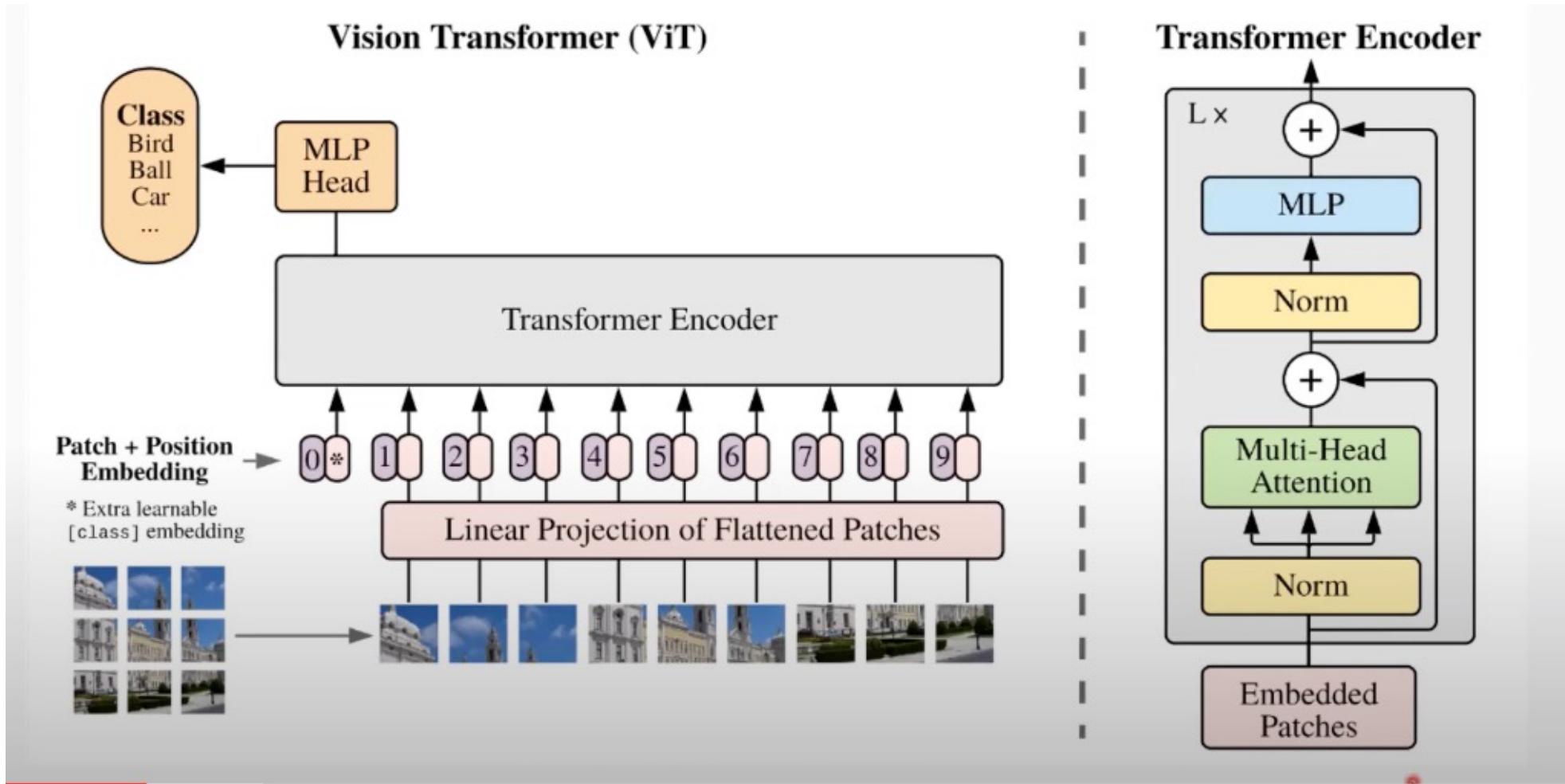


DEEP LEARNING FOR COMPUTER VISION

Week11

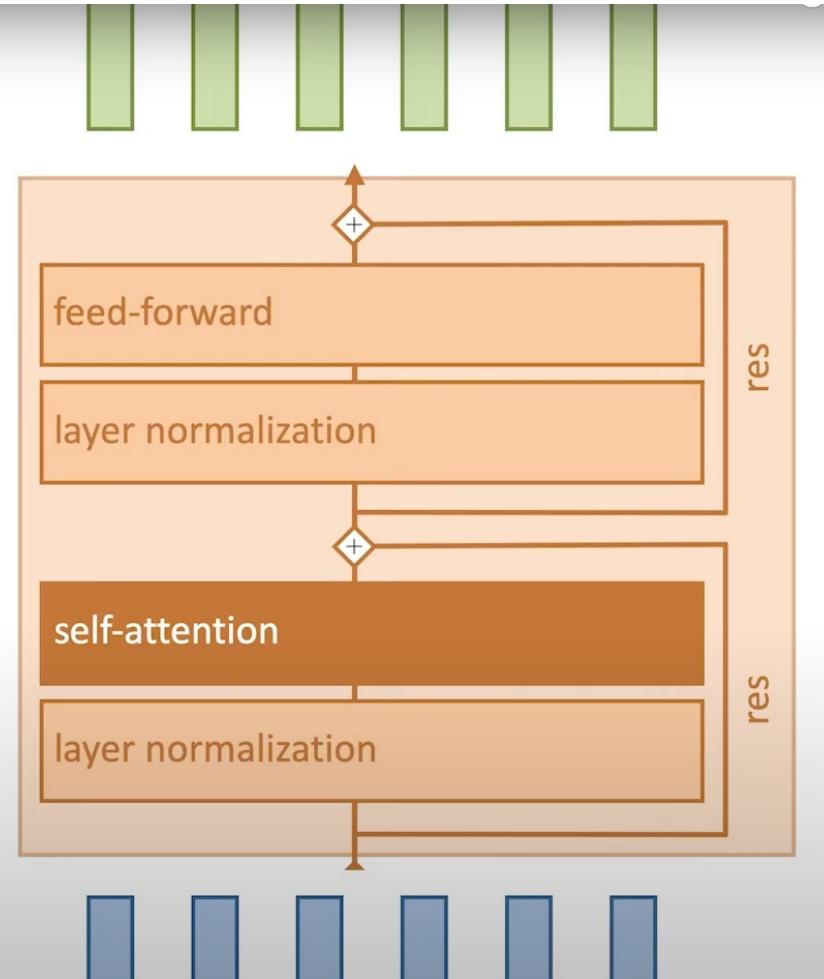


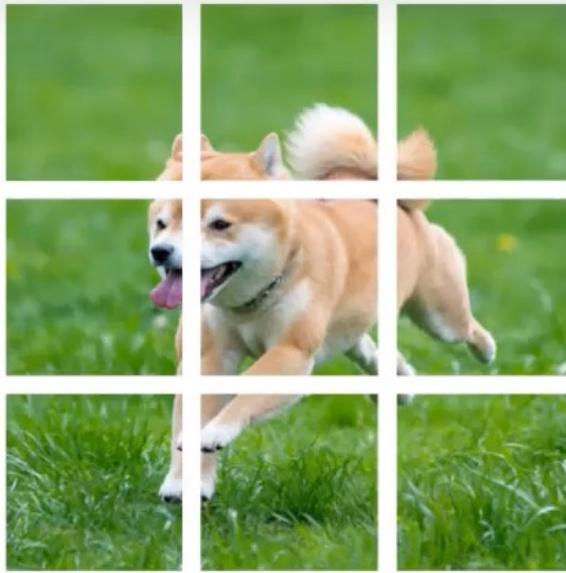
Dr. Tuchsanai. PloySuwan



TRANSFORMER BLOCK

```
class Block(nn.Module):  
  
    def forward(self, x):  
  
        y = self.layernorm(y)  
  
        y = self.attention(x)  
  
        x = x + y  
  
        y = self.layernorm(x)  
  
        y = self.linear(x)  
  
    return x + y
```





Positional encodings

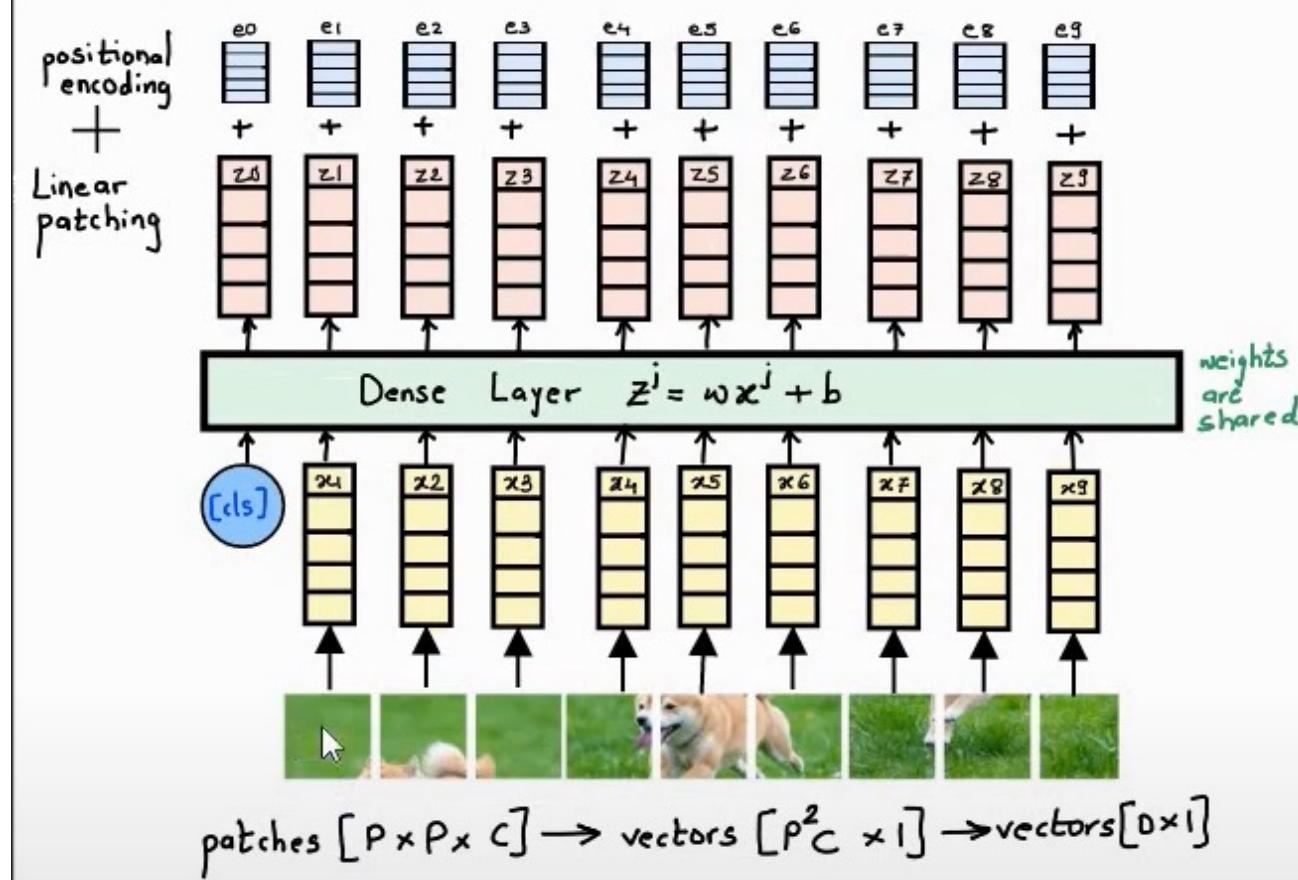
Linear patching

Share weights

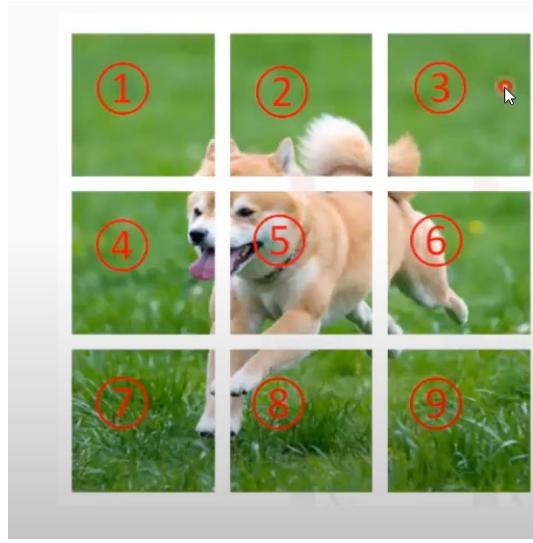
- The linear representation of the patches is learned by the network. The patches do not overlap in the paper but you can overlap them if you want to.
- There is no activation function applied to the Dense Layer that is the activation function is linear.

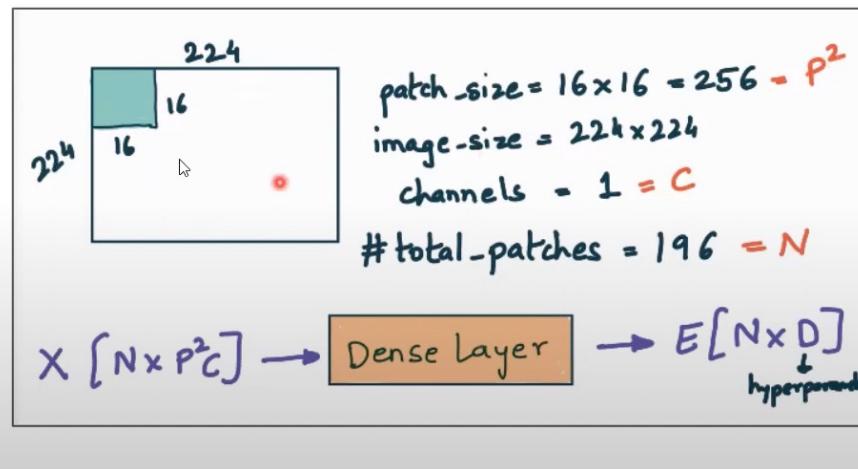
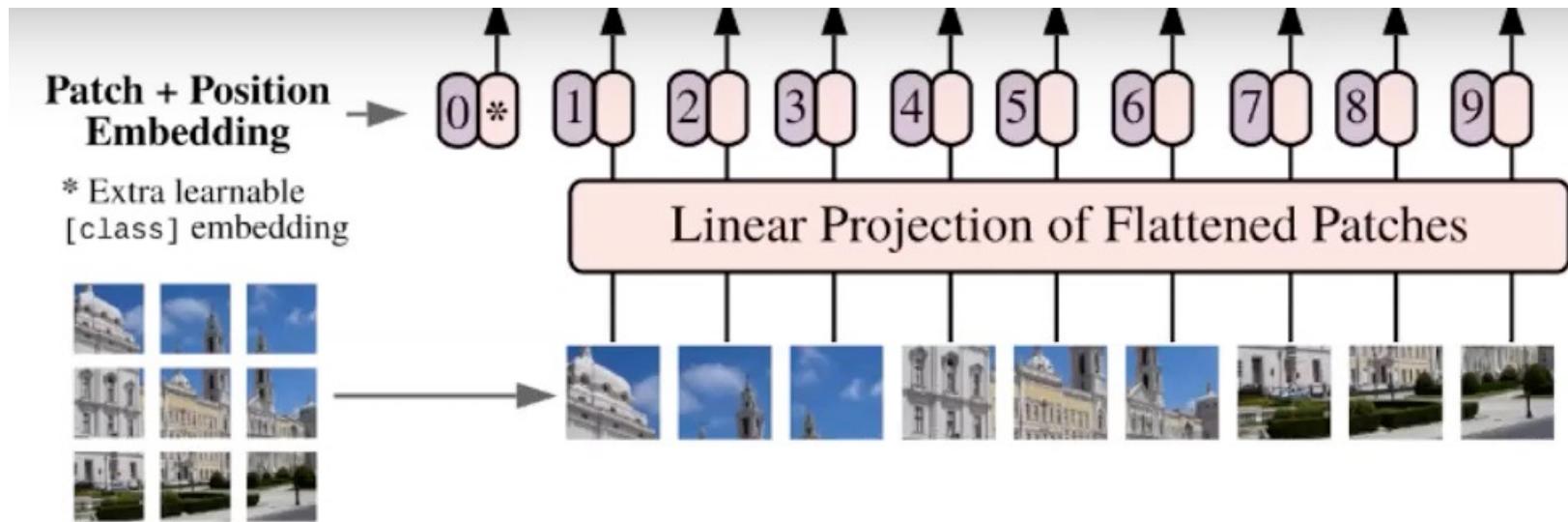
patch_size: 16x16

embed_dim : D



patches $[P \times P \times C] \Rightarrow$ vectors $[P^2 C \times 1] \Rightarrow$ vectors $[D \times 1]$

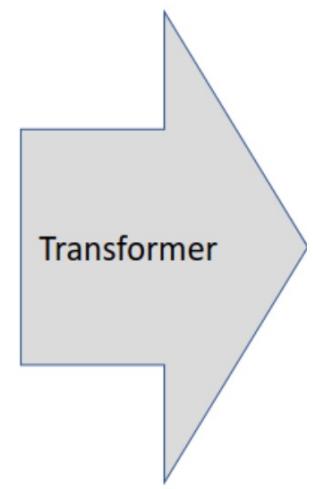
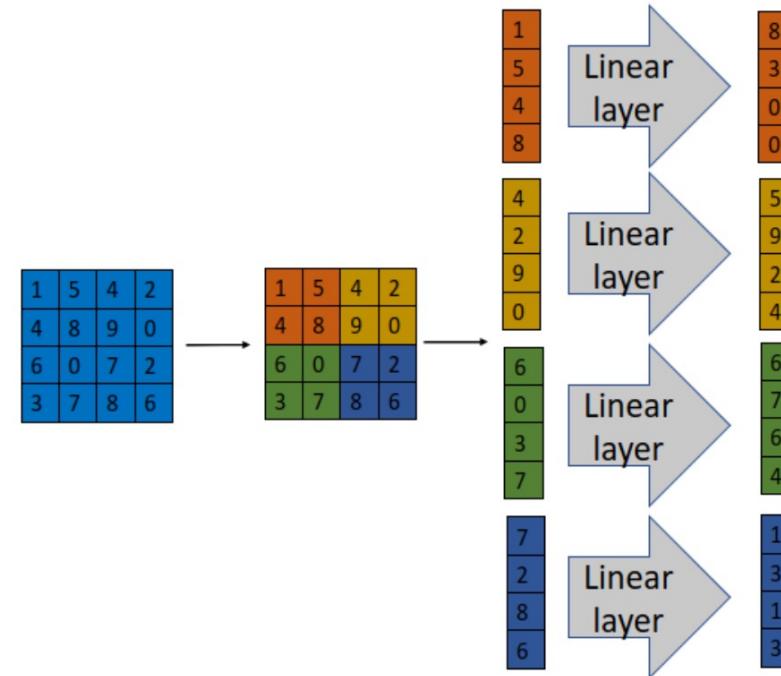
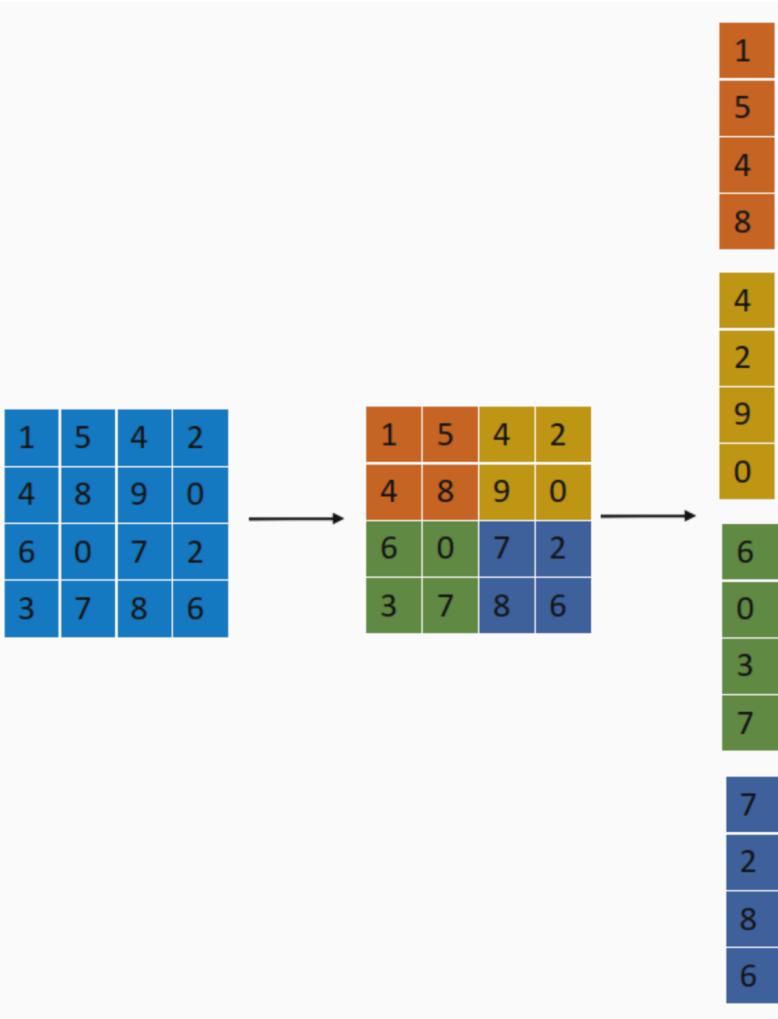


