Decoder with Atrous Convolution
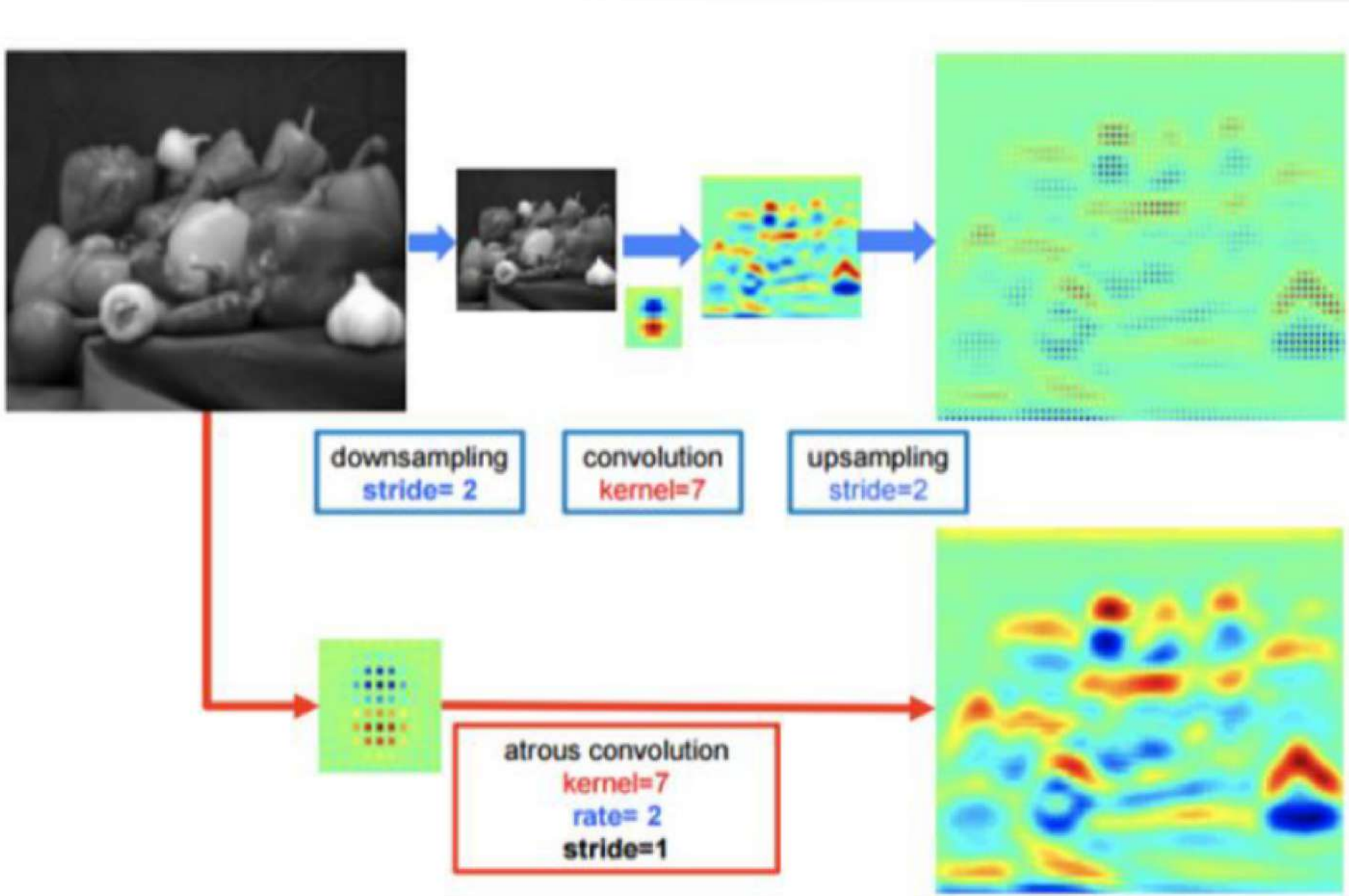
## Atrous convolution



(a) Going deeper without atrous convolution.

**output_stride**

(b) Going deeper with atrous convolution. Atrous convolution with $rate > 1$ is applied after block3 when $output\_stride = 16$.

**Smaller output feature map is more efficient in the segmentation task.**



Atrous Spatial Pyramid Pooling (ASPP)

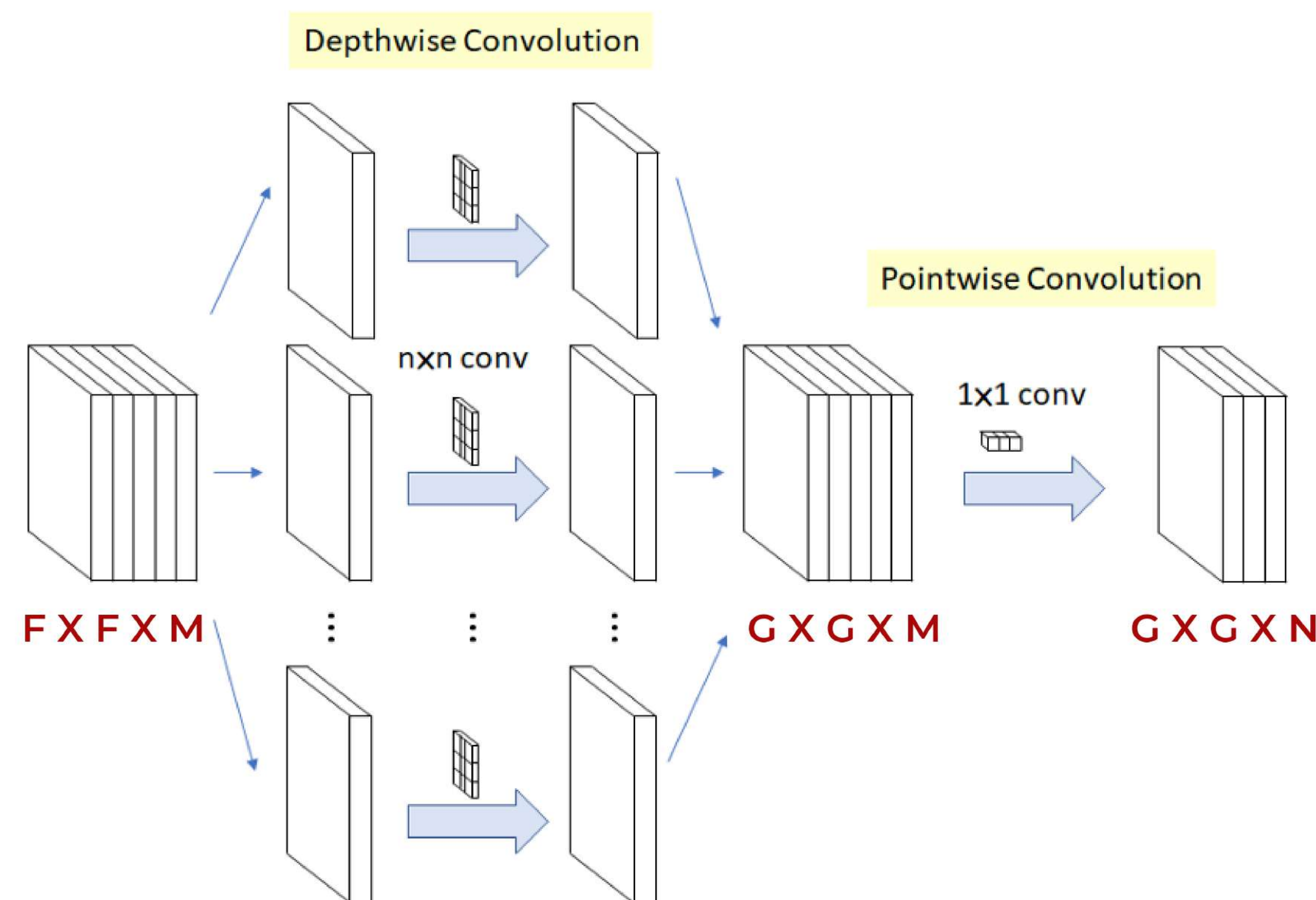**downsampling - convolution - upsampling**



**atrous convolution**

**Using atrous convolution shows better performance than downsampling - convolution - upsampling.**
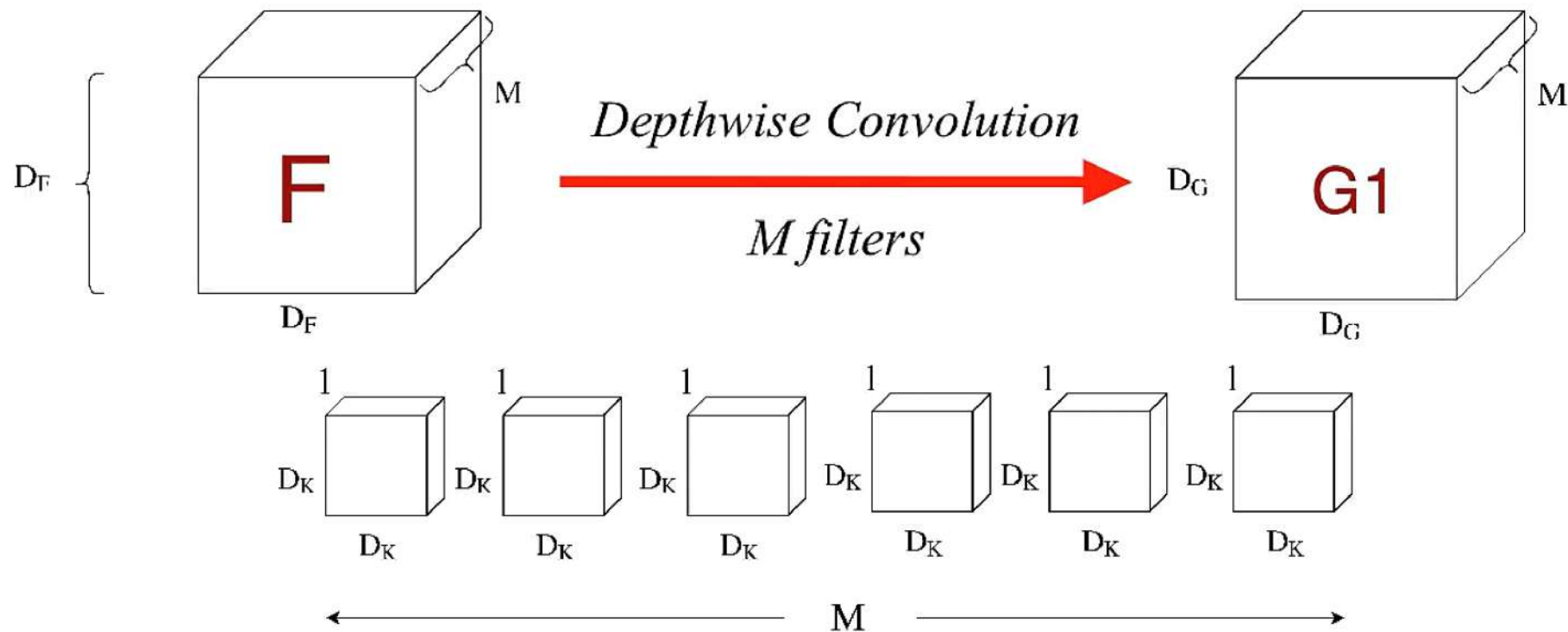
# Encoder-Decoder with Atrous Convolution
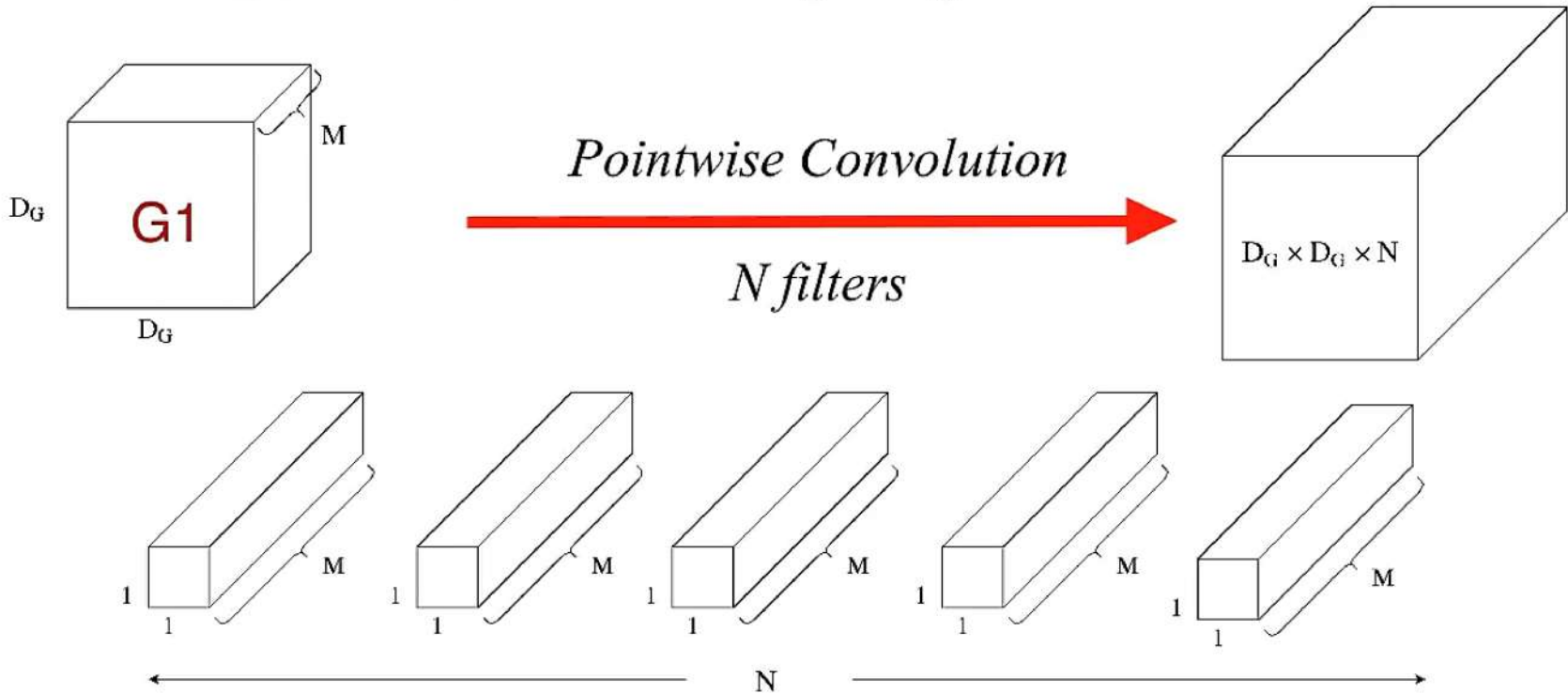
## Depthwise separable convolution