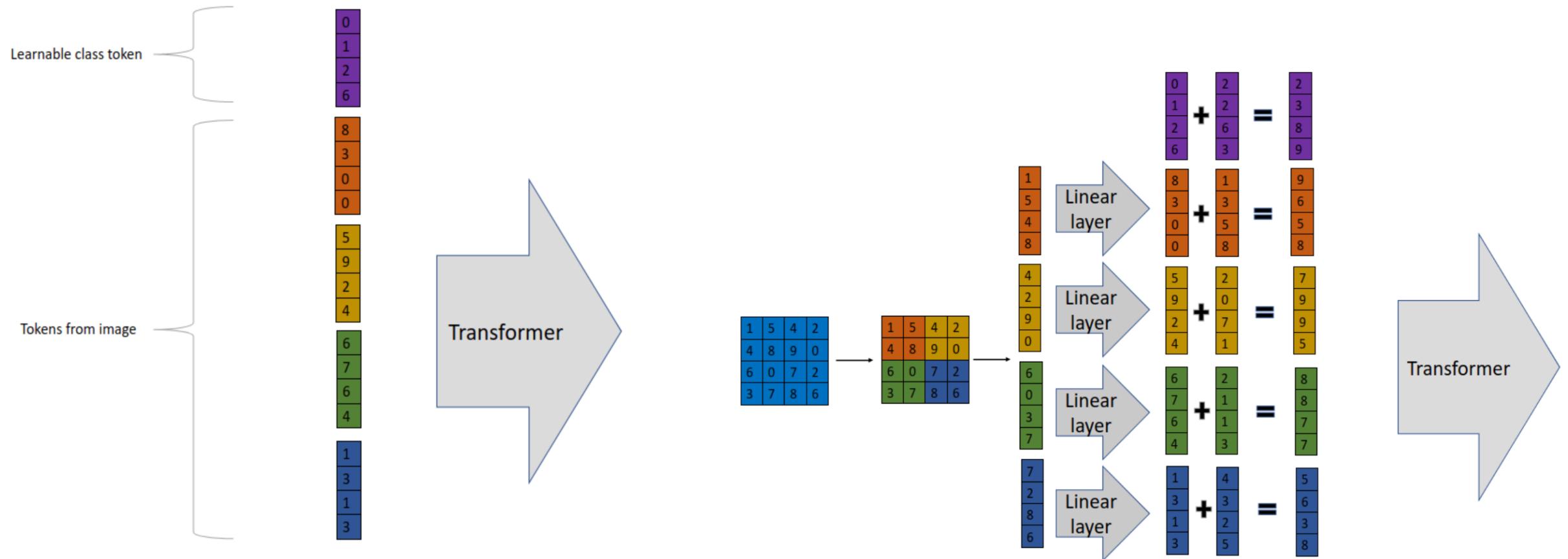


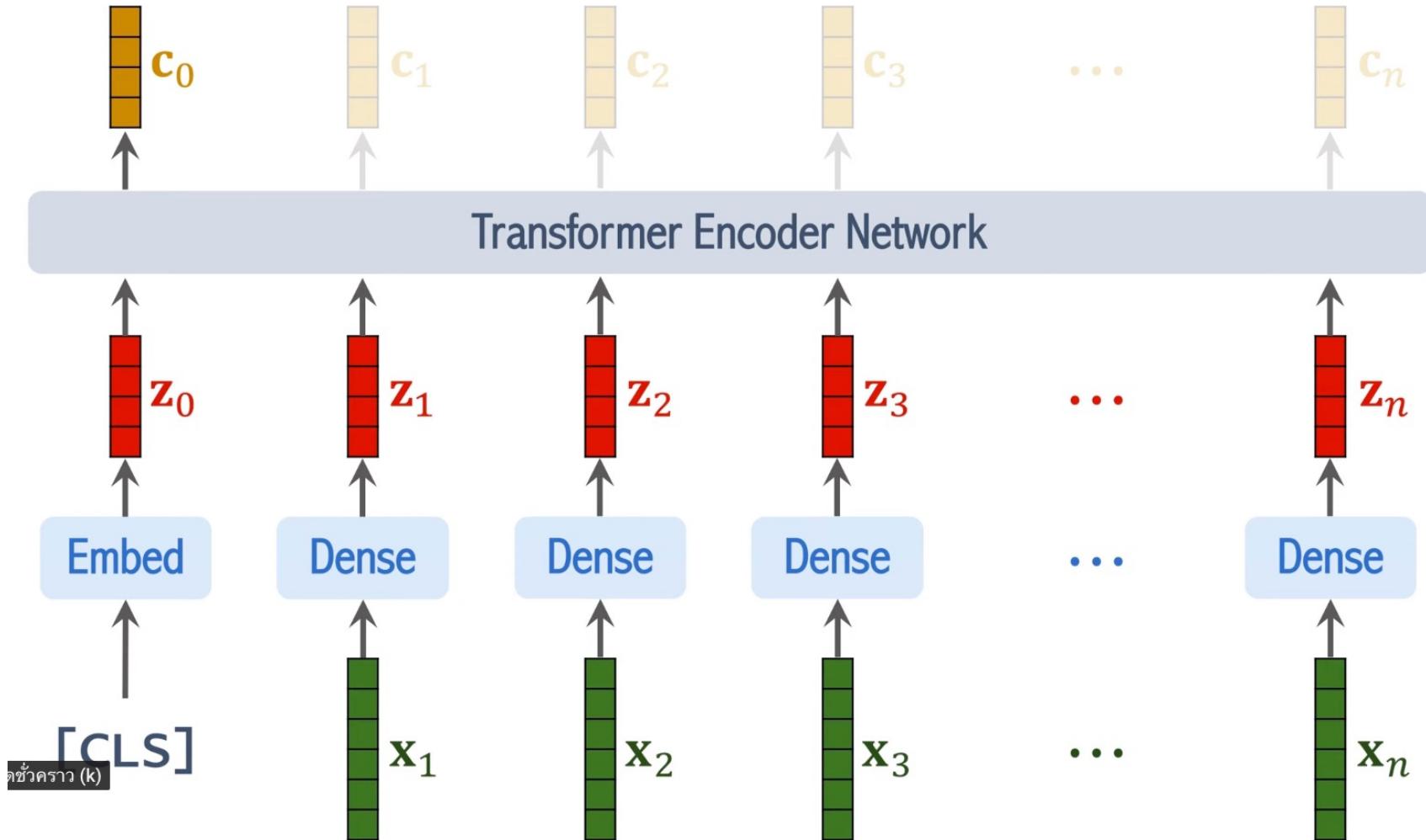
Append class embedding $\Rightarrow [1 \times D]$ randomly initialized
 $E [(N+1) \times D]$

Add positional encoding $\Rightarrow E_{\text{pos}} [(N+1) \times D]$

$E + E_{\text{pos}} = Z$





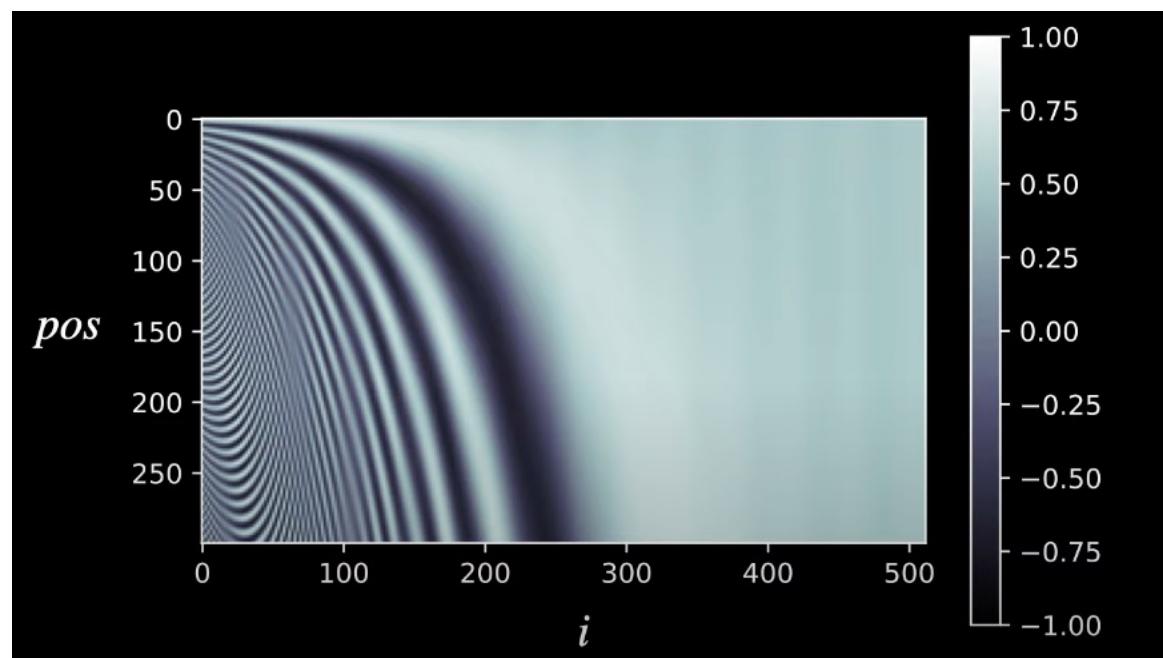


Positional encoding



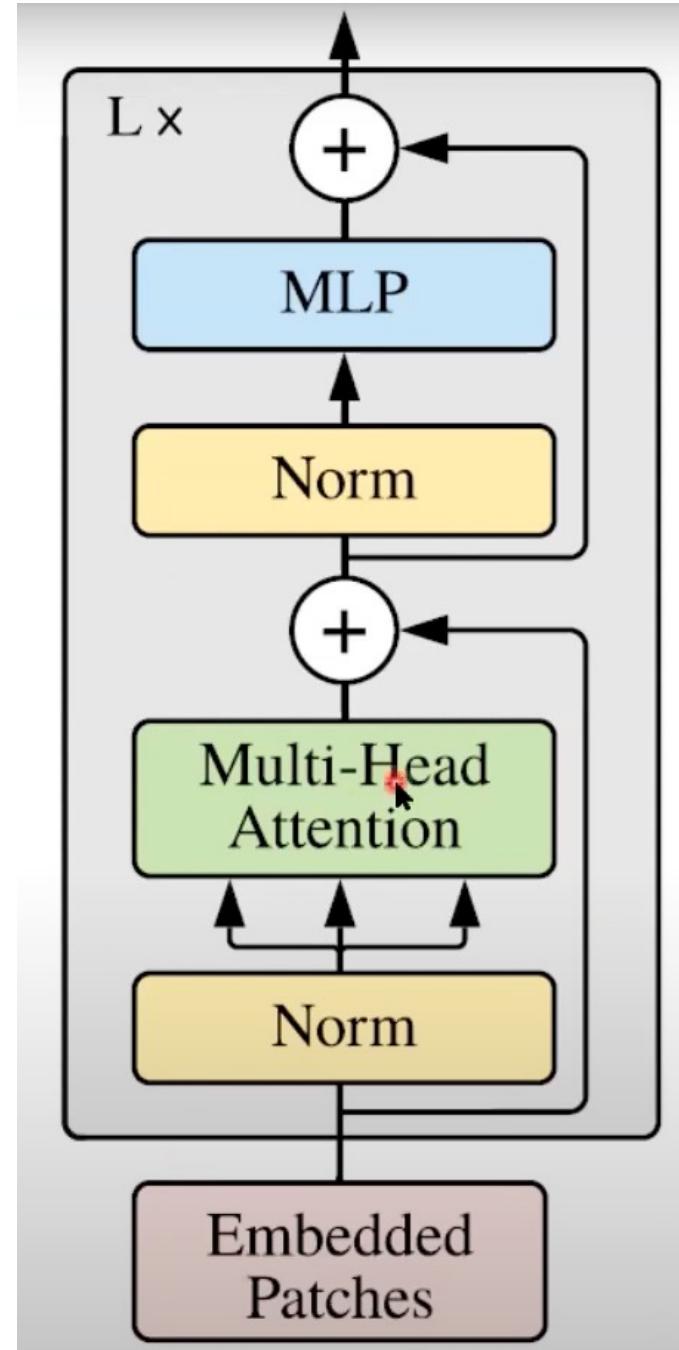
$$PE_{(pos,2i)} = \sin(pos/10000^{2i/D})$$

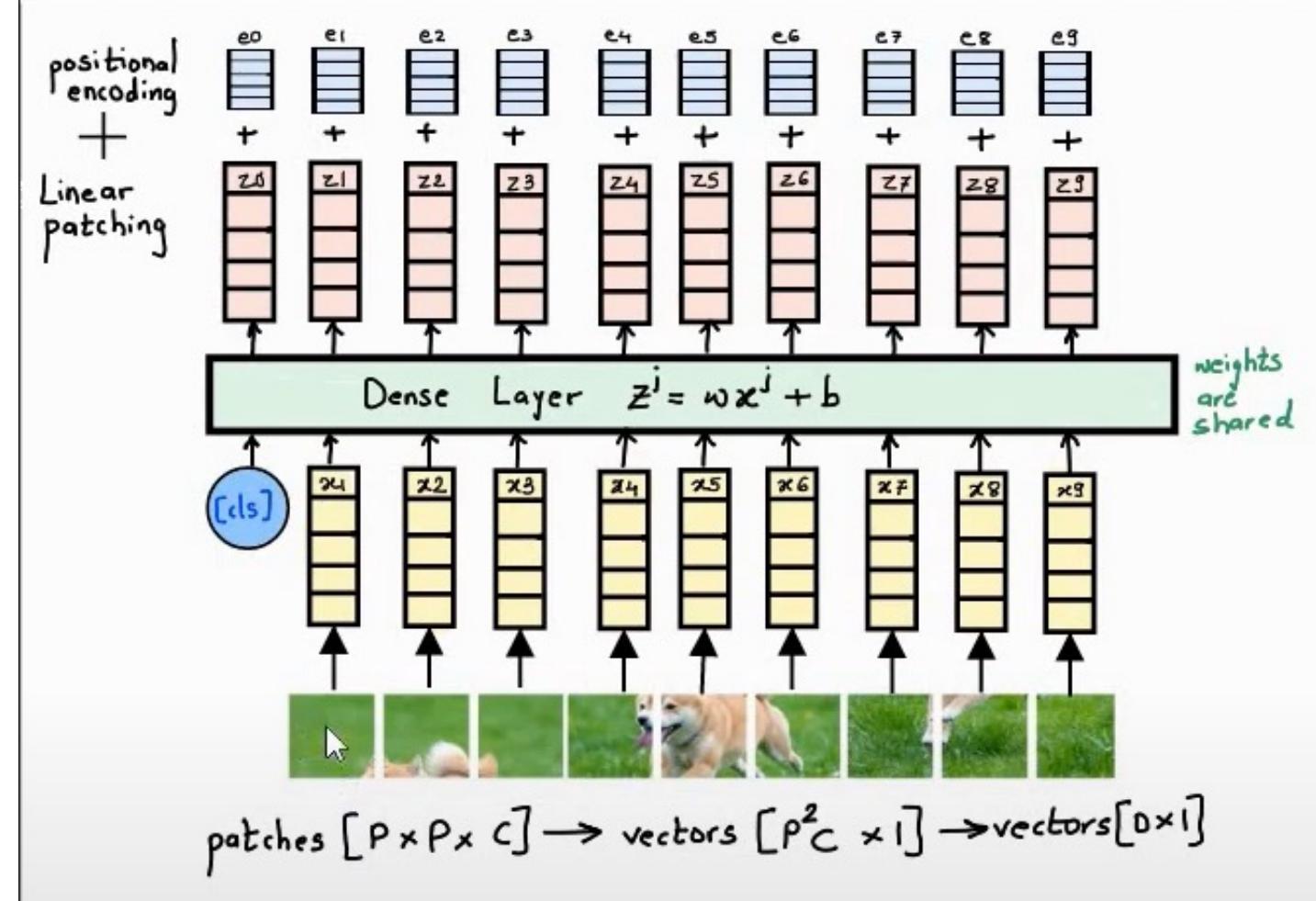
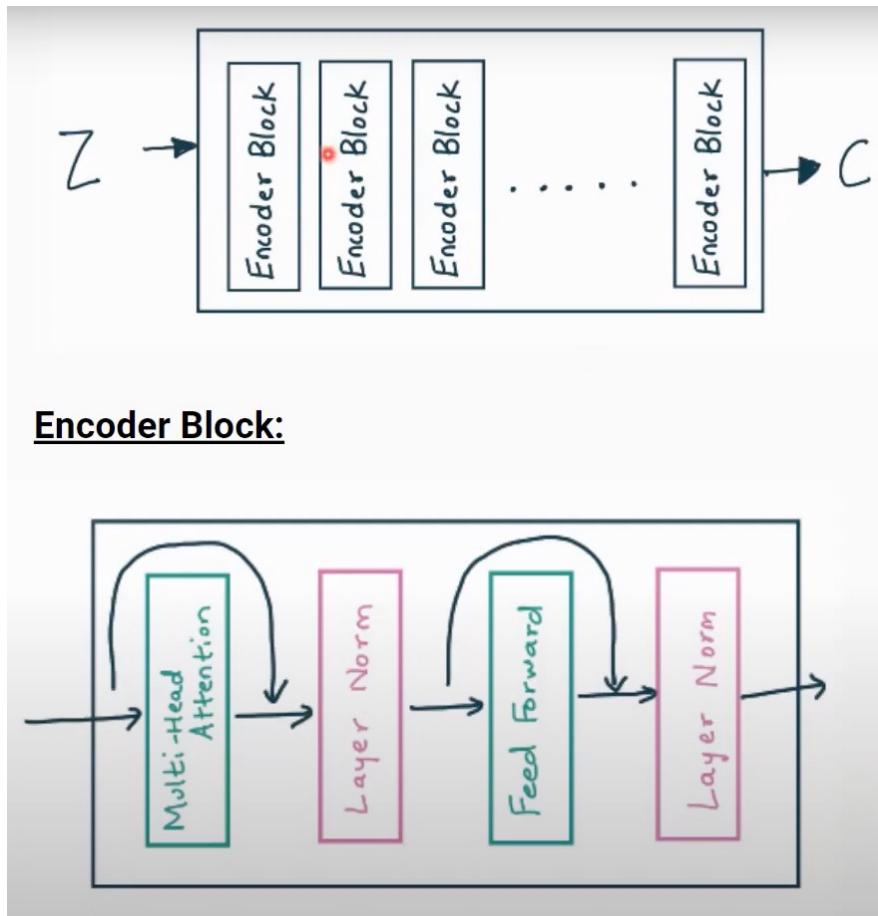
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/D})$$



Step 2: Transformer Encoder

- We then pass the embedded patches and their positional encodings through the transformer encoder architecture.