# Encoder-Decoder with Atrous Convolution

## Standard vs Depthwise

### Convolution



Mults once $= D_K^2 \times M$

Mults per Kernel $= D_G^2 \times D_K^2 \times M$

Mults N Kernels $= N \times D_G^2 \times D_K^2 \times M$

**N * G² * K² * M**

### Depthwise Separable Convolution

#### 1. Depthwise Convolution: Filtering Stage



**G² * K² * M**

### Depthwise Separable Convolution

#### 2. Pointwise Convolution: Filtering Stage



**N * G² * M**

$$\frac{No.\, Mults\ in\ Depthwise\ Separable\ Conv}{No.\, Mults\ in\ Standard\ Conv} = \frac{M \times D_G^2\, (D_K^2 + N)}{N \times D_G \times D_G \times D_K \times D_K \times M}$$
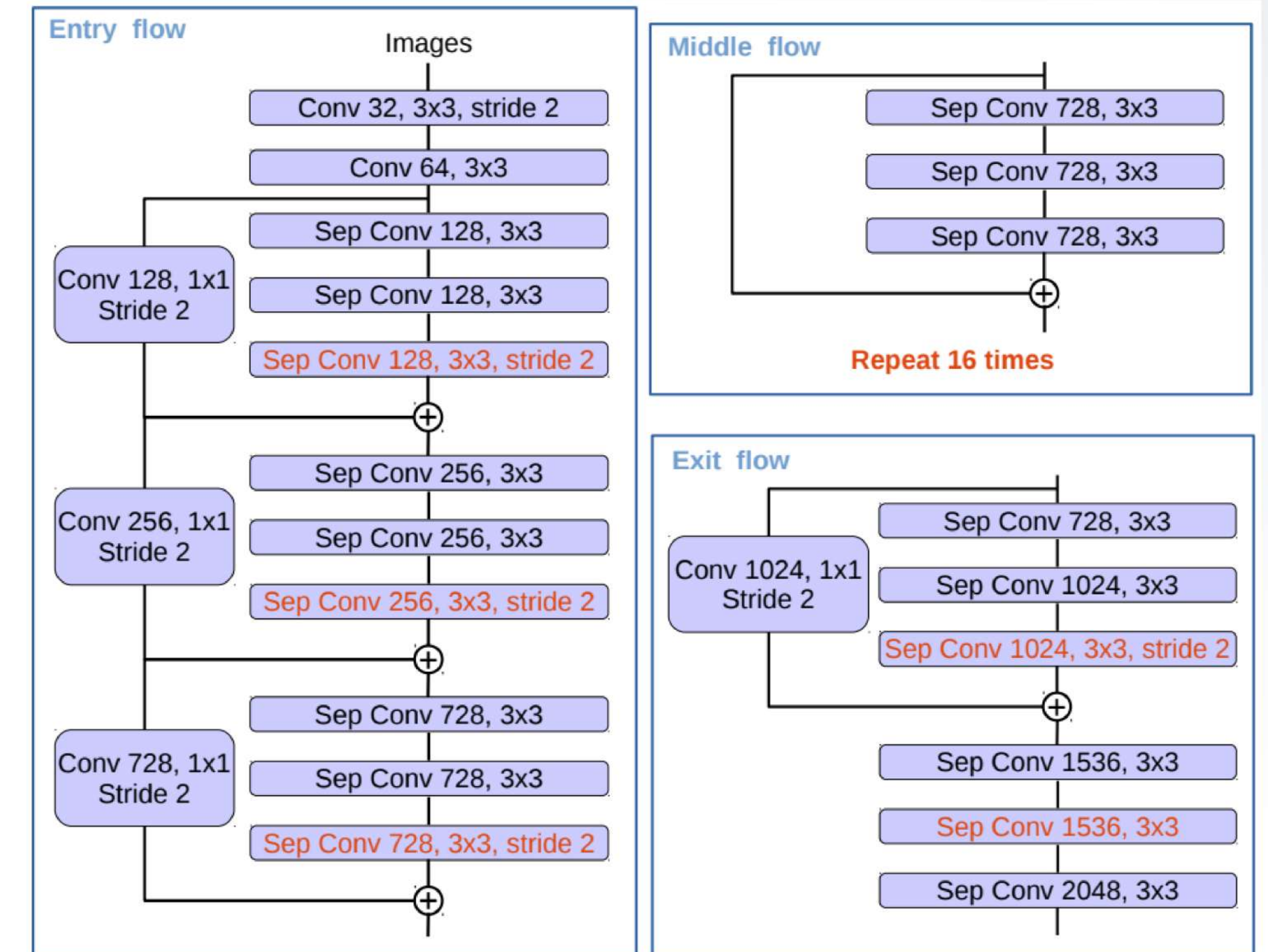
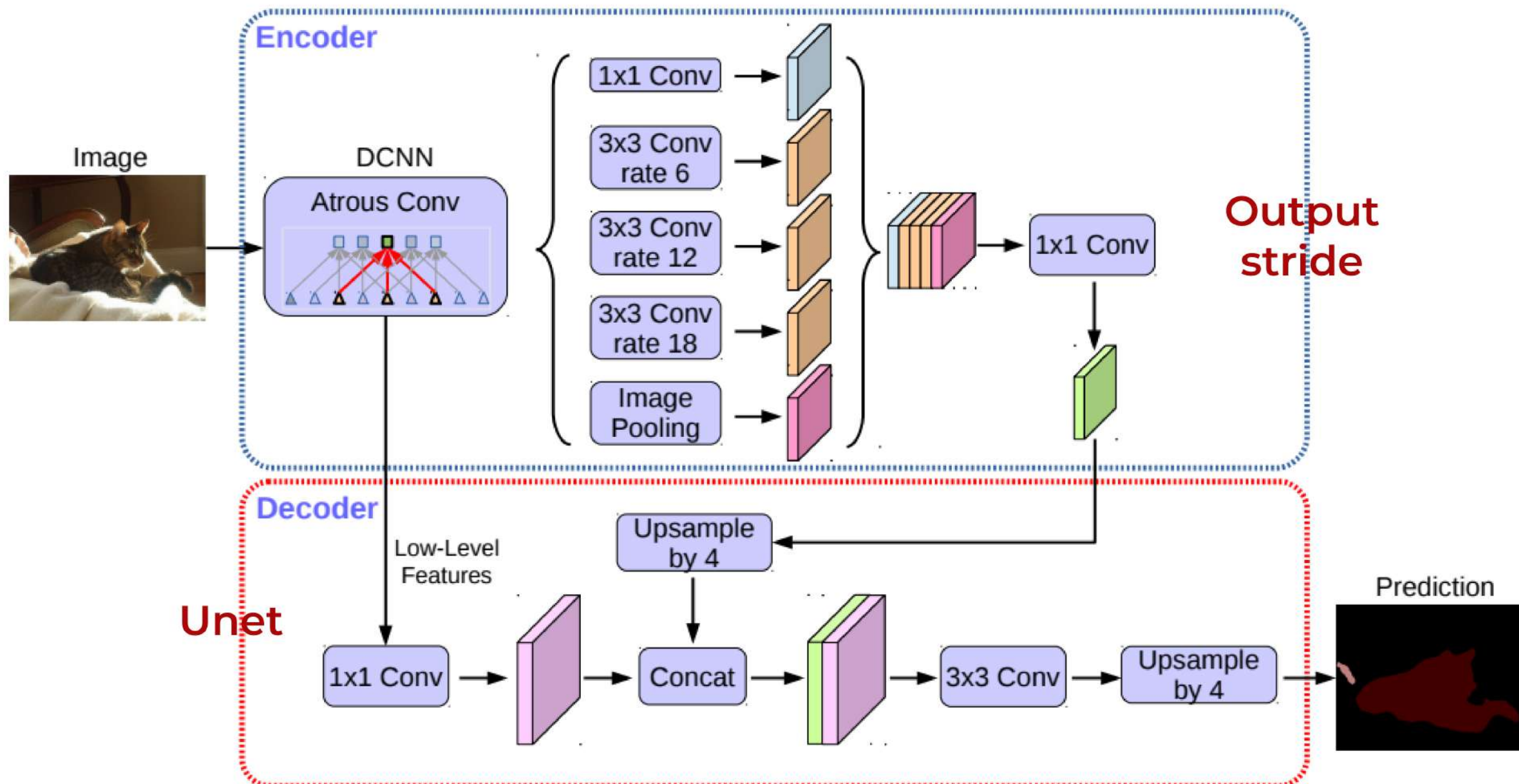$$\frac{No.\, Mults\ in\ Depthwise\ Separable\ Conv}{No.\, Mults\ in\ Standard\ Conv} = \frac{D_K^2 + N}{(D_K^2 \times N)} = \frac{1}{N} + \frac{1}{D_K^2}$$