The top portion of the image shows the YouTube search interface. A search bar contains the text "k-means clustering". To the right of the search bar are several icons: a magnifying glass for search, a camera for upload, a grid for library, and a bell for notifications. In the top right corner, there is a small circular icon with an 'i' inside.

This image is identical to the one above it, but a large blue rectangular box highlights the search bar area. Inside this blue box, the text "Query (Q)" is written in white capital letters. The rest of the interface elements remain the same.



Introduction to Clustering and K-means Algorithm

Kanza Batool Haider

60K views • 5 years ago



Man Spends 30 Years Turning Degraded Land into Massive...

Happen Films

1.1M views • 1 year ago



The Reality of Reality: A Tale of Five Senses

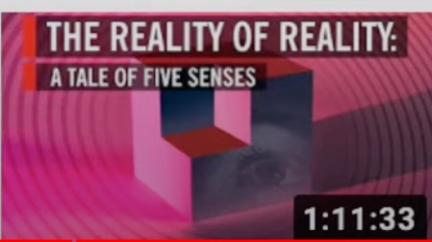
World Science Festival

198K views • 1 year ago



Query (Q)

Most Similar

**Introduction to Clustering and K-means Algorithm**Kanza Batool Haider
60K views • 5 years agoKey (K_1)**Man Spends 30 Years Turning Degraded Land into Massive...**Happen Films
1.1M views • 1 year agoKey (K_2)**The Reality of Reality: A Tale of Five Senses**World Science Festival
198K views • 1 year agoKey (K_3)



YouTube

k-means clustering

Query (Q)



Introduction to Clustering and K-means Algorithm

Kanza Batool Haider

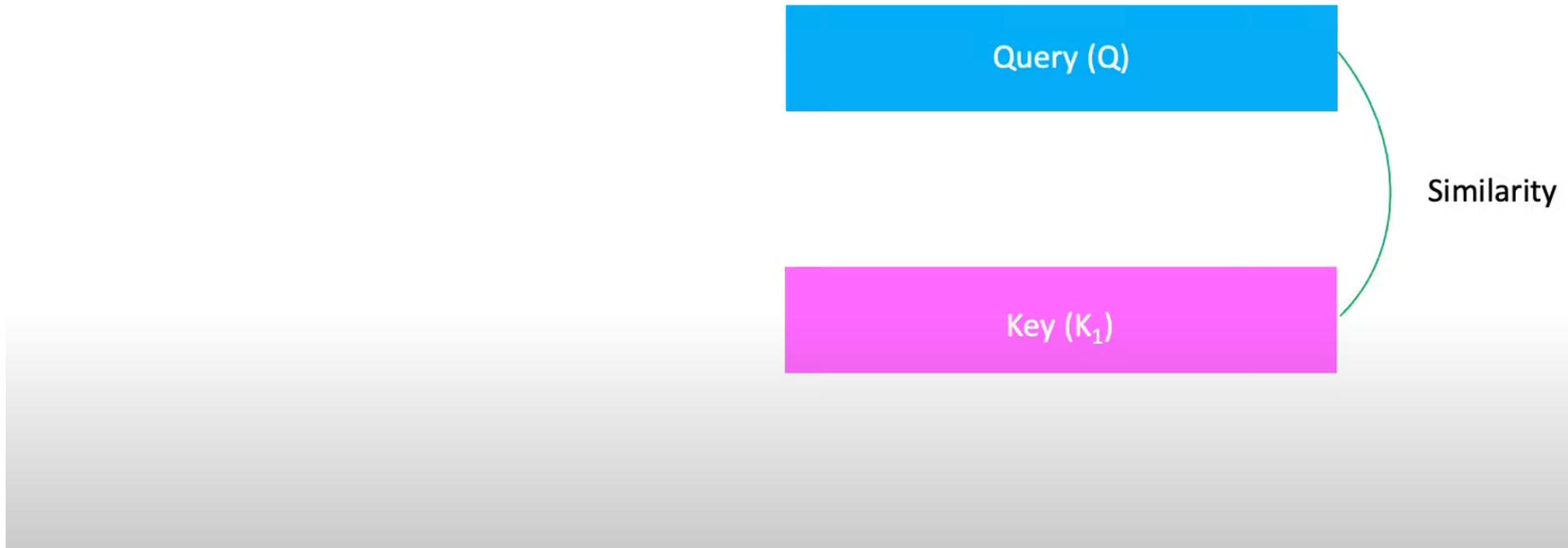
60K views • 5 years ago

Key (K_1)

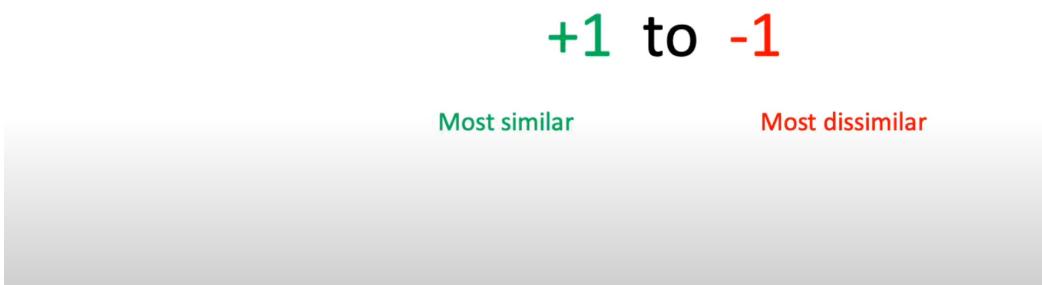
Video Contents

Value (V_1)

Similarity?



Cosine Similarity



Cosine Similarity



Cosine Similarity



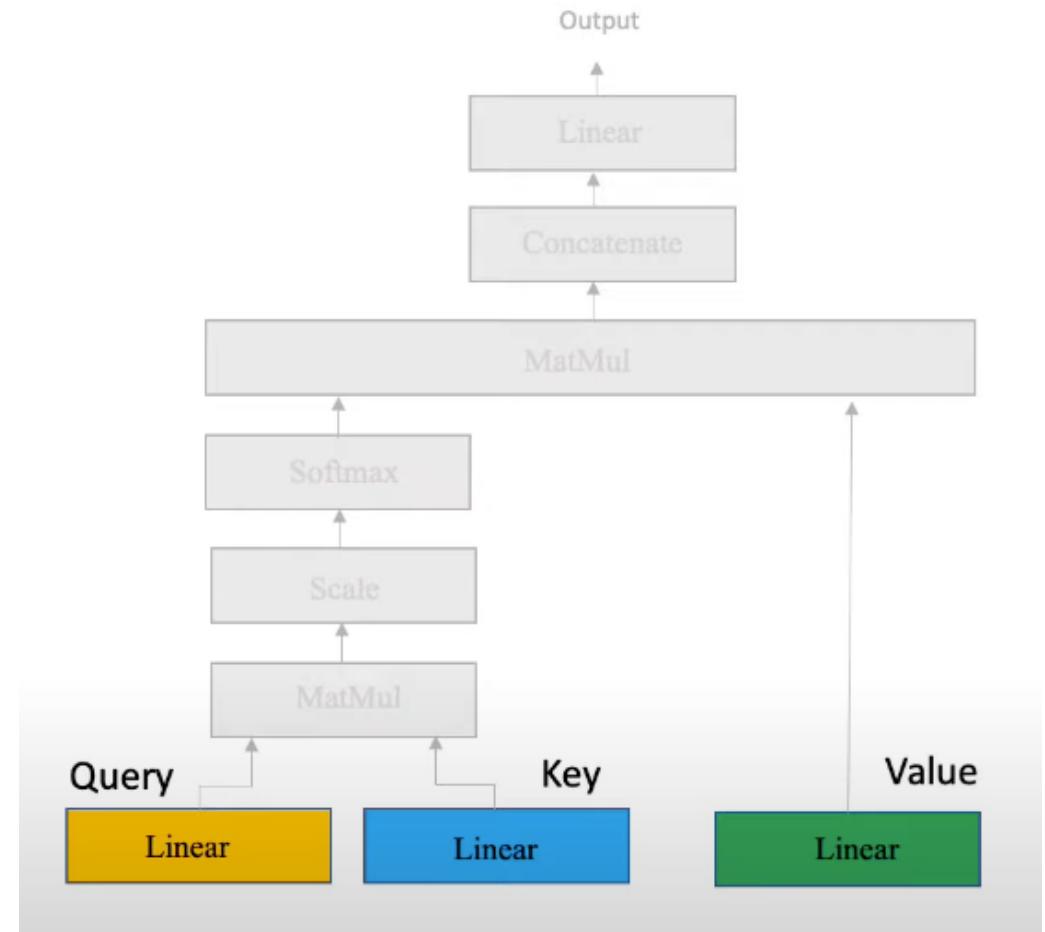
$$\text{Cos}(A, B) = \frac{A \cdot B}{|A| |B|}$$

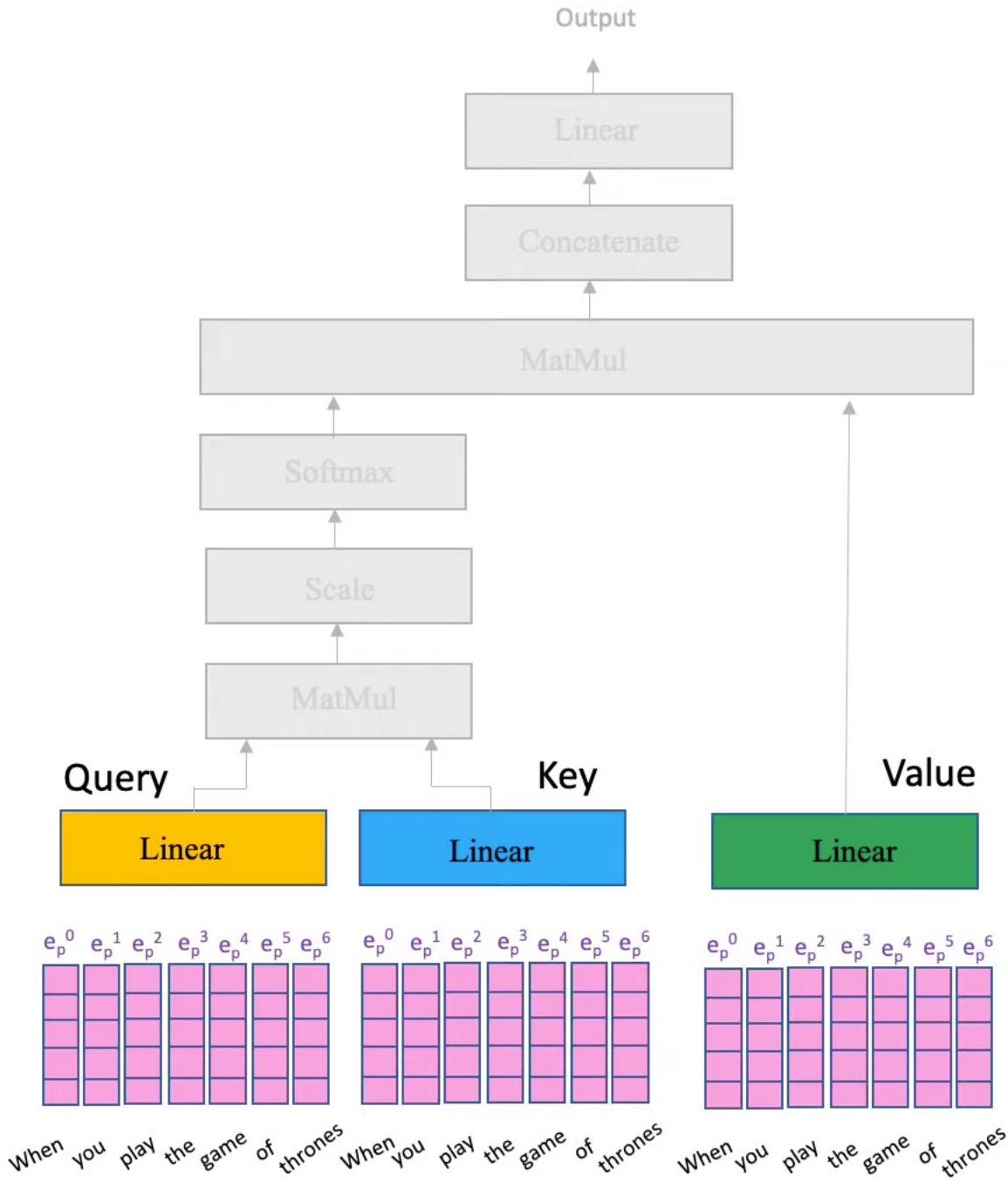
Similarity b/w Matrices

$$\text{similarity } (A, B) = \frac{A \cdot B^T}{\text{scaling}}$$

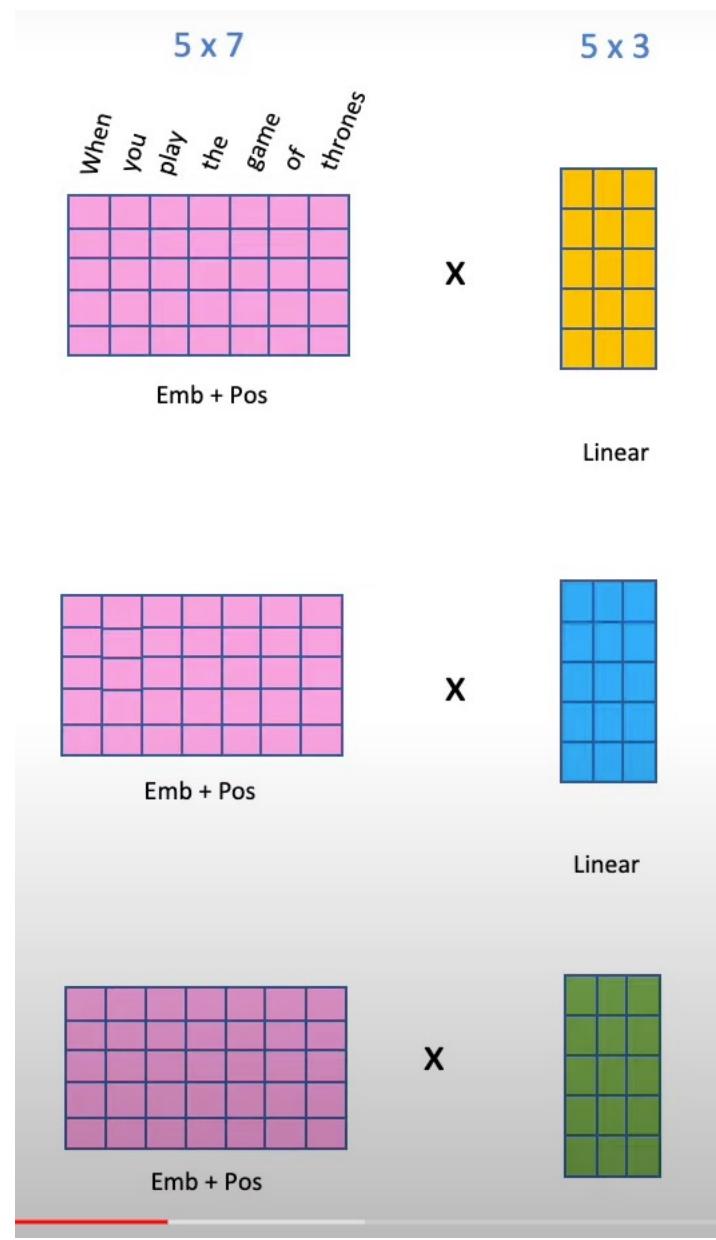
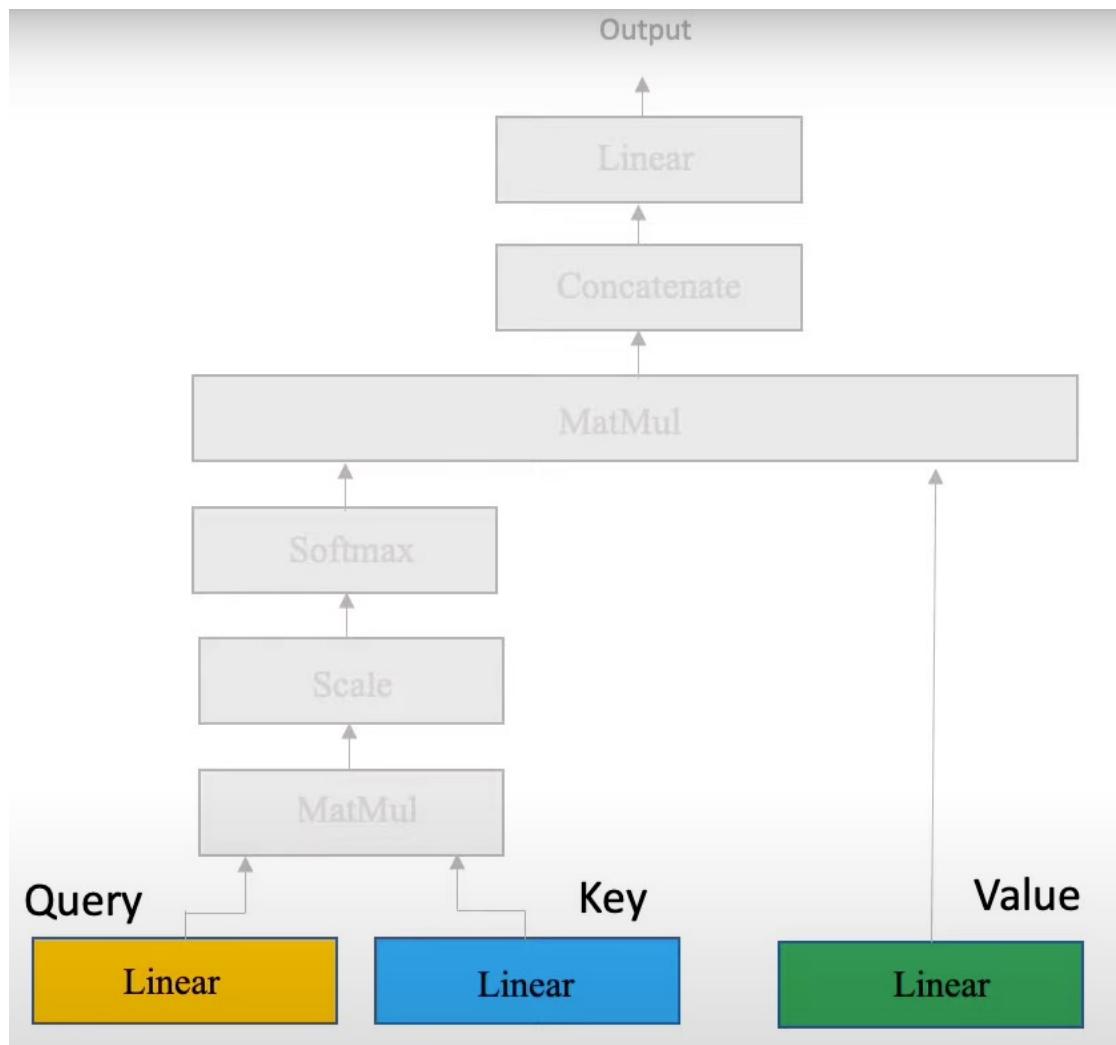
Similarity b/w **Query** and **Key**

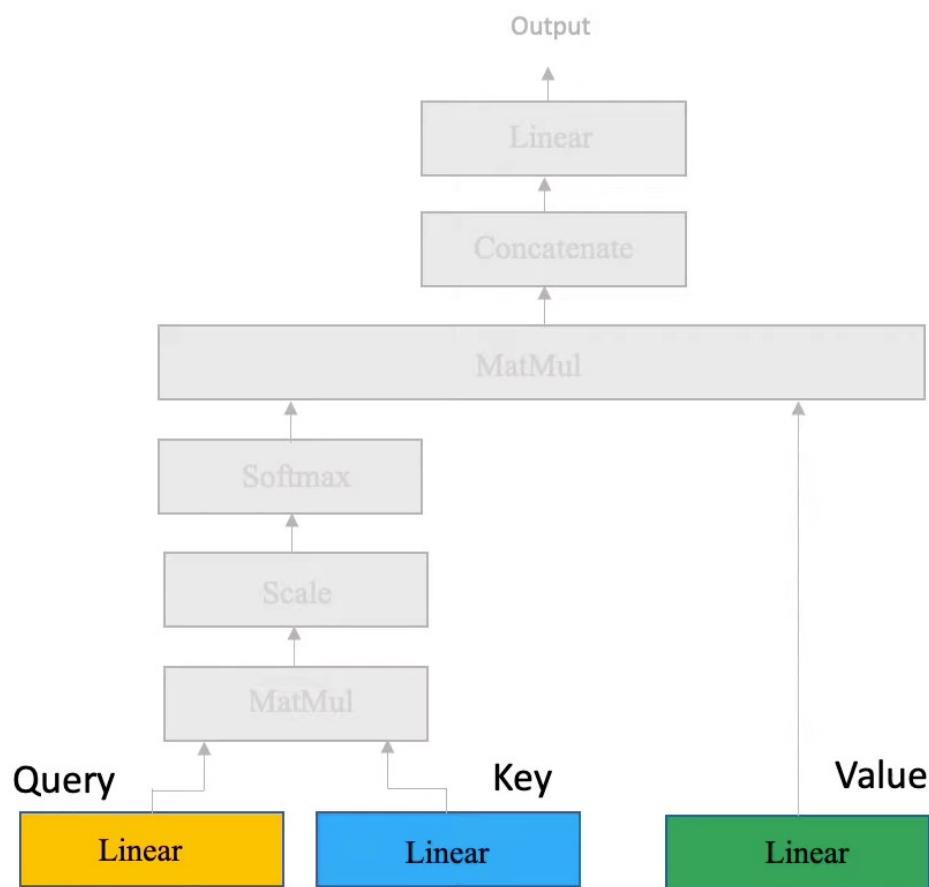
$$\text{similarity } (Q, K) = \frac{Q \cdot K^T}{\text{scaling}}$$