**If N=1024, K=3,**

$$\frac{No.\, Mults\ in\ Depthwise\ Separable\ Conv}{No.\, Mults\ in\ Standard\ Conv} = \frac{1}{1024} + \frac{1}{3^2} = 0.112$$
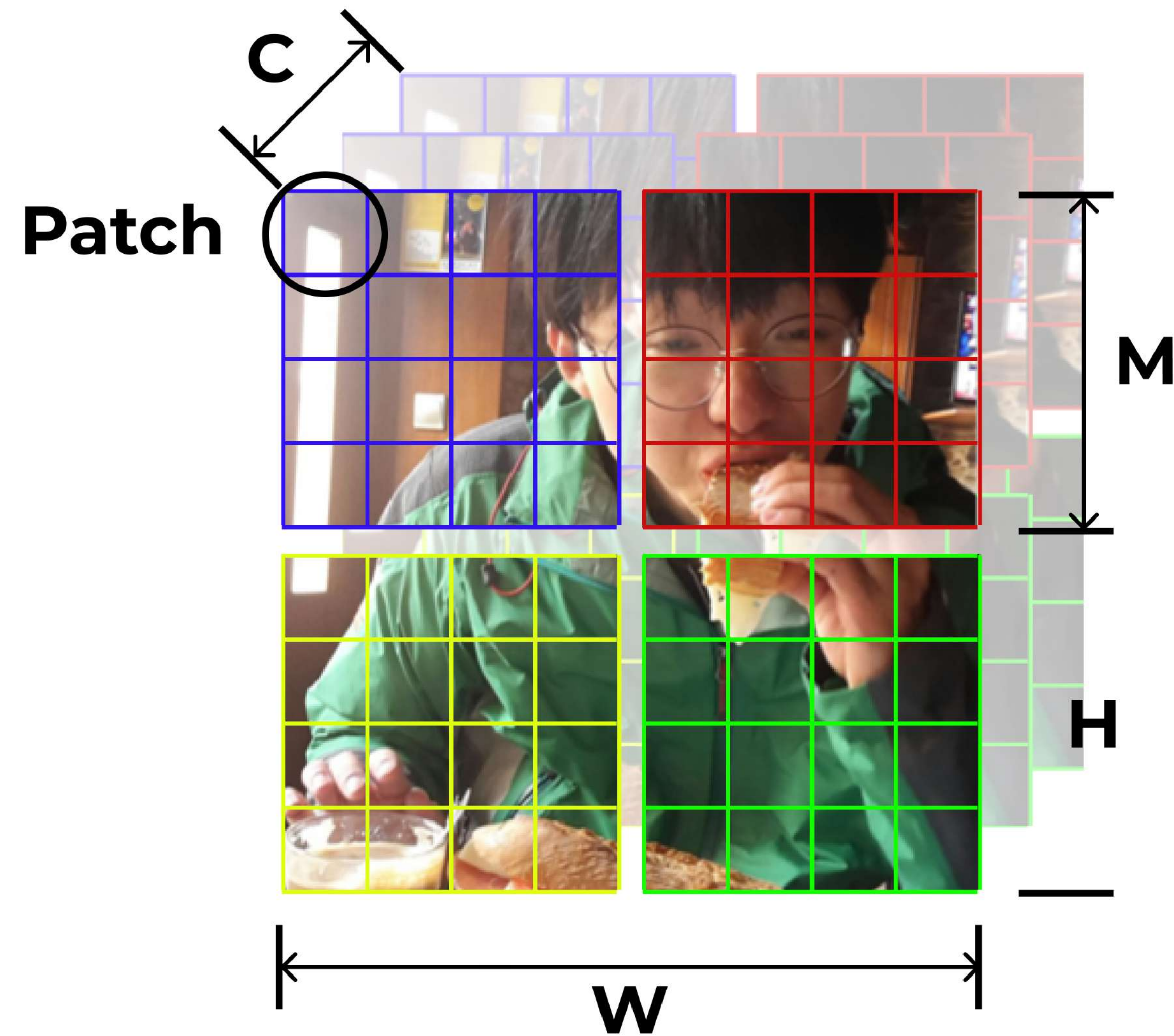
# DeepLabV3+