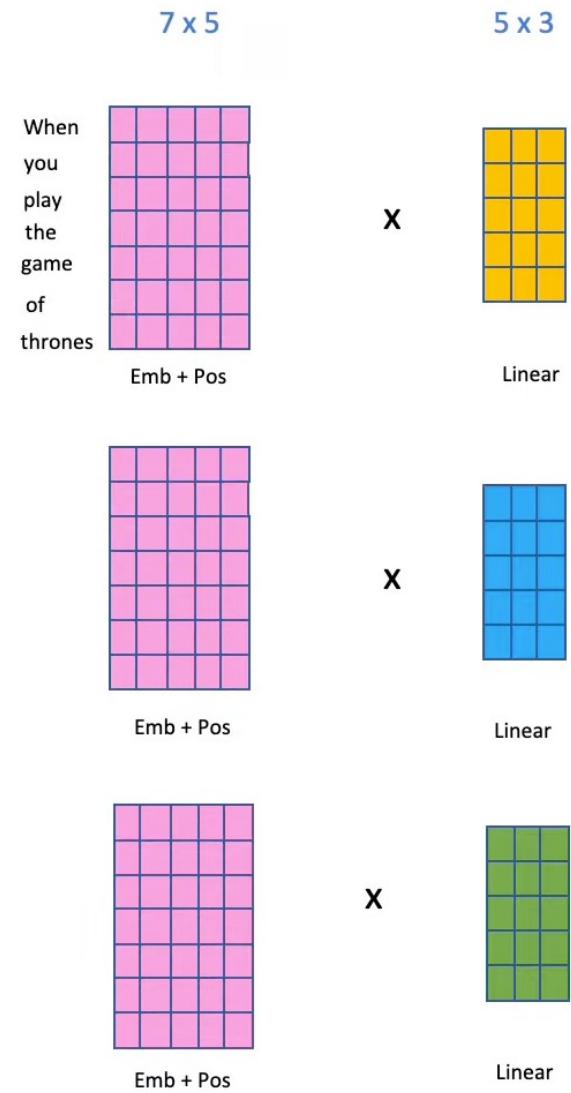


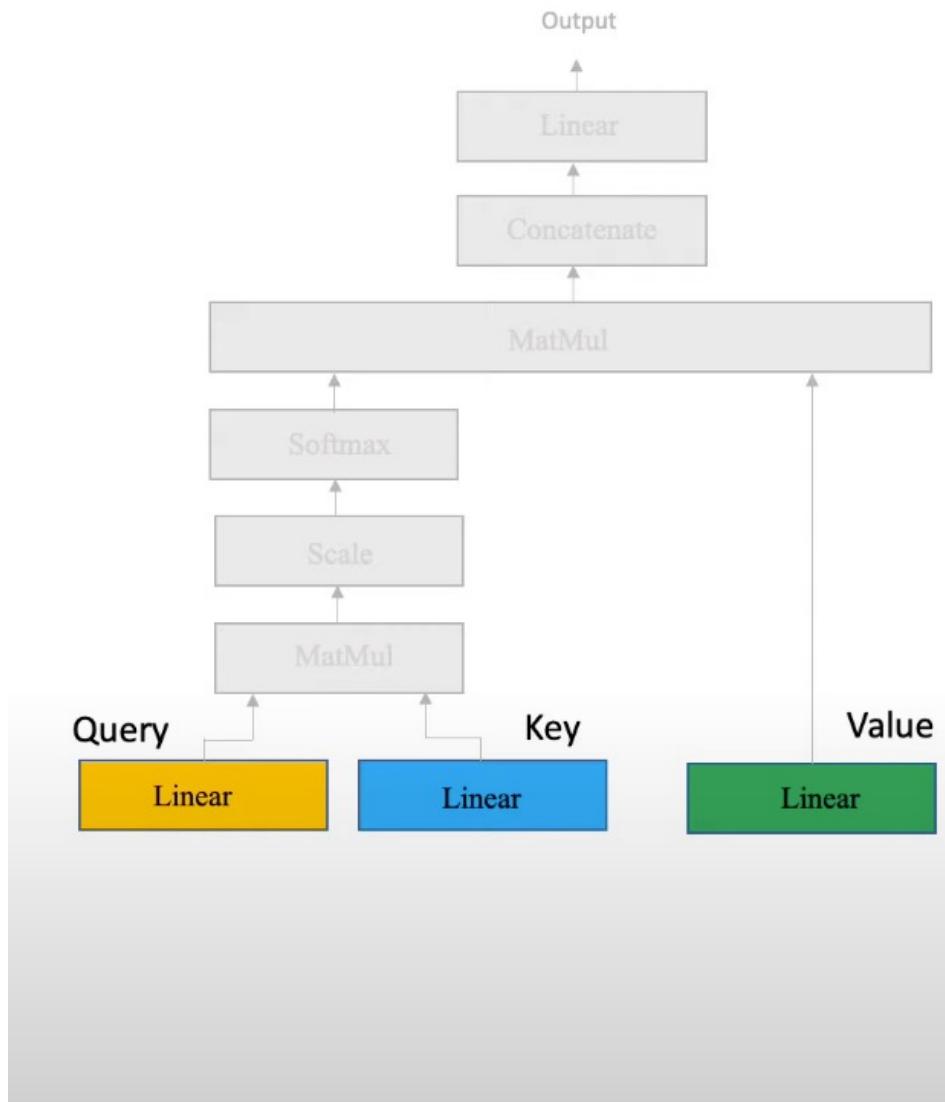
Multi-Head Attention



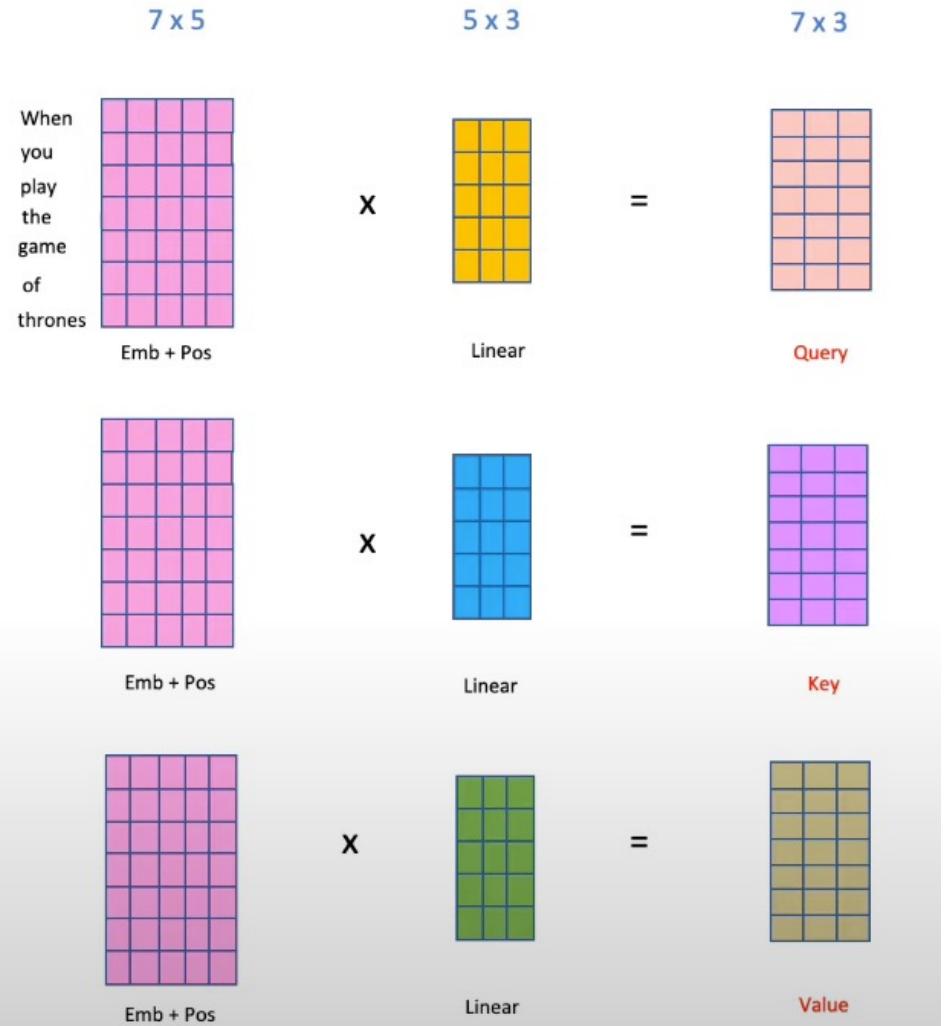


Multi-Head Attention

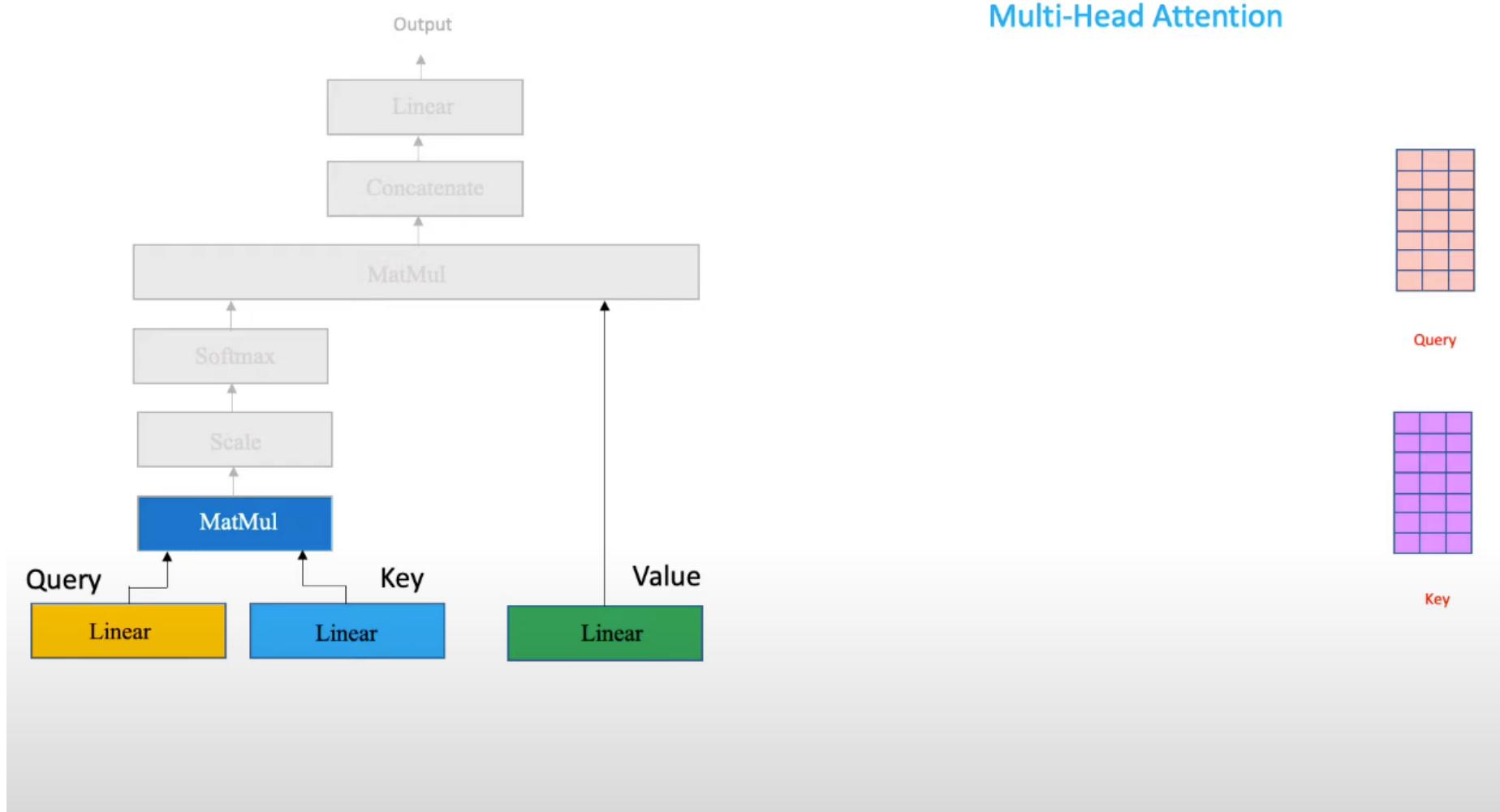




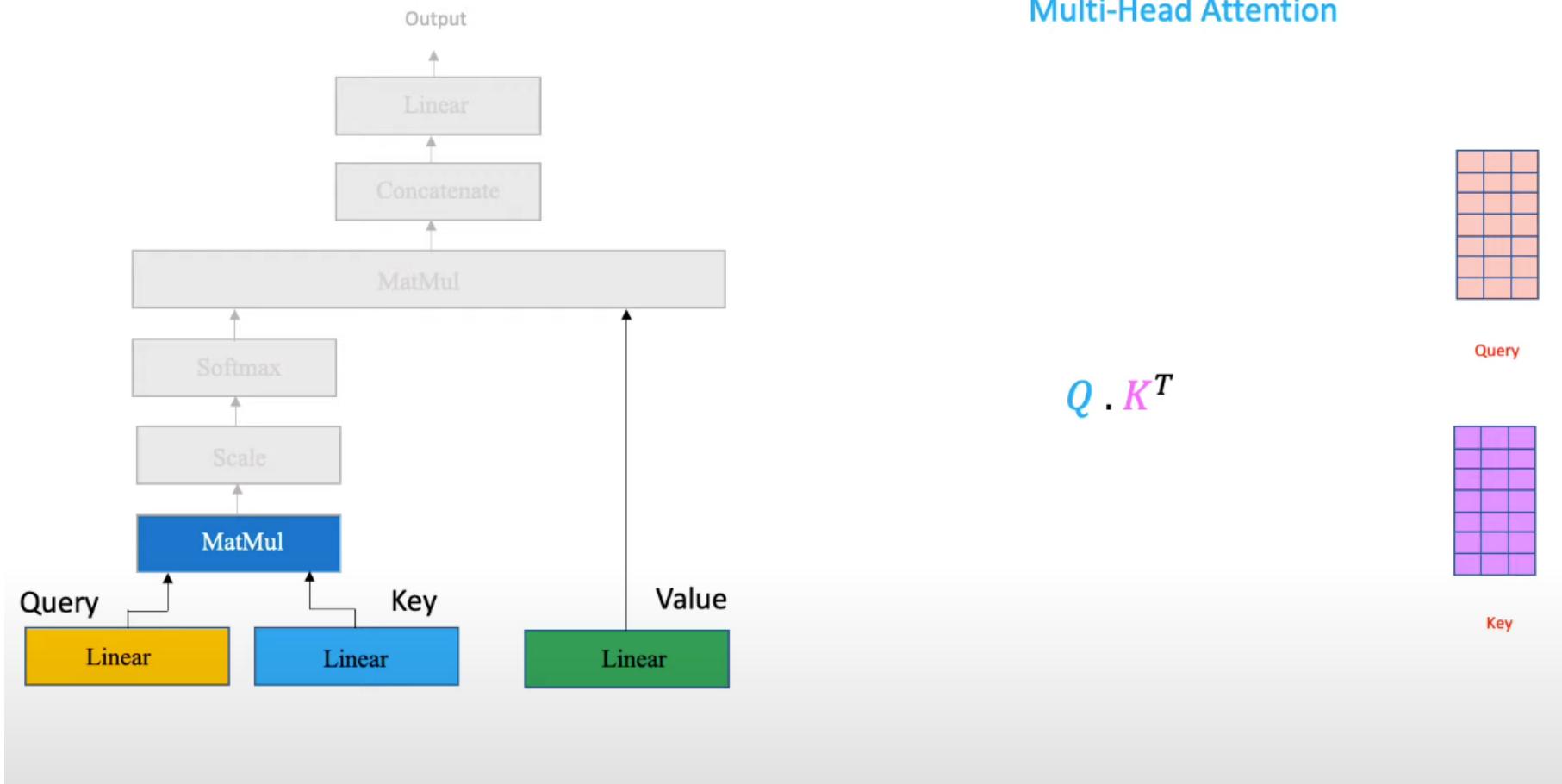
Multi-Head Attention

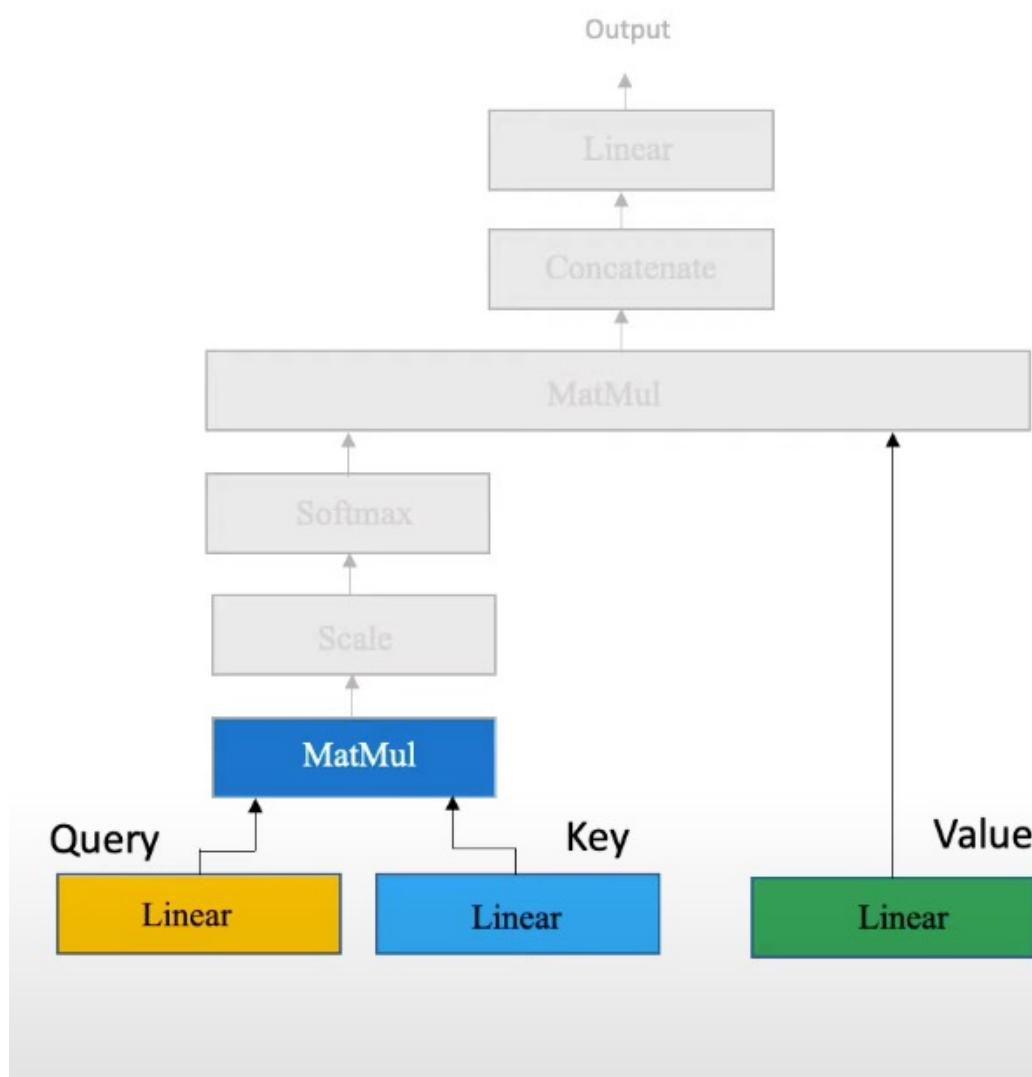


Multi-Head Attention



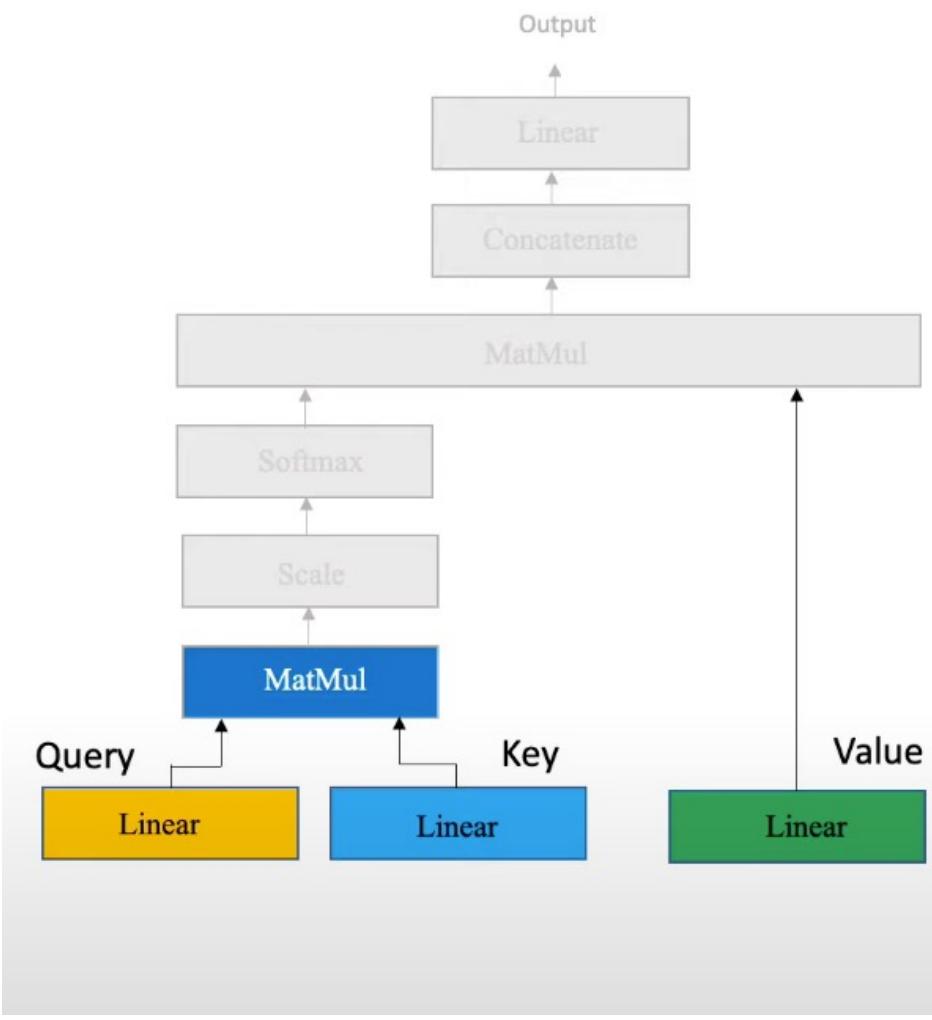
Multi-Head Attention





Multi-Head Attention

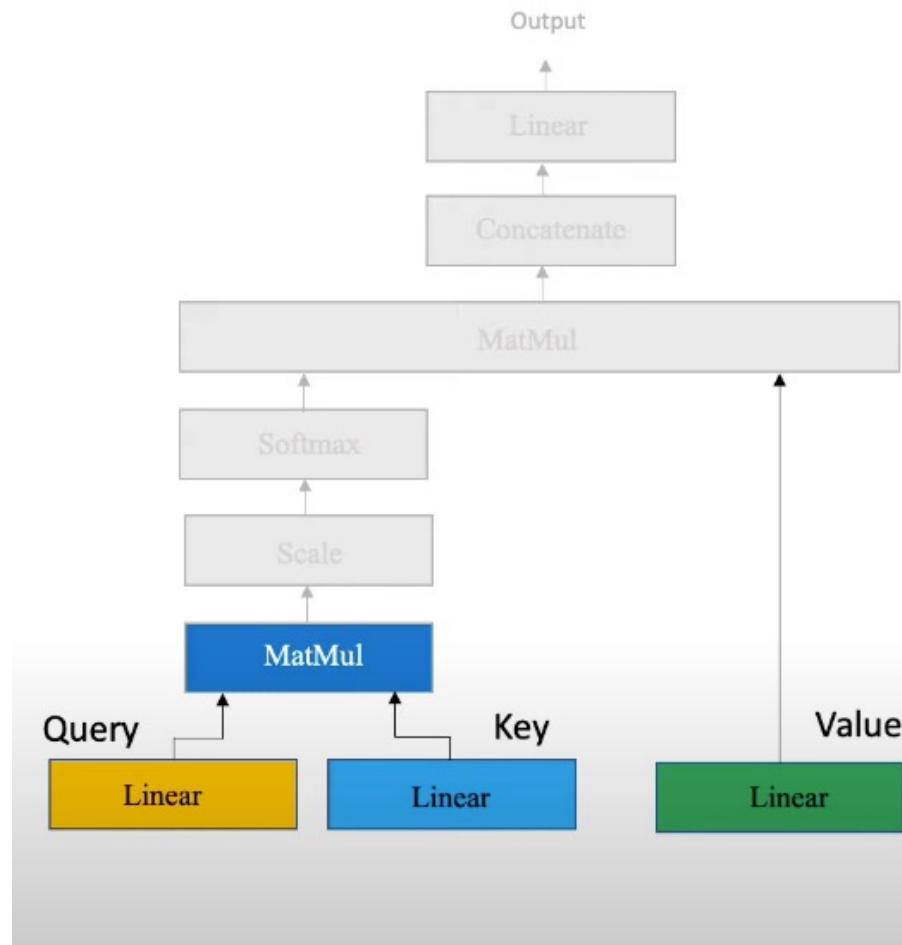
$$\begin{array}{c}
 7 \times 3 \quad 3 \times 7 \\
 \begin{matrix} \text{Query} \\ \text{Key}^T \end{matrix} * \begin{matrix} \text{Query} \\ \text{Key}^T \end{matrix}
 \end{array}$$



Multi-Head Attention

$$\begin{array}{ccc}
 7 \times 3 & \times & 3 \times 7 \\
 \text{Query} & * & \text{Key}^T \\
 & & = \\
 & & 7 \times 7
 \end{array}$$

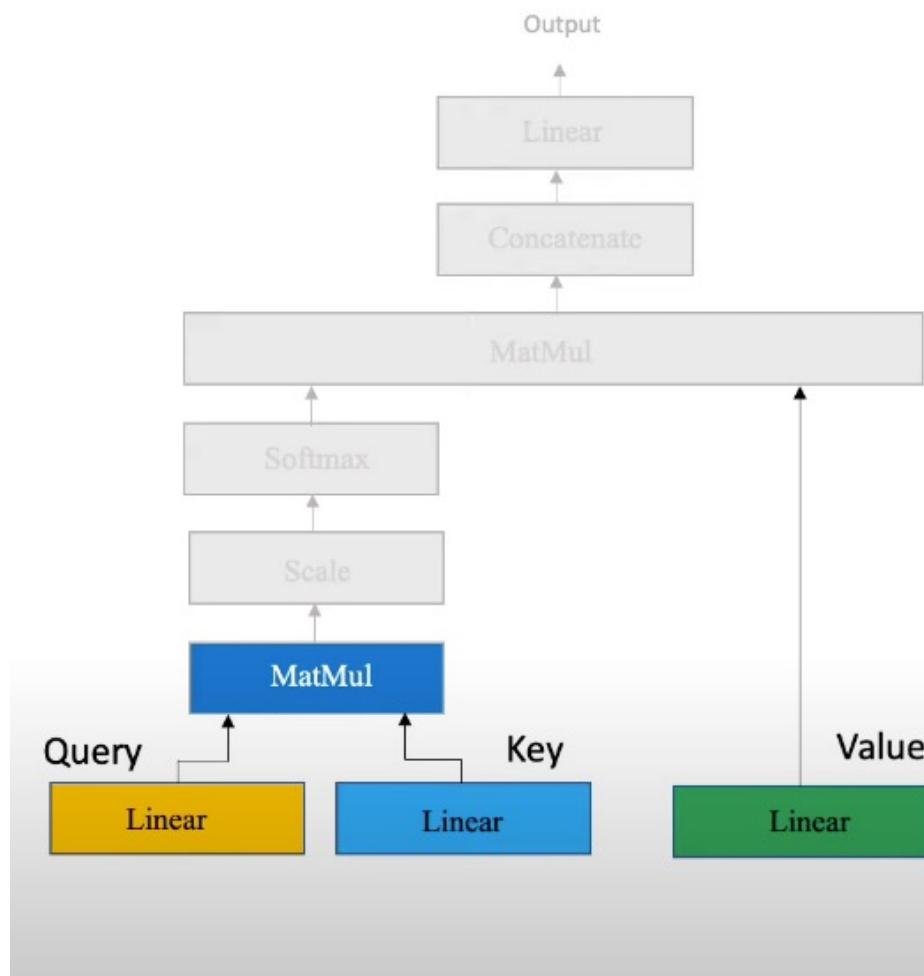
Attention Filter



Multi-Head Attention

7 x 7

	When	you	play	the	game	of	thrones
When	89	20	41	10	55	10	59
you	20	90	81	22	70	15	72
play	41	81	95	10	90	30	92
the	10	22	10	92	88	40	89
game	55	70	90	88	98	44	87
of	10	15	30	40	44	85	59
thrones	59	72	92	90	95	59	99

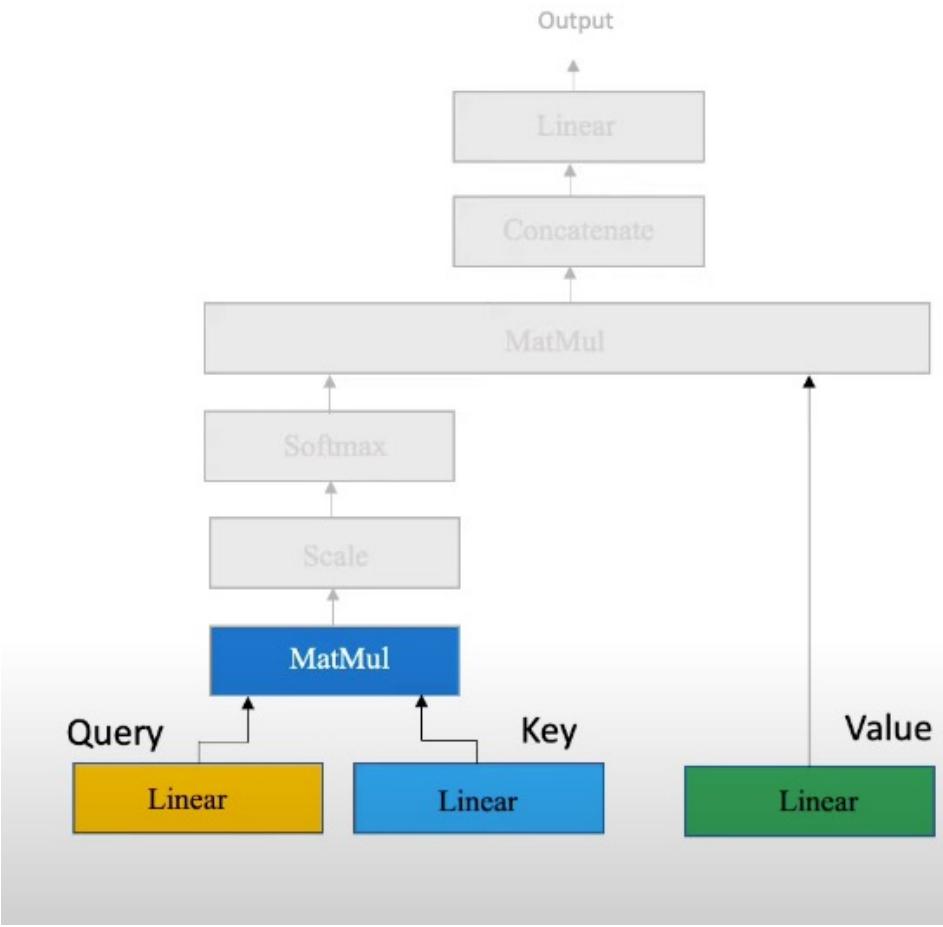


Multi-Head Attention

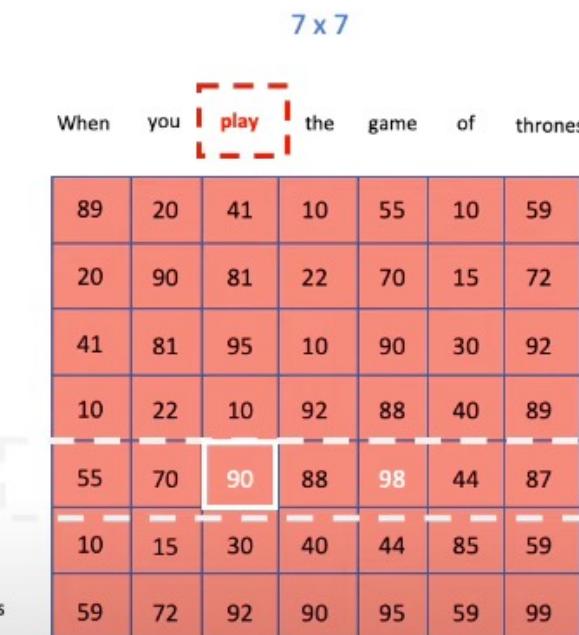
7 x 7

When you play the game of thrones

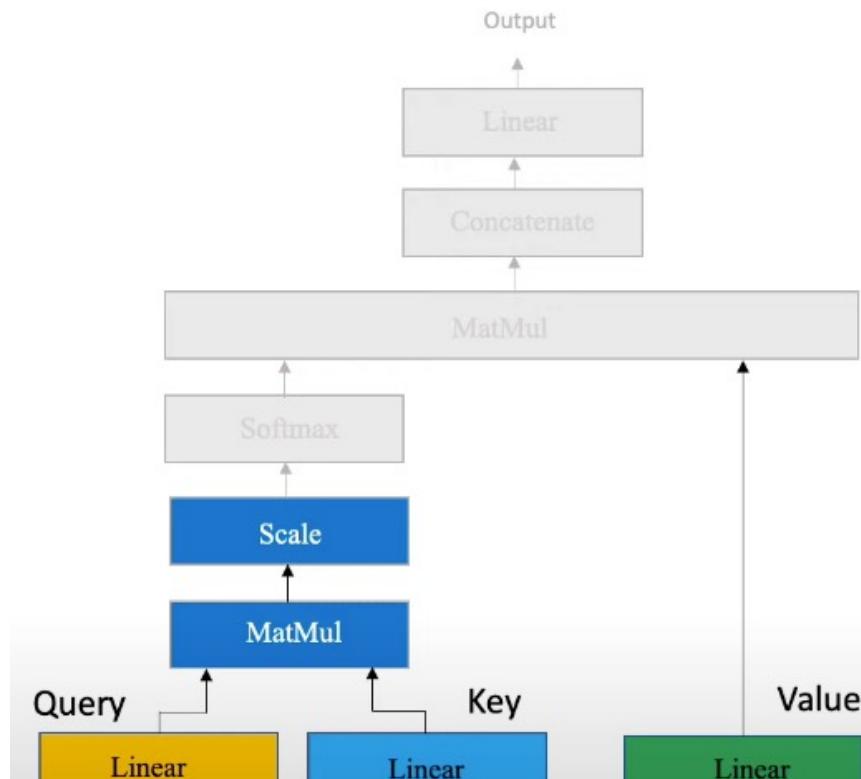
When	89	20	41	10	55	10	59
you	20	90	81	22	70	15	72
play	41	81	95	10	90	30	92
the	10	22	10	92	88	40	89
game	55	70	90	88	98	44	87
of	10	15	30	40	44	85	59
thrones	59	72	92	90	95	59	99