## Structure



**Backbone - modified Xception**

1. More layer
2. All max pooling operations are replaced by Depthwise separable convolutions
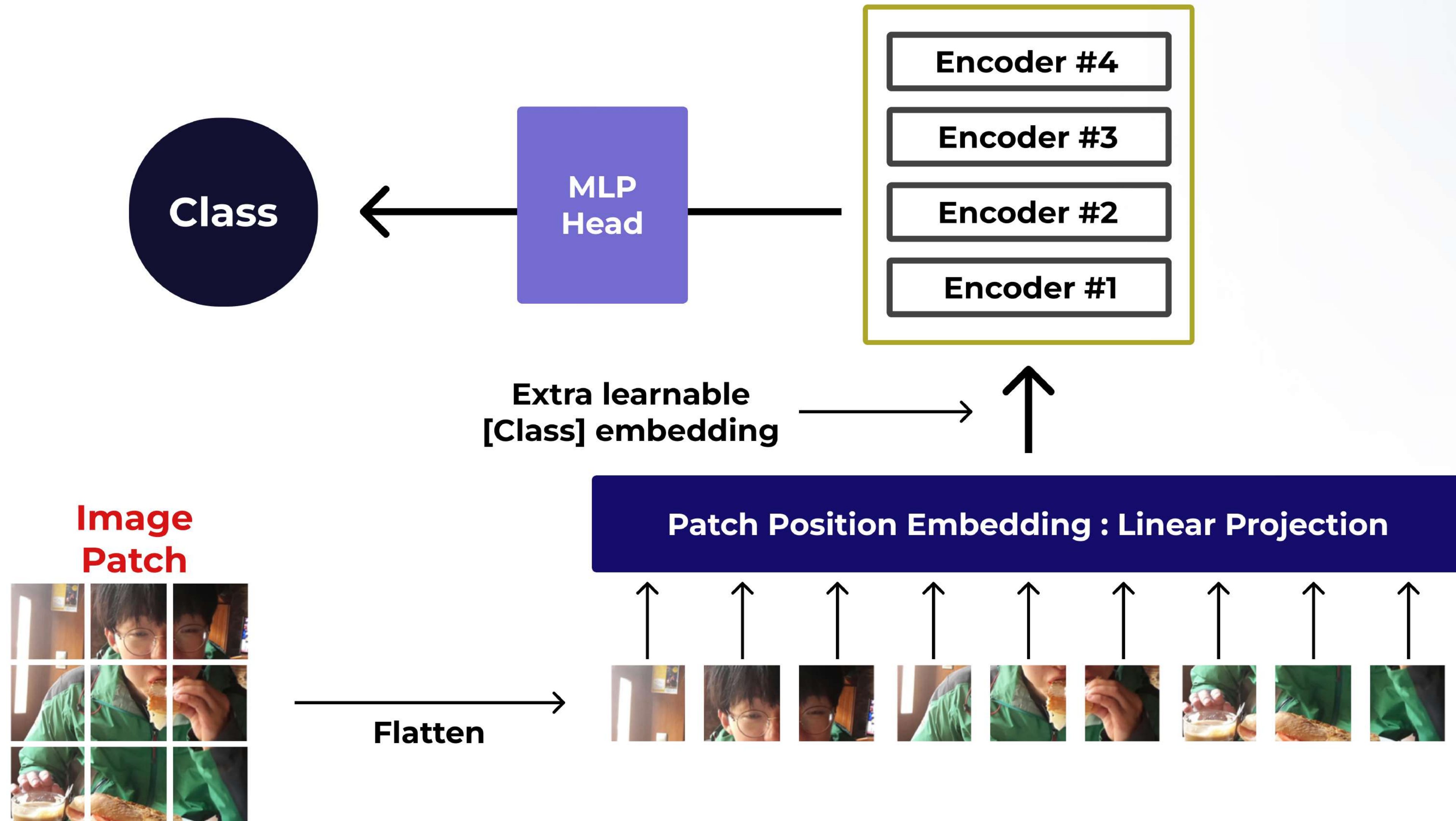3. Batch normalization and ReLU ard added

# Vision Transformer

AN IMAGE IS WORTH 16X16 WORDS:
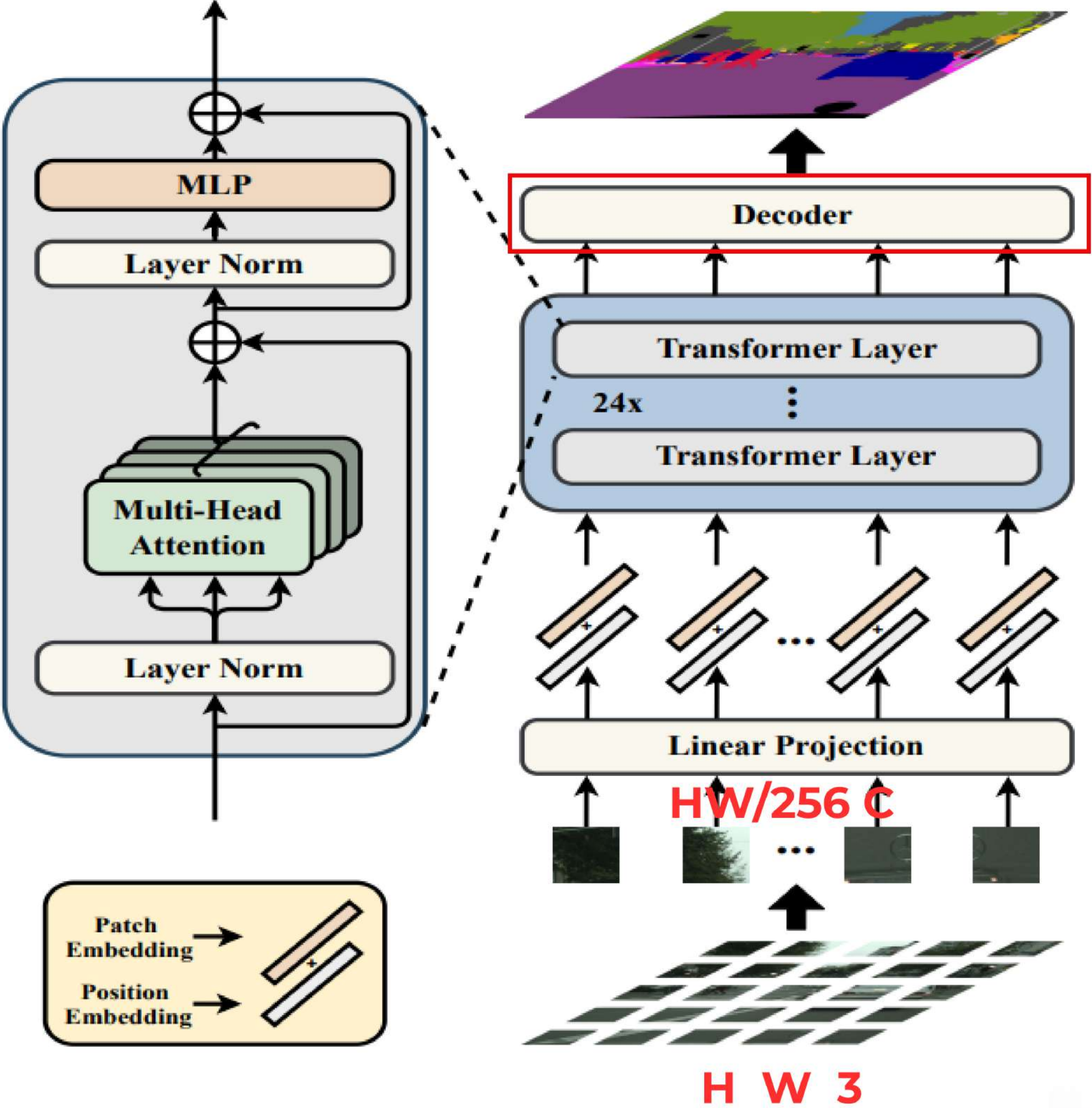TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

# Notice

Patch

C

M

H

W

# SETR Summary

## Structure



**HW/256 C**

**H W 3**

## Decoder

1. Naive - simply bilinearly upsampled

2. PUP - Progressive UPsampling



reshape    conv→2x    conv→2x    conv→2x    conv→2x

$\frac{HW}{256} \times 1024$    $\frac{H}{16} \times \frac{W}{16} \times 1024$    $\frac{H}{8} \times \frac{W}{8} \times 256$    $\frac{H}{4} \times \frac{W}{4} \times 256$    $\frac{H}{2} \times \frac{W}{2} \times 256$    $H \times W \times 19$