Multi-Head Attention



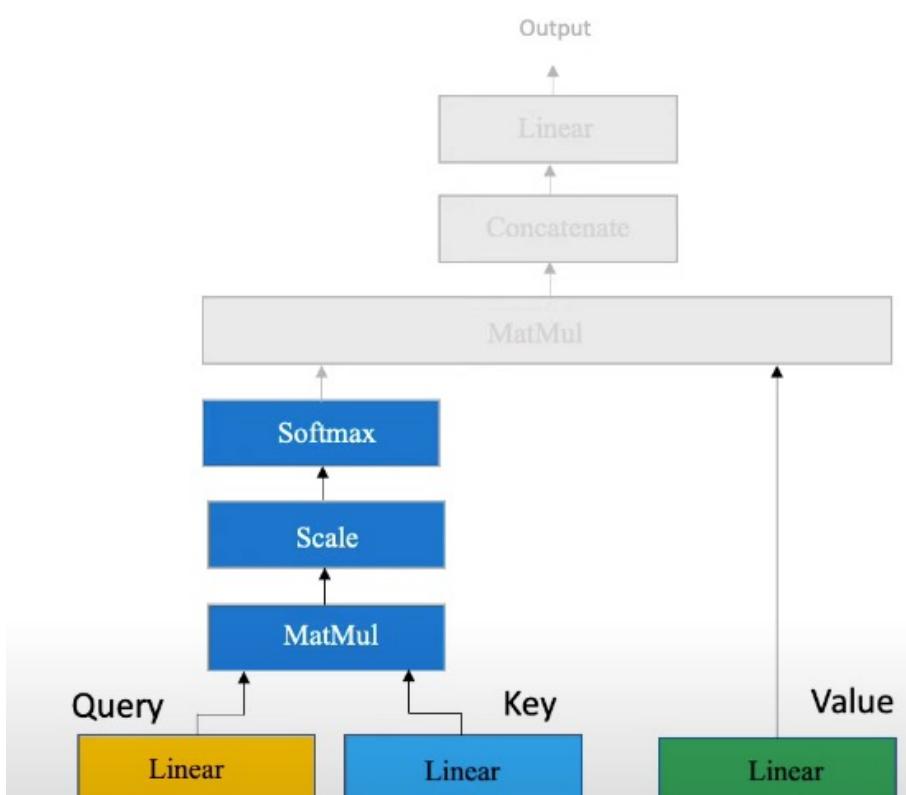
Multi-Head Attention



7×7

	When	you	play	the	game	of	thrones
When	89	20	41	10	55	10	59
you	20	90	81	22	70	15	72
play	41	81	95	10	90	30	92
the	10	22	10	92	88	40	89
game	55	70	90	88	98	44	87
of	10	15	30	40	44	85	59
thrones	59	72	92	90	95	59	99

$$\sqrt{d_k}$$

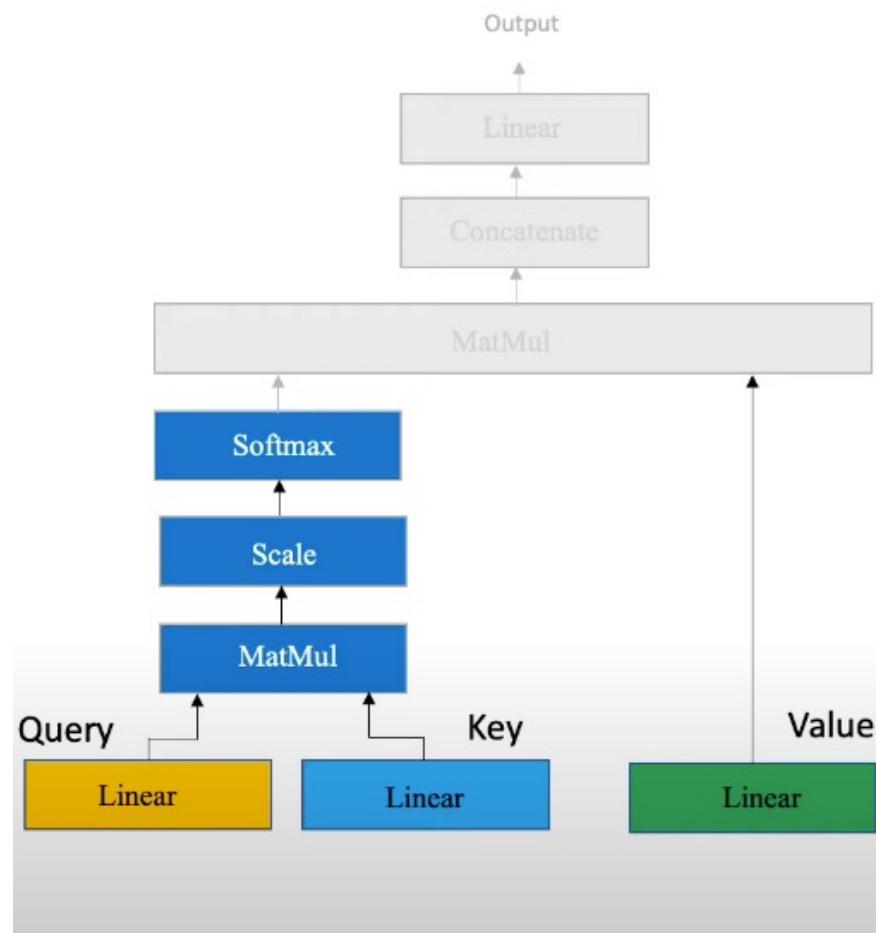


Multi-Head Attention

7 x 7

	When	you	play	the	game	of	thrones
When	33.6	7.6	15.5	3.8	20.8	3.8	22.3
you	7.6	34.0	30.6	8.3	26.5	5.7	27.2
play	15.5	30.6	35.9	3.8	34.0	11.3	34.8
the	3.8	8.3	3.8	34.8	33.3	15.1	33.6
game	20.8	26.5	34.0	33.3	37.0	16.6	35.9
of	3.8	5.7	11.3	15.1	16.6	32.1	22.3
thrones	22.3	27.2	34.8	34.0	35.9	22.3	37.4

$$\text{softmax } (x_i) = \frac{\exp(x_i)}{\sum_j \exp(x_i)}$$



Multi-Head Attention

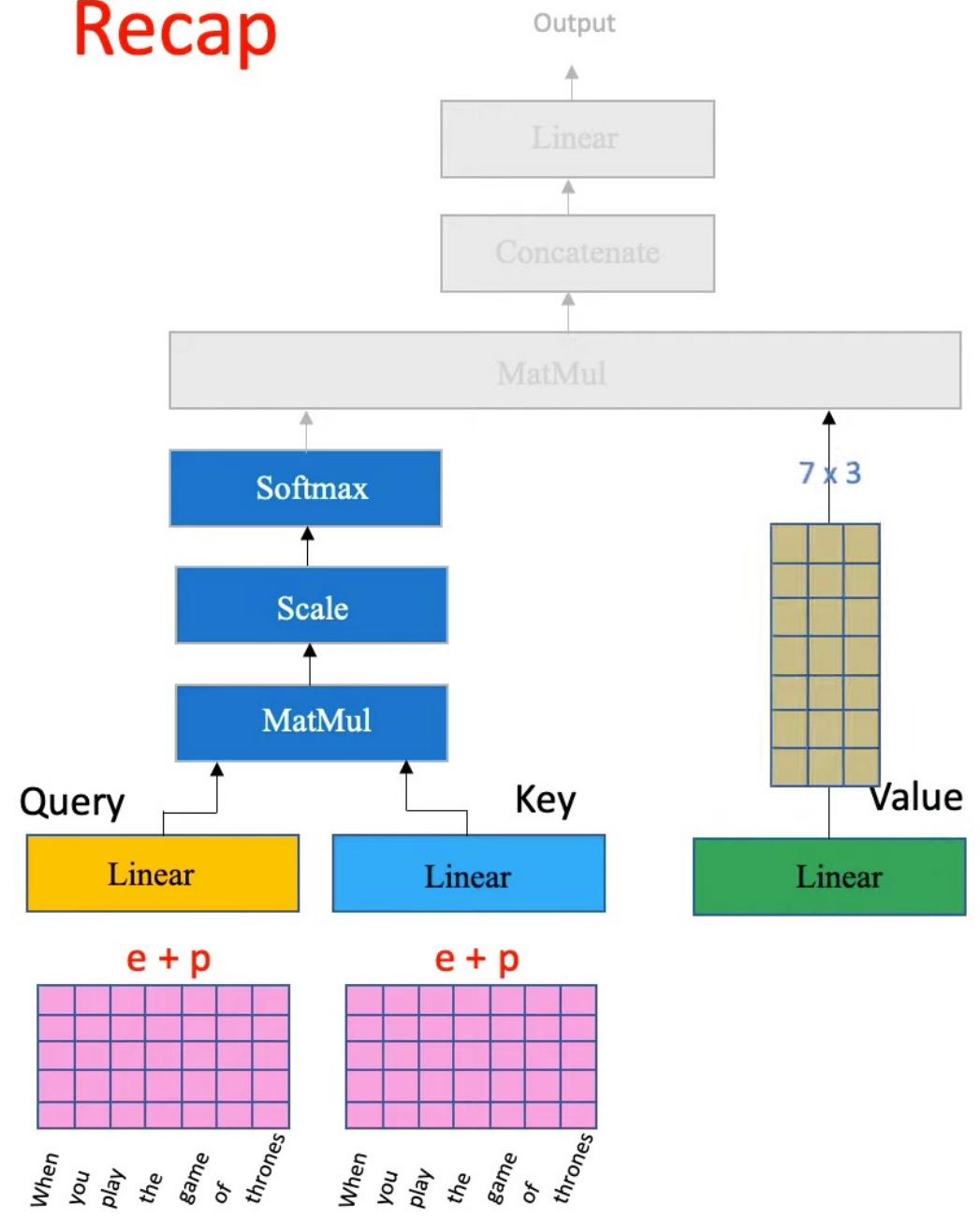
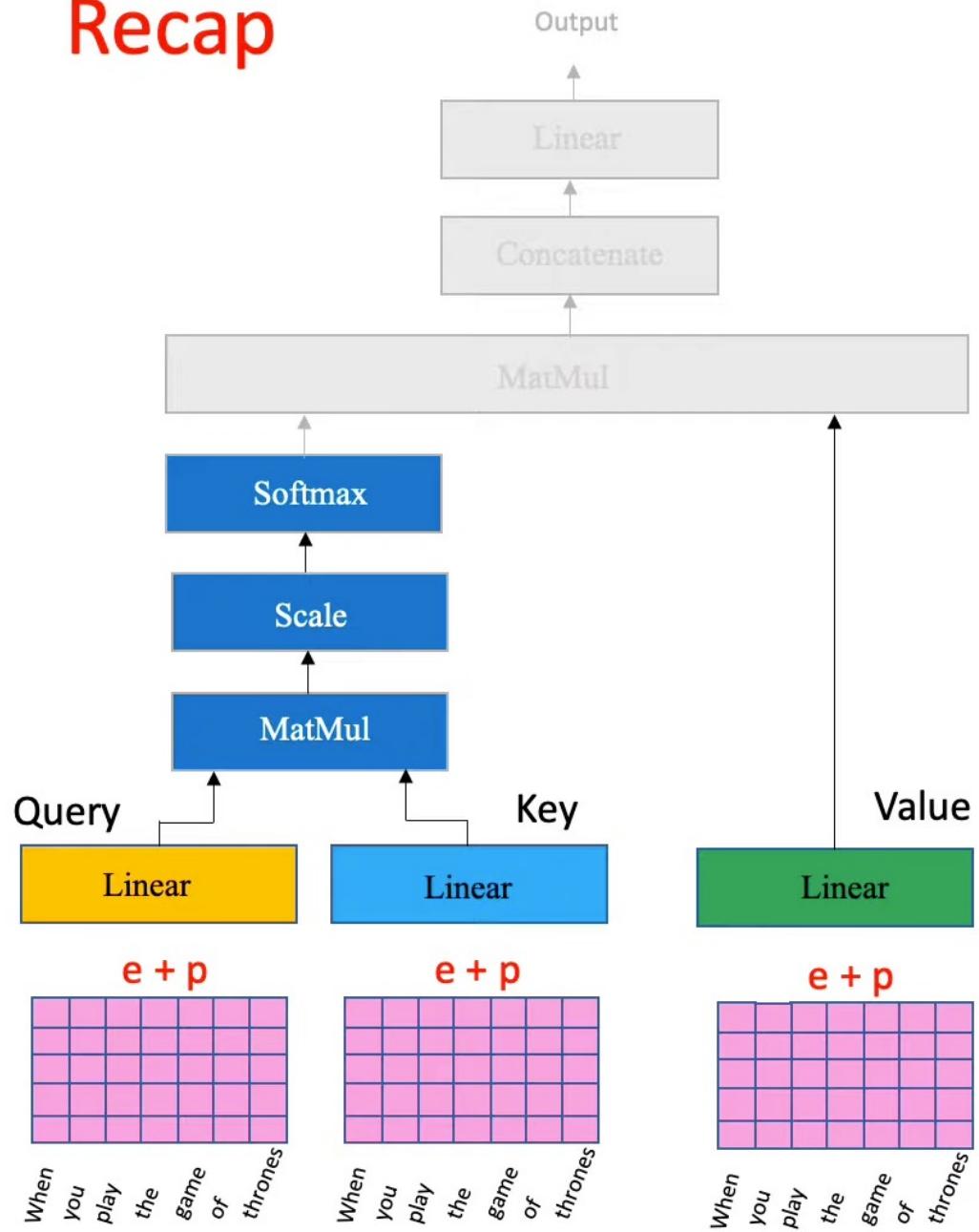
7×7

When you play the game of thrones							
When	1.00	0.00	0.00	0.00	0.00	0.00	0.00
you	0.00	0.97	0.03	0.00	0.00	0.00	0.00
play	0.00	0.00	0.68	0.00	0.10	0.00	0.22
the	0.00	0.00	0.00	0.65	0.14	0.00	0.21
game	0.00	0.00	0.03	0.02	0.72	0.00	0.23
of	0.00	0.00	0.00	0.00	0.00	1.00	0.00
thrones	0.00	0.00	0.05	0.03	0.17	0.00	0.75

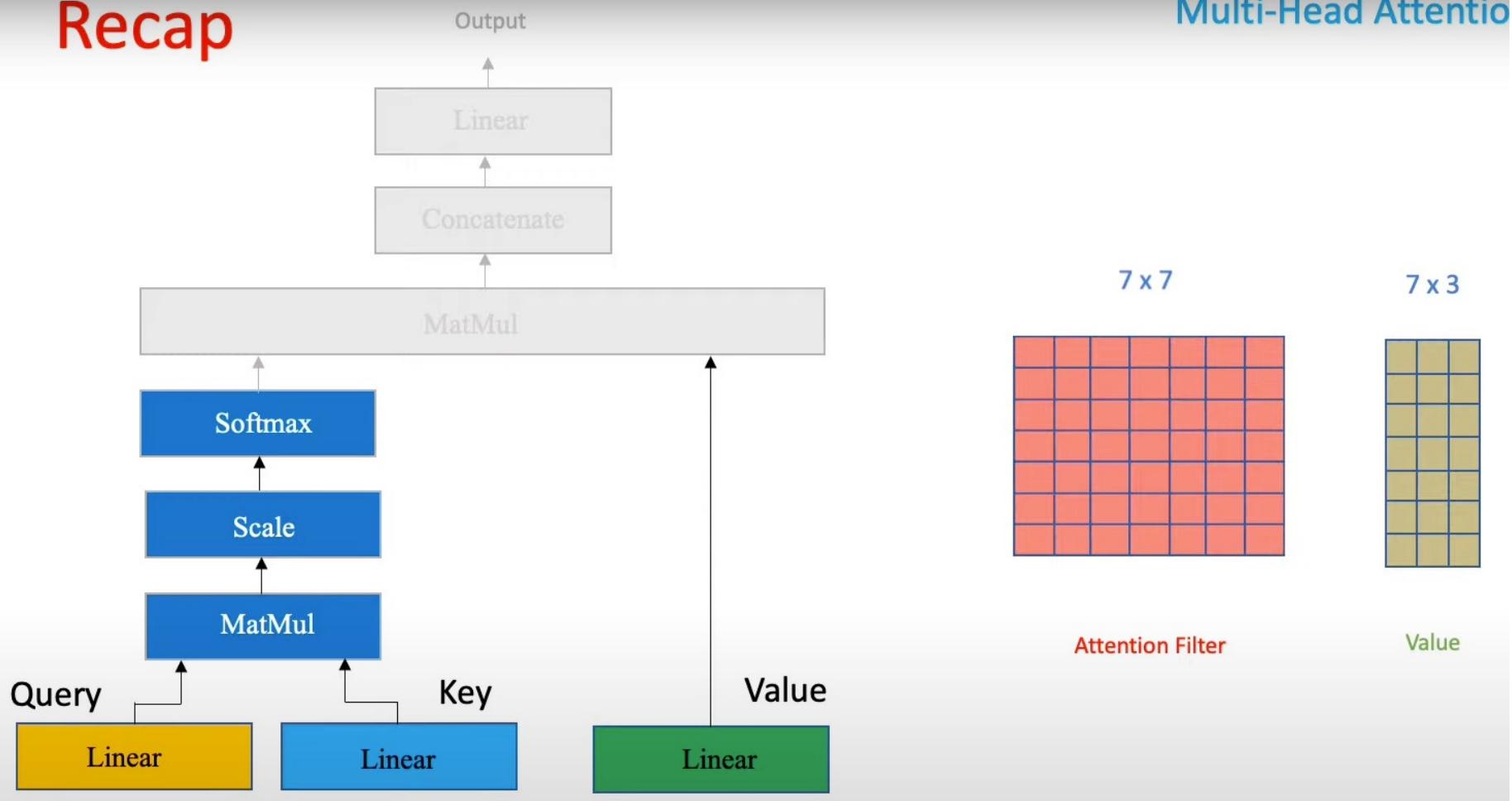
Attention Filter

Recap

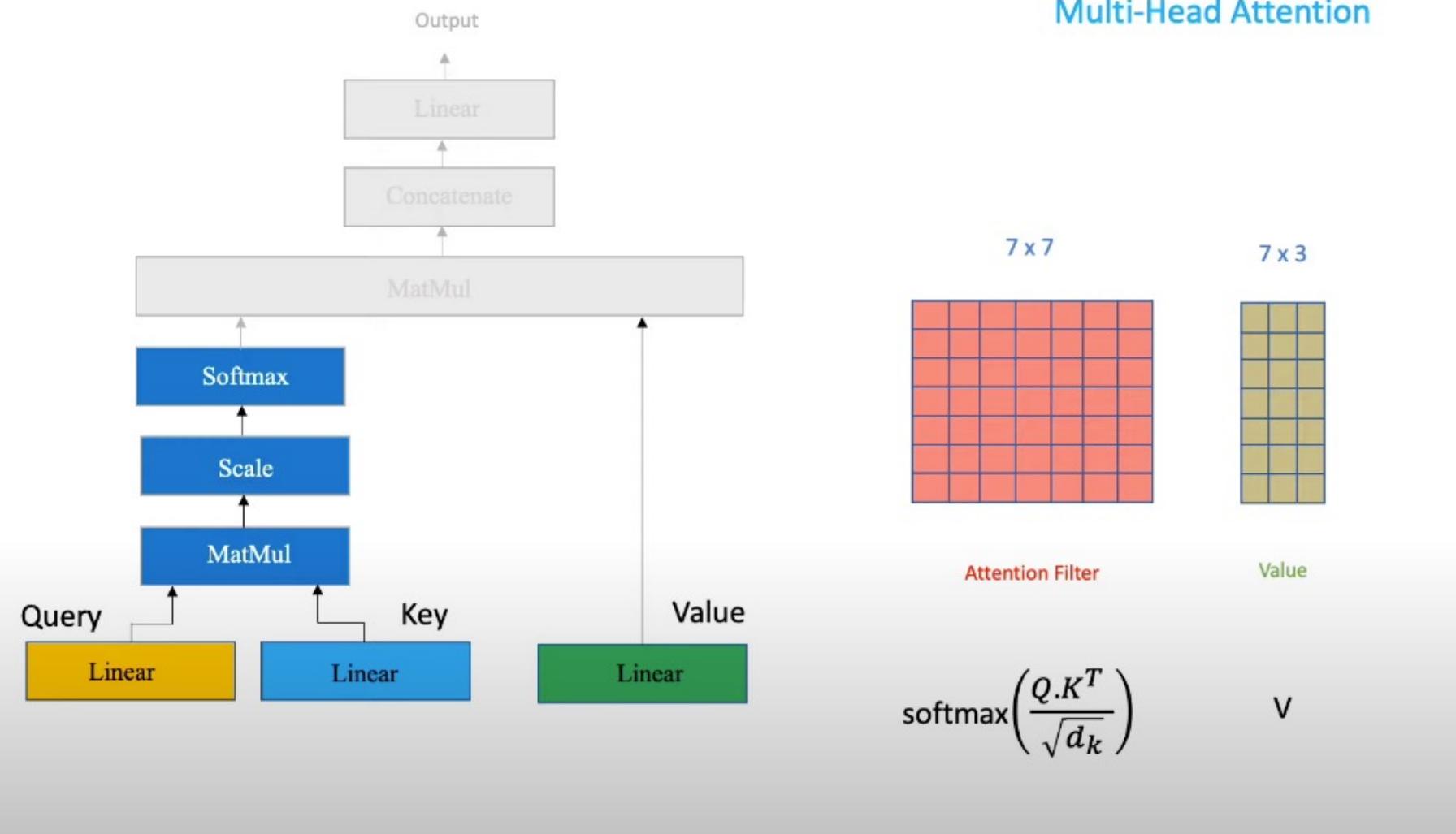
Recap

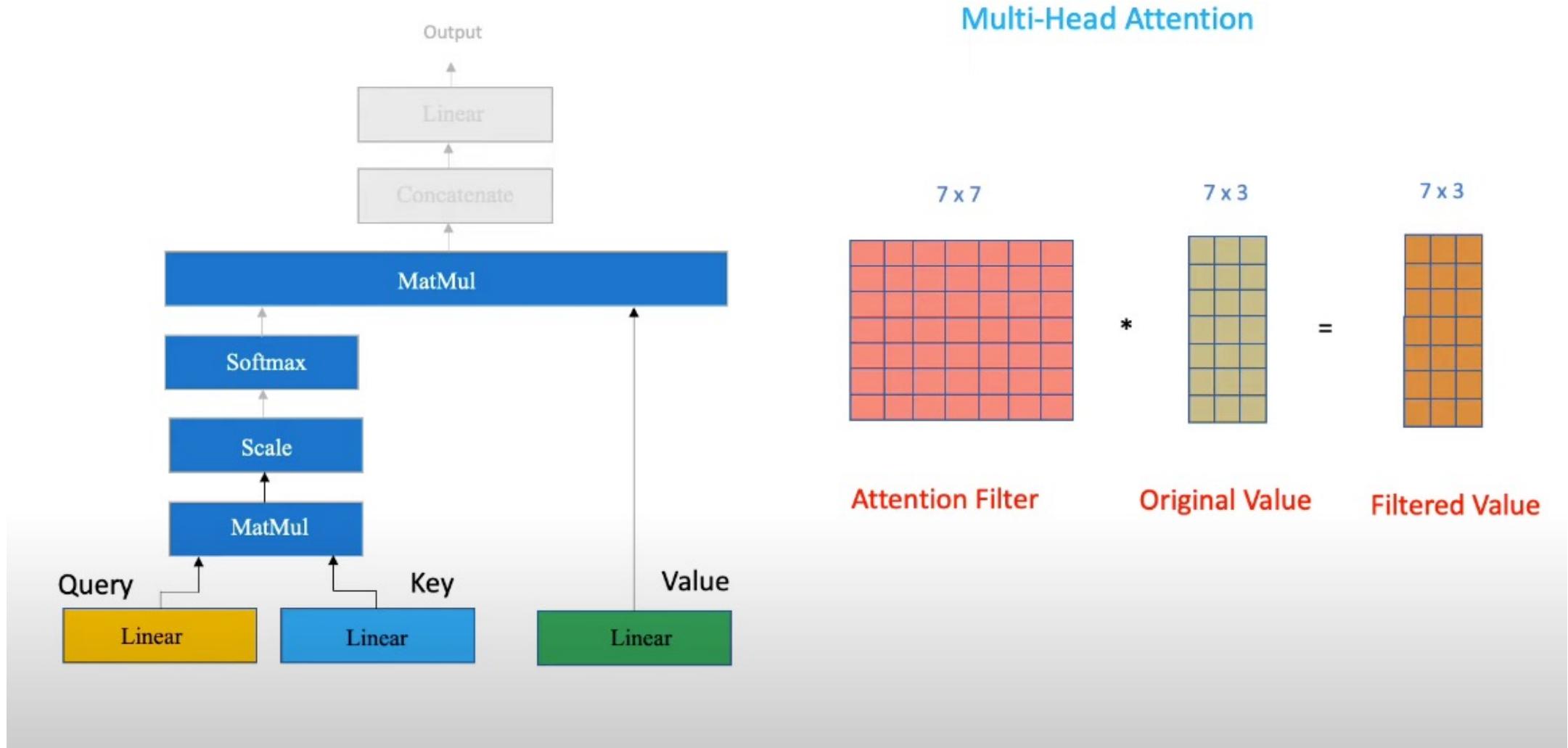


Recap

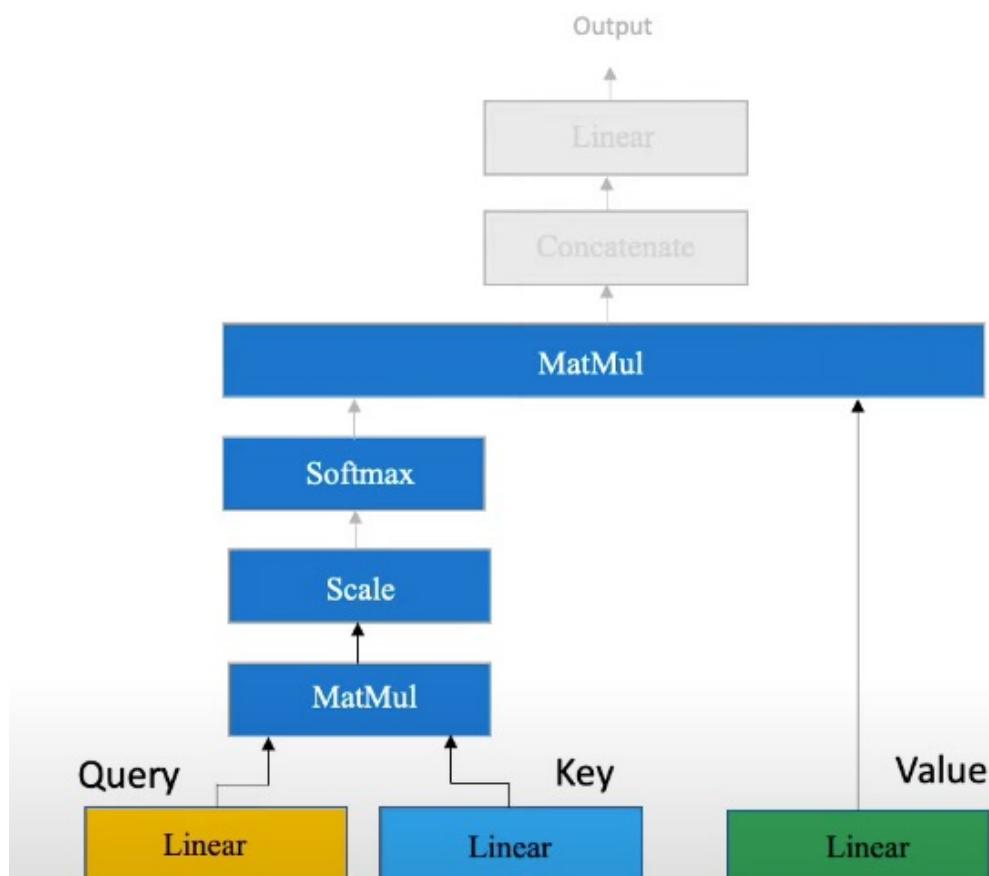


Multi-Head Attention





Multi-Head Attention



7 x 3

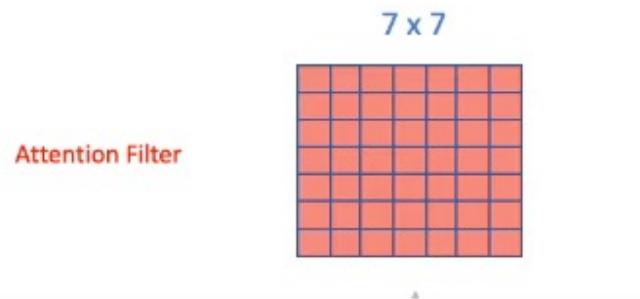


Filtered Value

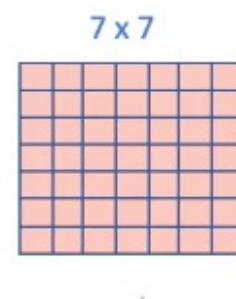
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_k}}\right) \cdot V$$

Multi-Head Attention

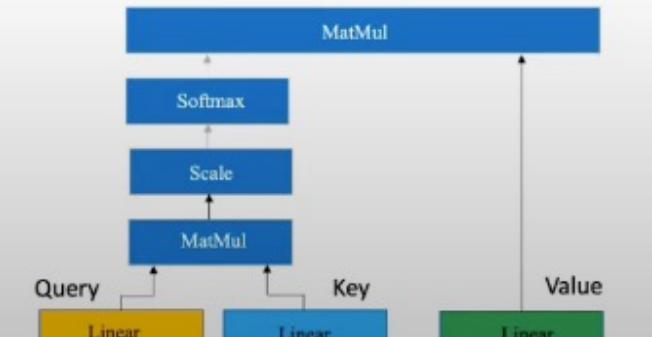
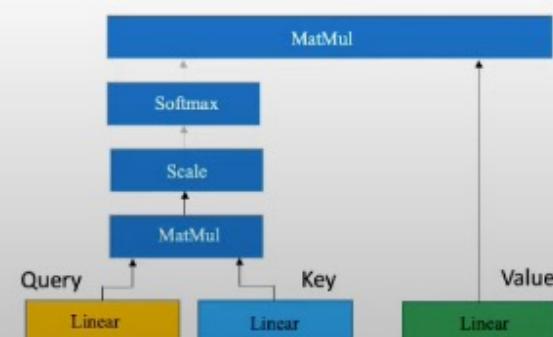
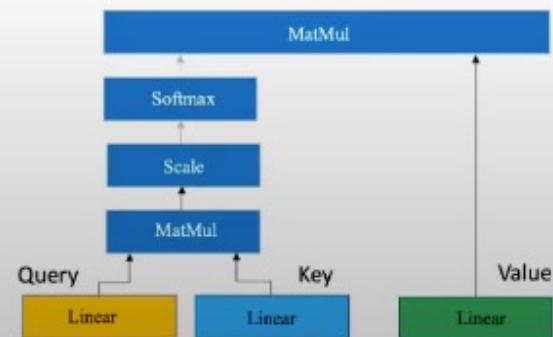
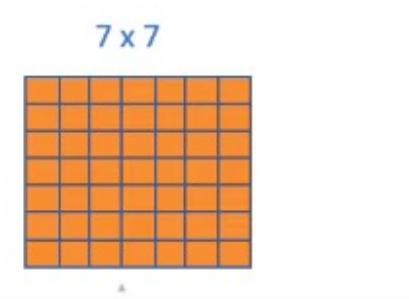
Attention Head 1