3. MLA - Multi-Level feature Aggregation



$z^{24}$

$z^{18}$

$z^{12}$

$z^6$

reshape-conv    conv-conv-4x    conv-4x

**bilinear operation**

# SegFormer

**Simple and Efficient Design
for Semantic Segmentation with Transformers**

# SegFormer - Summary

## Structure
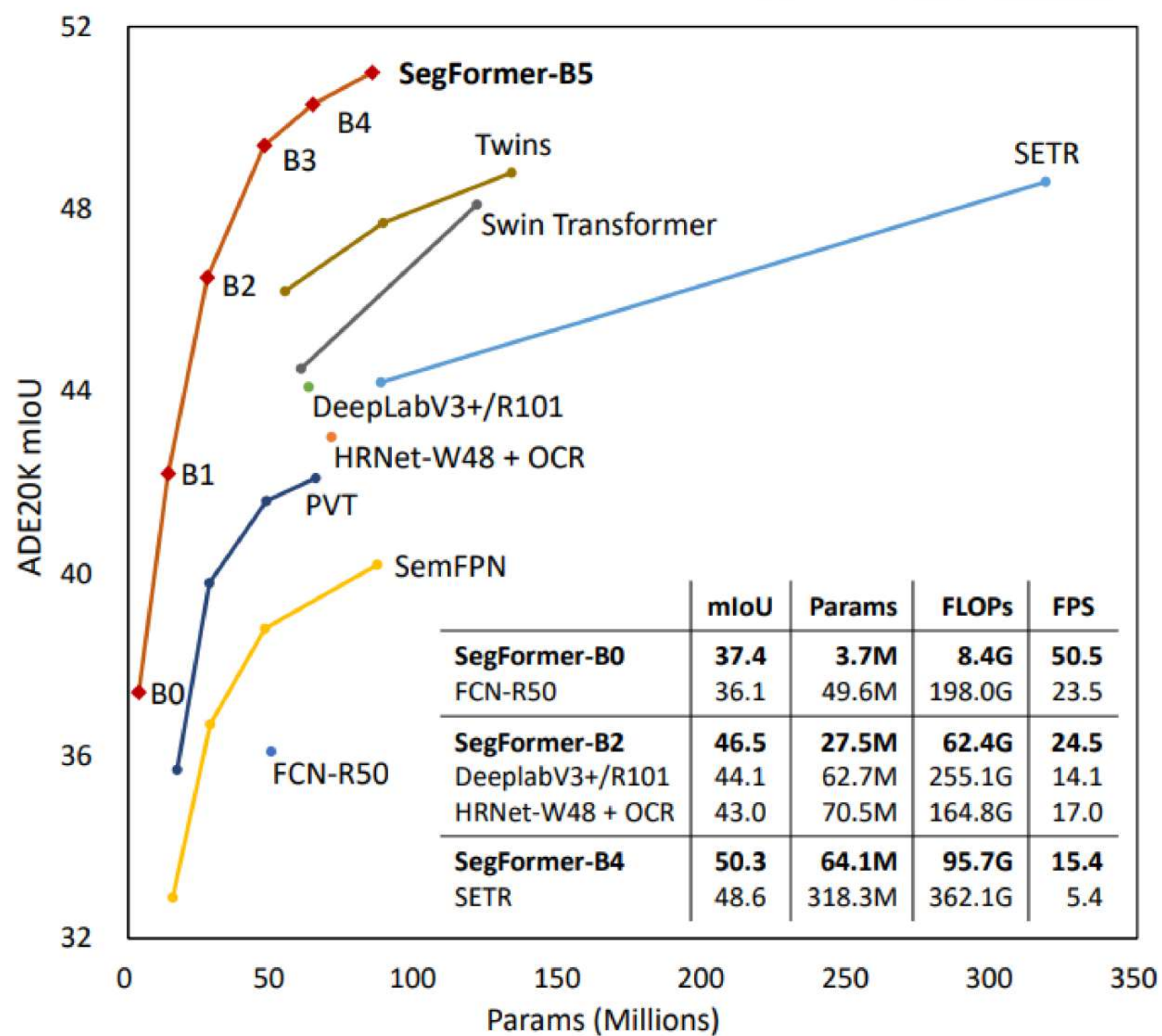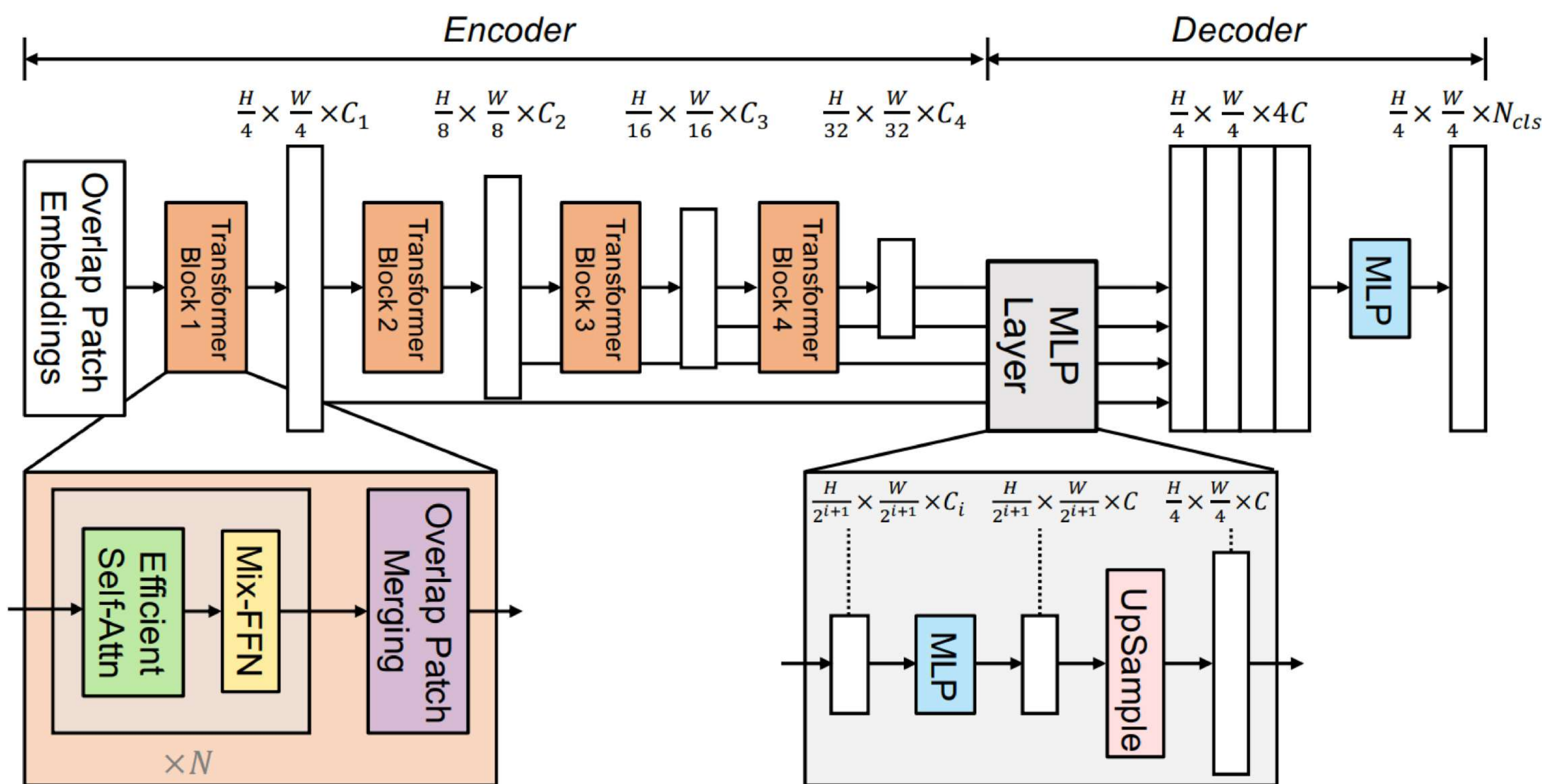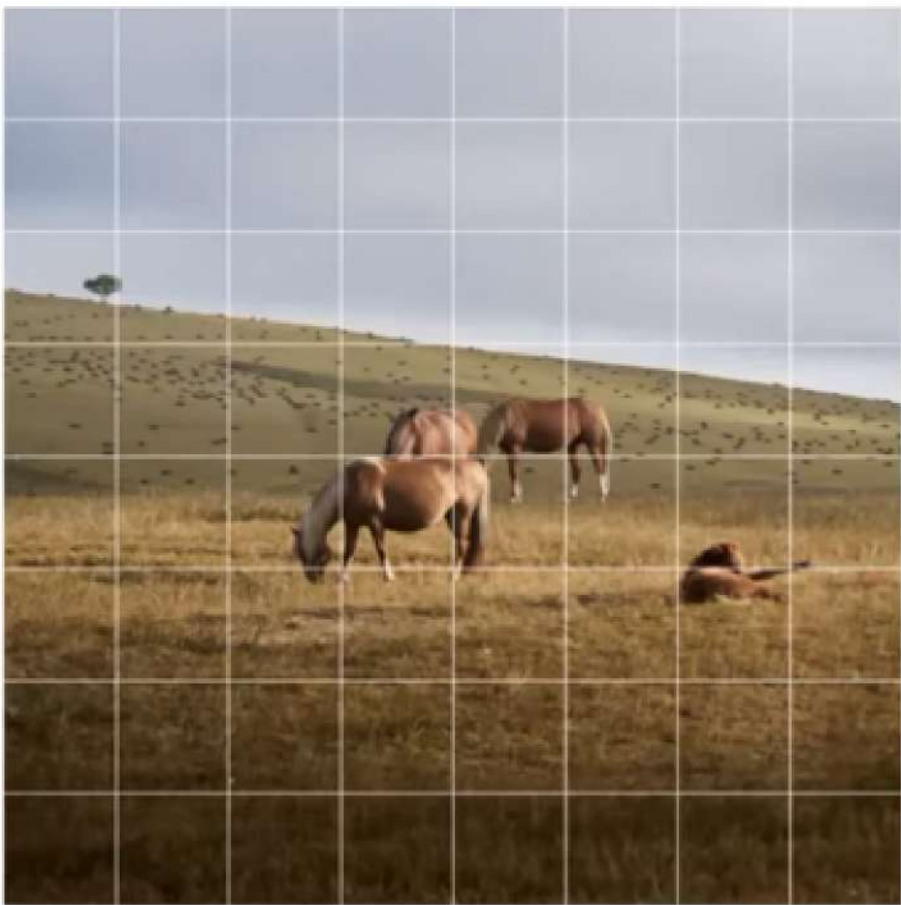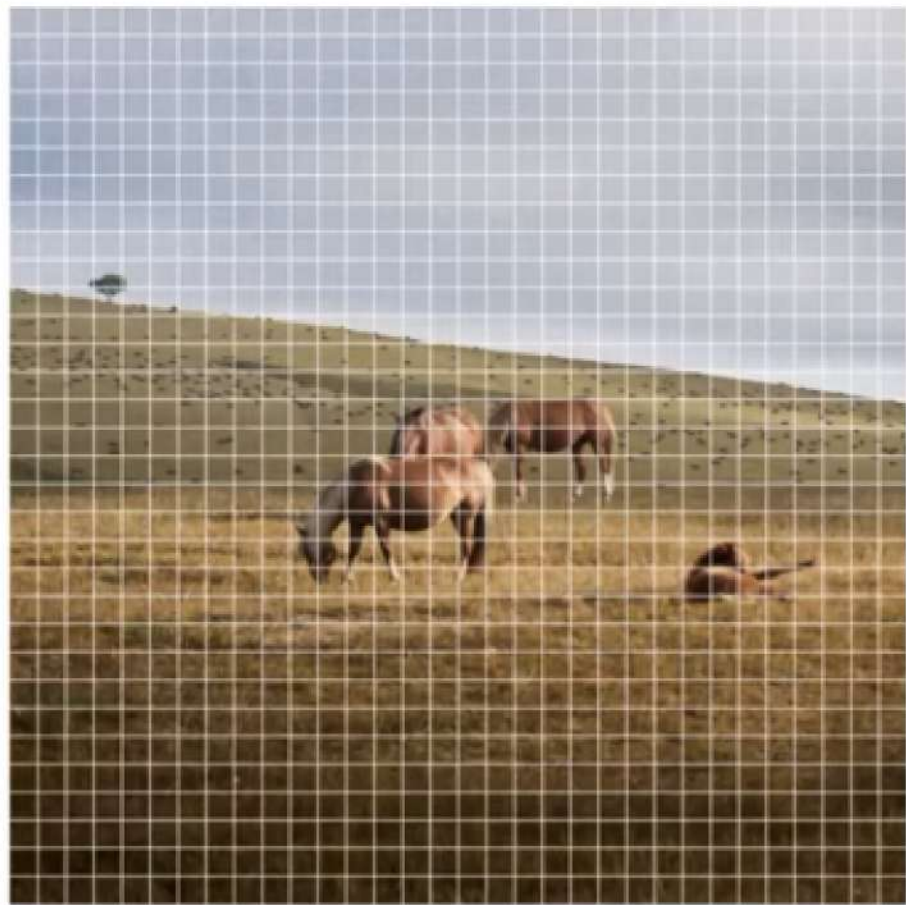
# Hierarchical Feature Representation

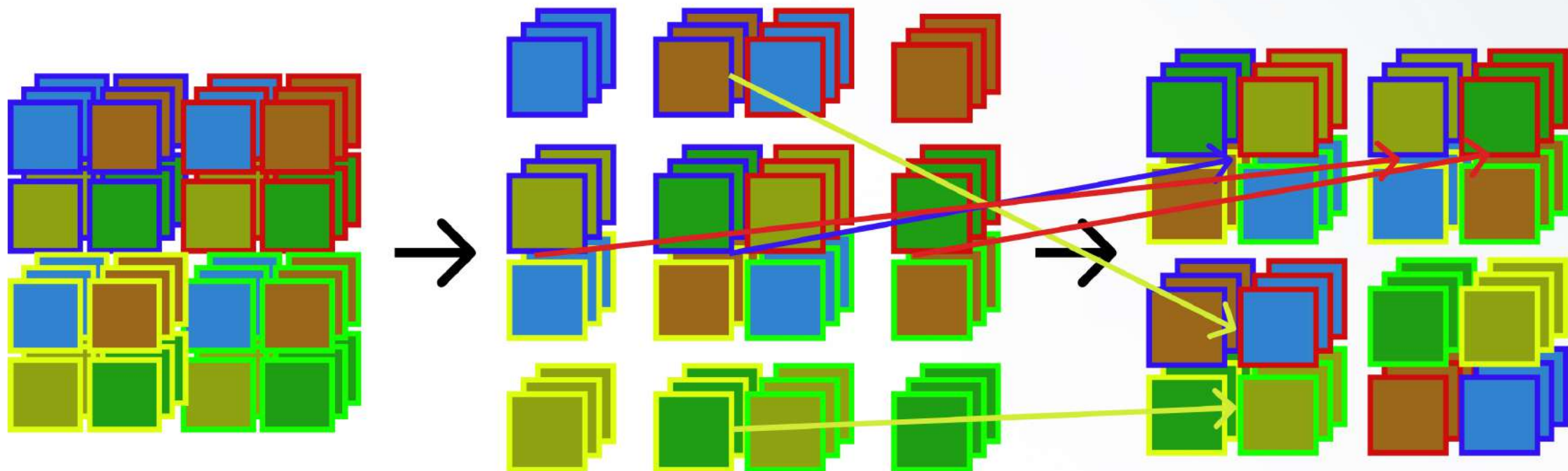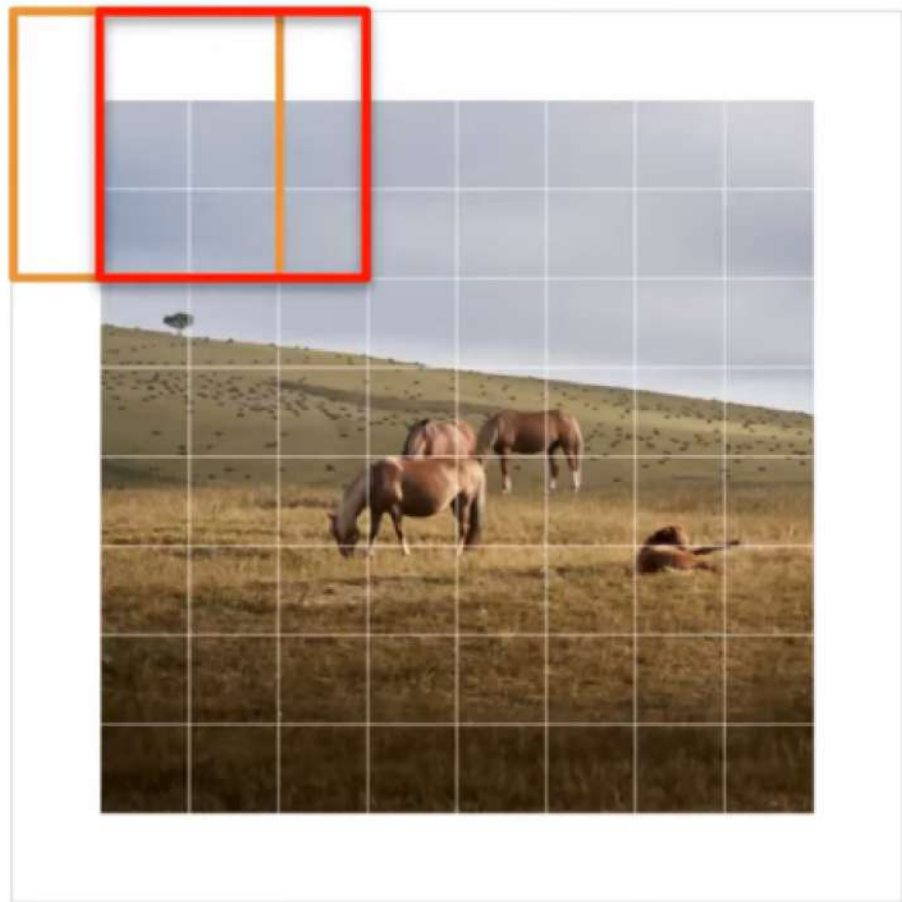## Hierarchical Feature Representation



ViT
(P=16)

Segformer
(P=4)

**Self Attention cost ?**
**local continuity ?**

Shifted Window ( Swin TR )



## Overlapping Patch Window ( SegFormer )



**Similar to how CNNs work**

Stage 1         (K=7, S=4, P=3)

Stage 2, 3, 4     (K=3, S=2, P=1)

# Hierarchical Feature Representation

## Efficient Self-Attention

$$\text{Attention}(Q, K, V) = \text{Softmax}(\frac{QK^\mathsf{T}}{\sqrt{d_{head}}})V.$$

**K : (N, C)**

**K : (N/R, C)**

**Sequence Reduction Process**

$$\hat{K} = \text{Reshape}(\frac{N}{R}, C \cdot R)(K)$$

$$K = \text{Linear}(C \cdot R, C)(\hat{K}),$$



## Mix-FFN

**positional encoding** $\longrightarrow$ **Fixed input resolution**

**The size must be matched through interpolation which causes performance degradation**

$$\mathbf{x}_{out} = \text{MLP}(\text{GELU}(\text{Conv}_{3\times3}(\text{MLP}(\mathbf{x}_{in})))) + \mathbf{x}_{in},$$

Patches

$$\text{Conv}_{3\times3}(\text{MLP}(\mathbf{x}_{in}))$$

- **ConV 3x3 layers use depth-wise convolution**
- **provide location information**

# Lightweight All-MLP Decoder

$$\hat{F}_i = \text{Linear}(C_i, C)(F_i), \forall i$$

$$\hat{F}_i = \text{Upsample}(\frac{W}{4} \times \frac{W}{4})(\hat{F}_i), \forall i$$

$$F = \text{Linear}(4C, C)(\text{Concat}(\hat{F}_i)), \forall i$$

$$M = \text{Linear}(C, N_{cls})(F),$$

1. All channels of multi-level features are integrated equally.

2. Integrate the feature size to 1/4 the size of the original image.

3. Concatenate the features and in this process restore the channel that was multiplied by a factor of 4.

4. Predict the final segmentation mask.
   (shape: B(batch) x N(num of classes) x H/4 x W/4)

1. Manual work and computational effort are not greatly required

2. It can have a larger effective field compared to CNN.

# Lightweight All-MLP Decoder

$$\hat{F}_i = \text{Linear}(C_i, C)(F_i), \forall i$$

$$\hat{F}_i = \text{Upsample}(\frac{W}{4} \times \frac{W}{4})(\hat{F}_i), \forall i$$

$$F = \text{Linear}(4C, C)(\text{Concat}(\hat{F}_i)), \forall i$$

$$M = \text{Linear}(C, N_{cls})(F),$$

1. All channels of multi-level features are integrated equally.

2. Integrate the feature size to 1/4 the size of the original image.

3. Concatenate the features and in this process restore the channel that was multiplied by a factor of 4.

4. Predict the final segmentation mask.
   (shape: B(batch) x N(num of classes) x H/4 x W/4)

1. Manual work and computational effort are not greatly required

2. It can have a larger effective field compared to CNN.