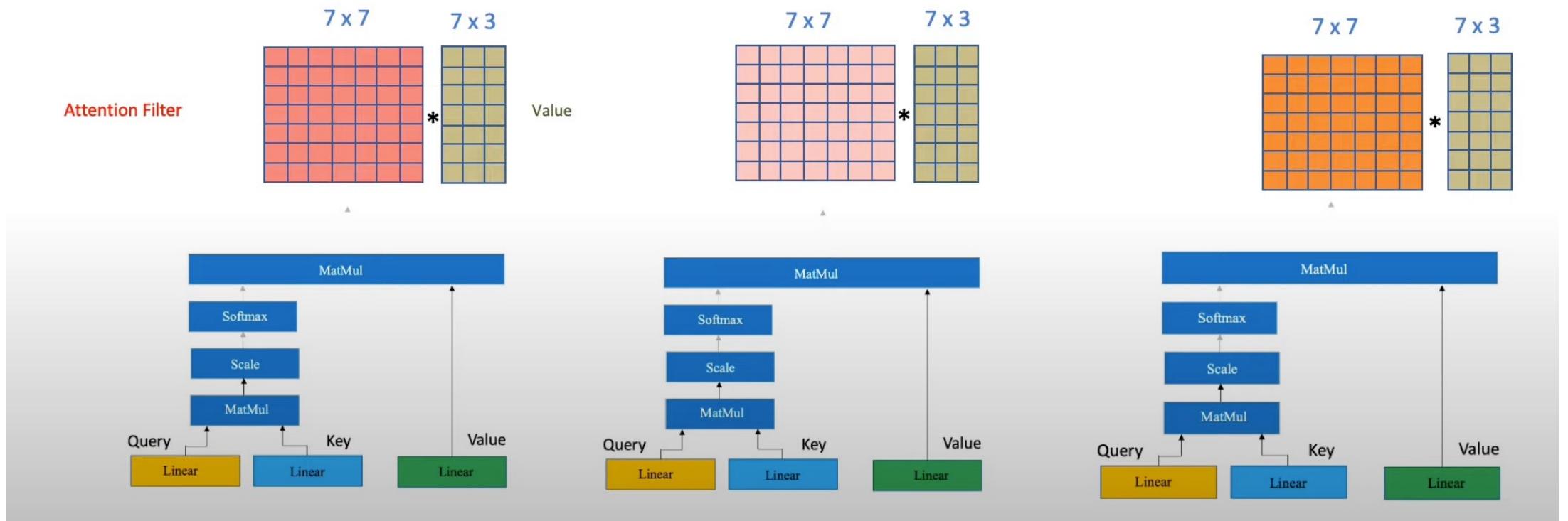


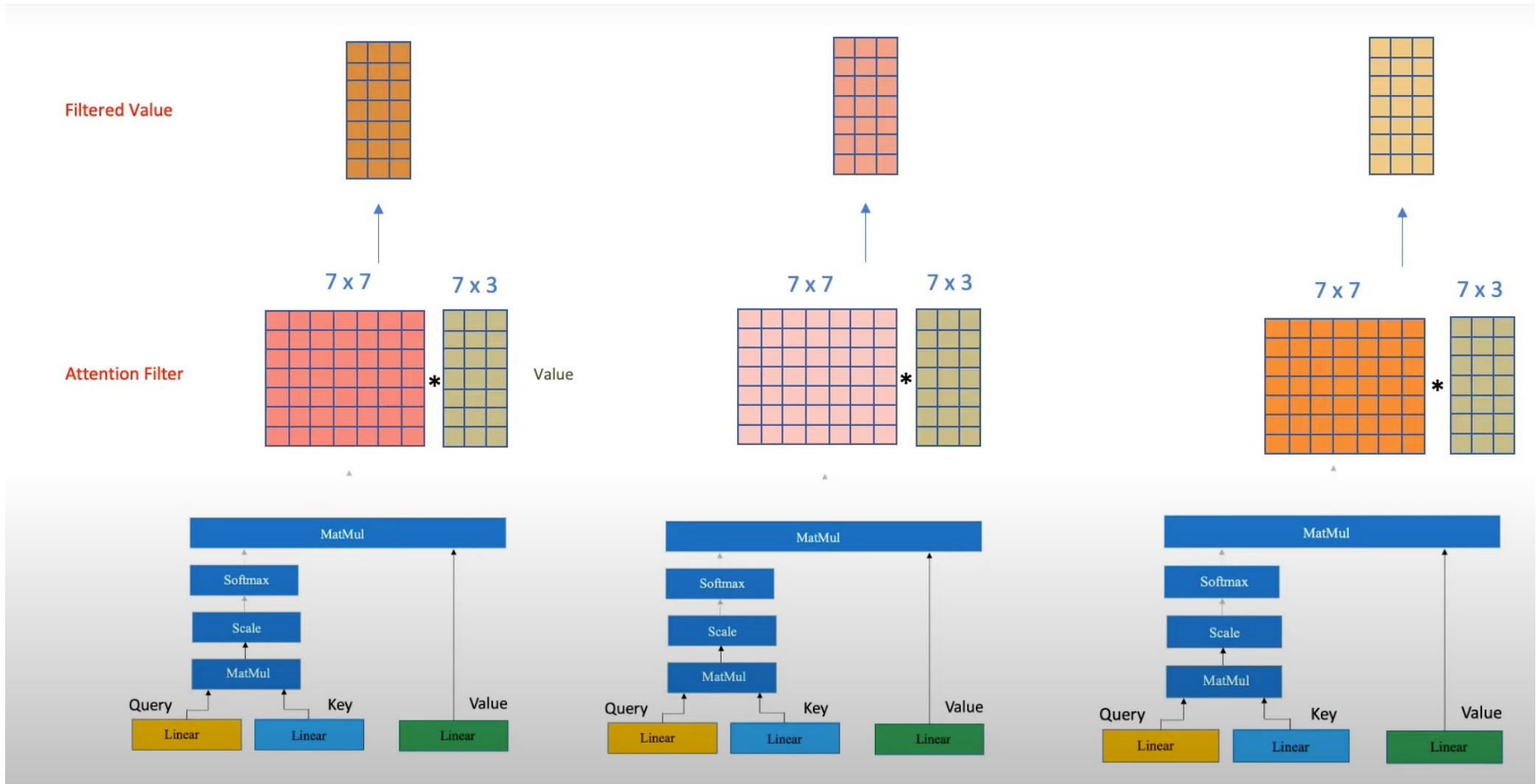
Attention Head 2

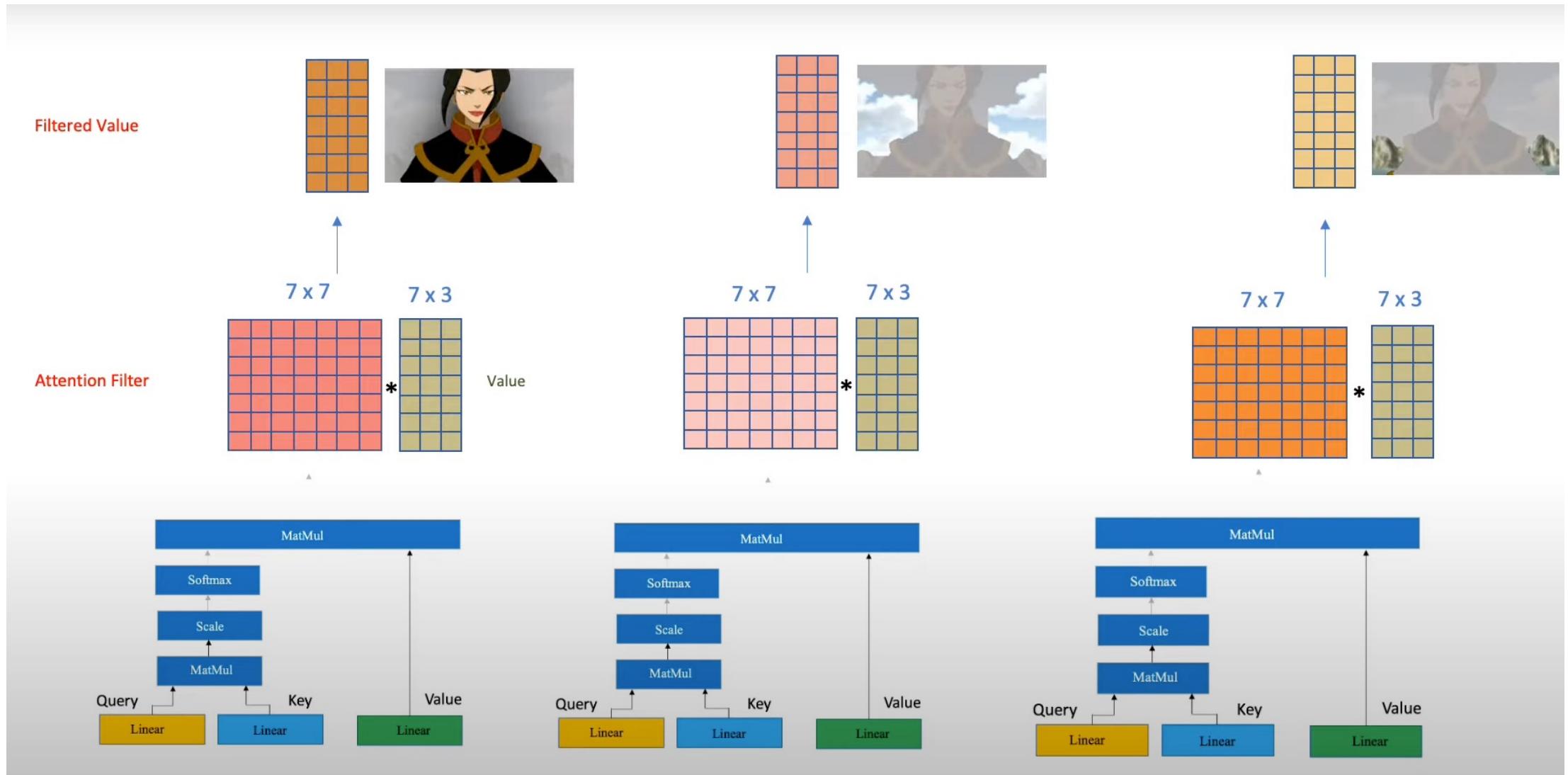


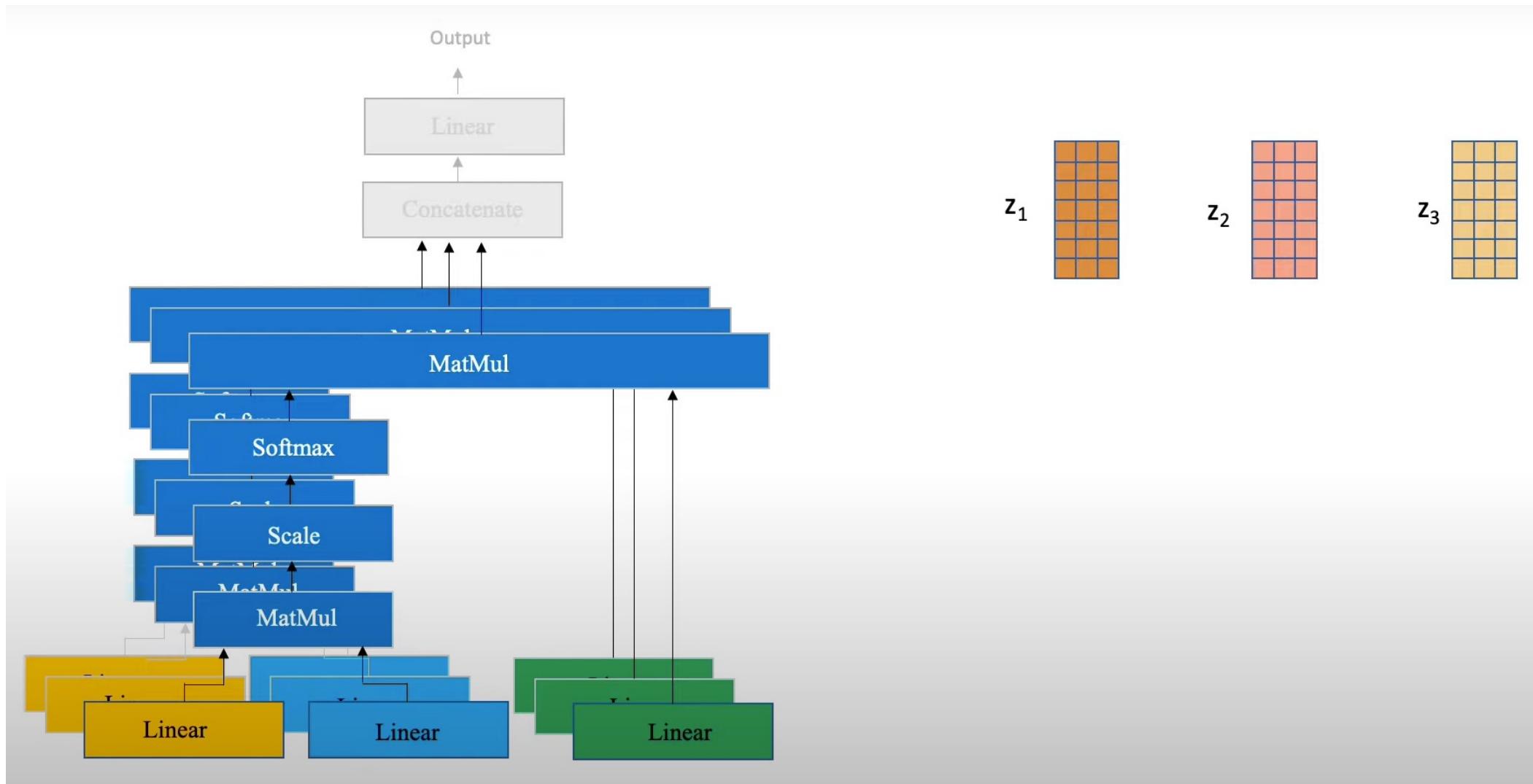
Attention Head 3

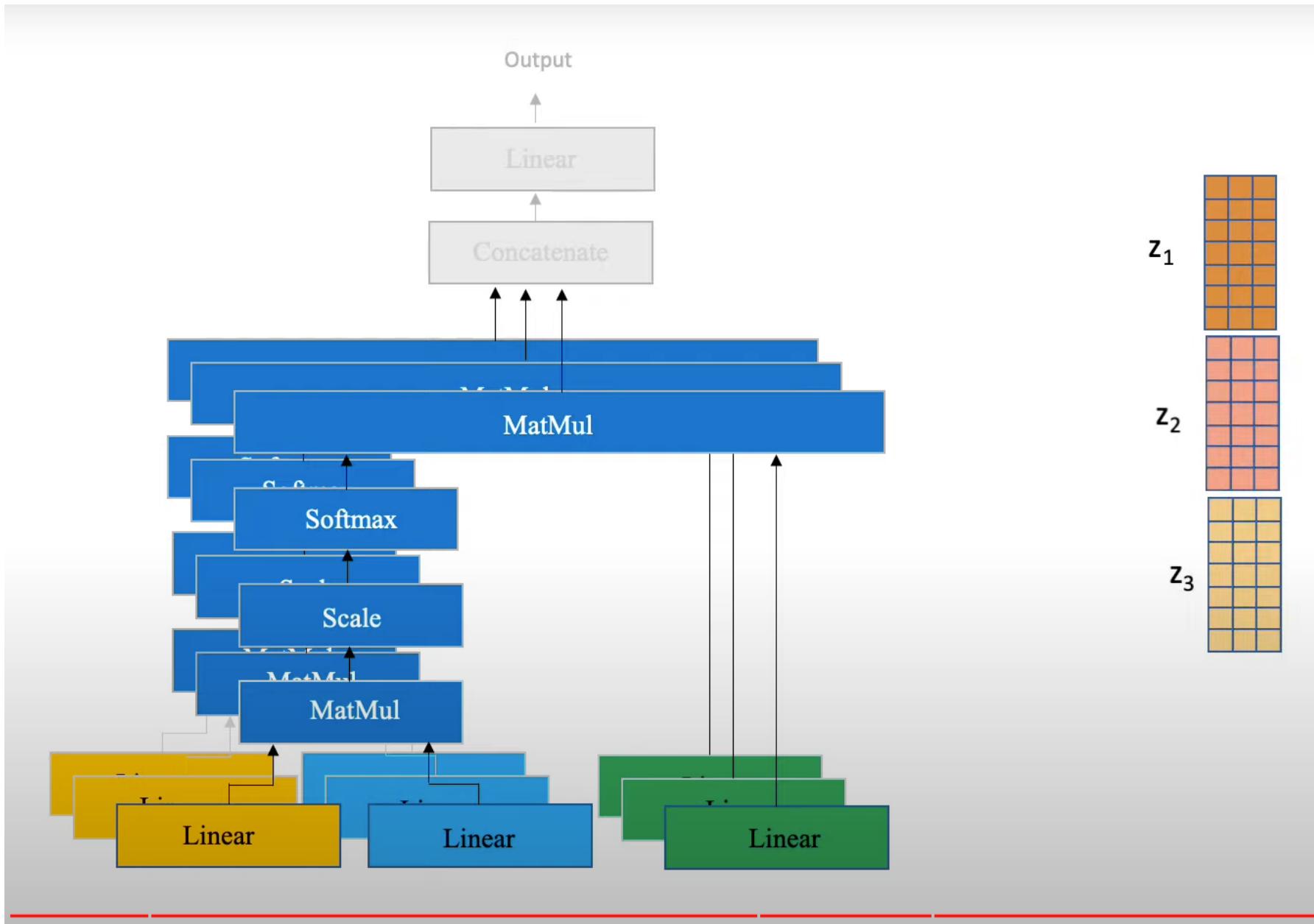


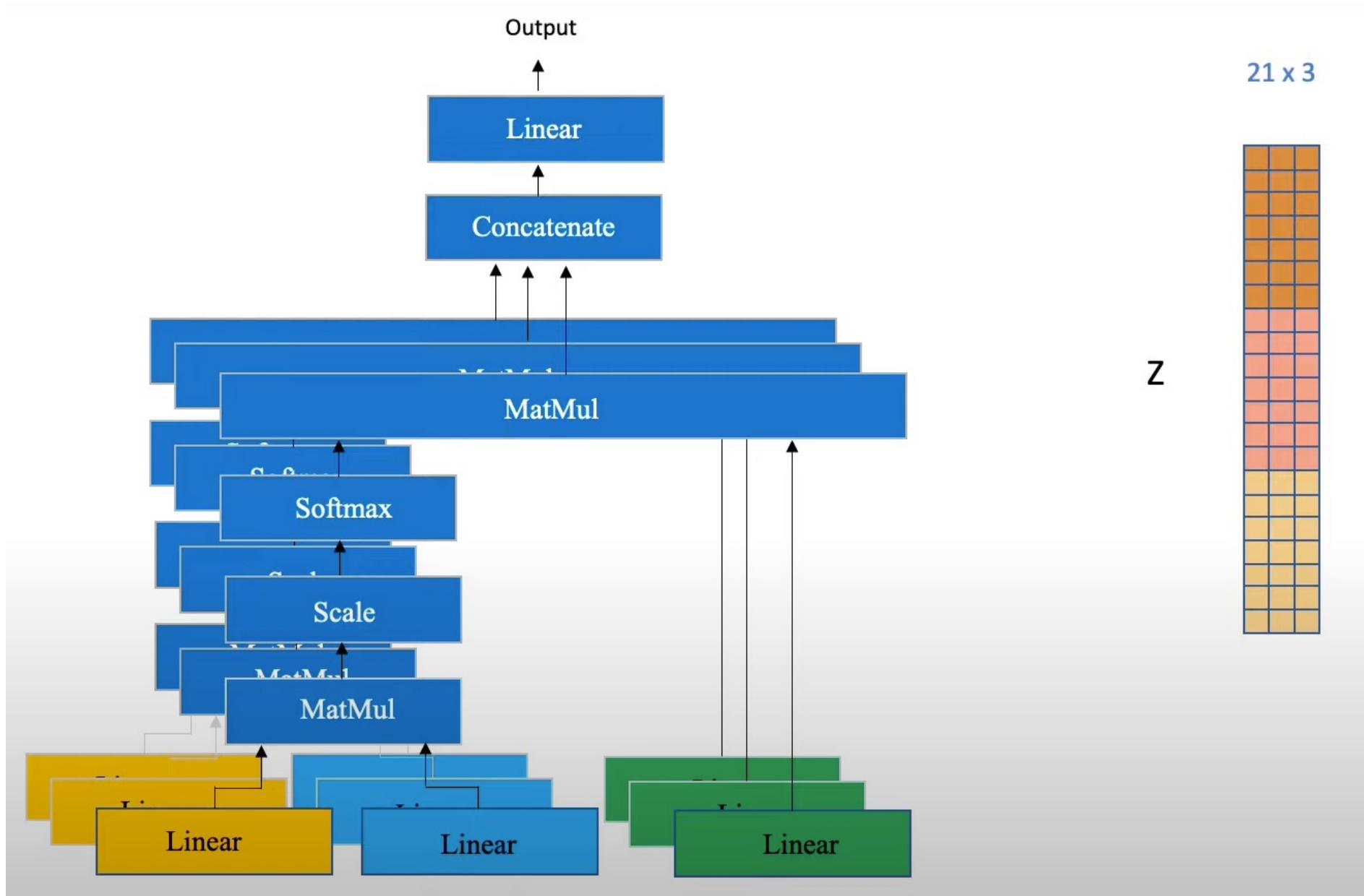


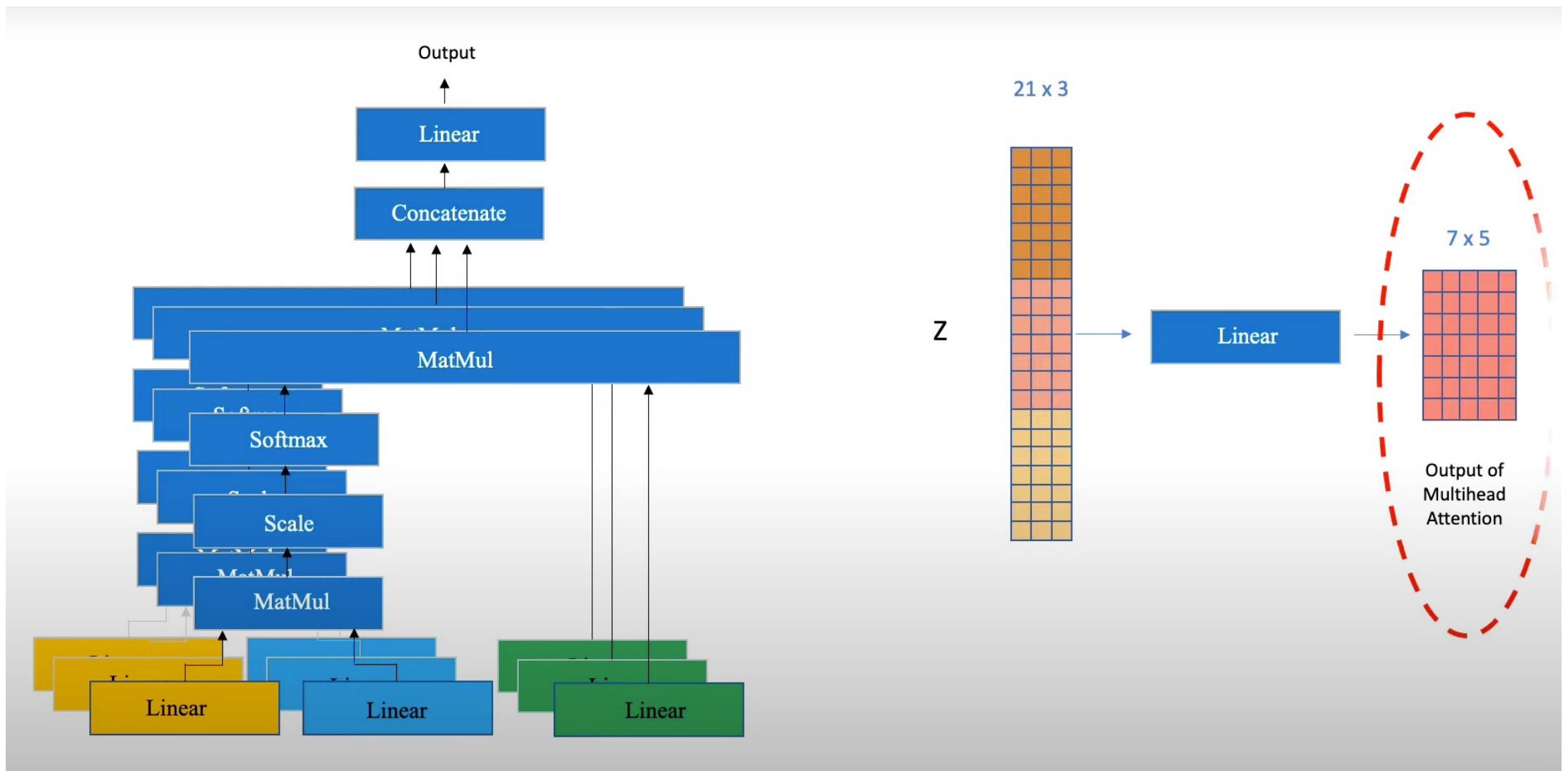












Reinventing Deep Learning Operation Via Einops

<https://einops.rocks/1-einops-basics/>

Reinventing Deep Learning Operation Via Einops

Flatten

```
img = image.imread("aki_dog.jpg")
```



Input tensor

`rearrange(img, "h w c -> (h w c)")`

Output shape:
Fuse height, width,
channel by multiplying
them together

Input shape:
height, width, channel

Load a batch of images to play with

```
In [2]: ims = numpy.load('./resources/test_images.npy', allow_pickle=False)
# There are 6 images of shape 96x96 with 3 color channels packed into tensor
print(ims.shape, ims.dtype)

(6, 96, 96, 3) float64
```

```
In [3]: # display the first image (whole 4d tensor can't be rendered)
ims[0]
```

Out[3]:



```
In [4]: # second image in a batch
ims[1]
```

Out[4]:



```
In [5]: # we'll use three operations
from einops import rearrange, reduce, repeat
```

```
In [6]: # rearrange, as its name suggests, rearranges elements
# below we swapped height and width.
# In other words, transposed first two axes (dimensions)
rearrange(ims[0], 'h w c -> w h c')
```

Out[6]:



In [7]:

```
# einops allows seamlessly composing batch and height to a new height dimension
# We just rendered all images by collapsing to 3d tensor!
rearrange(ims, 'b h w c -> (b h) w c')
```

Out[7]:

e
i
n
o
p
s

Decomposition of axis

```
In [11]: # decomposition is the inverse process - represent an axis as a combination of new axes
# several decompositions possible, so b1=2 is to decompose 6 to b1=2 and b2=3
rearrange(ims, '(b1 b2) h w c -> b1 b2 h w c ', b1=2).shape
```

```
Out[11]: (2, 3, 96, 96, 3)
```

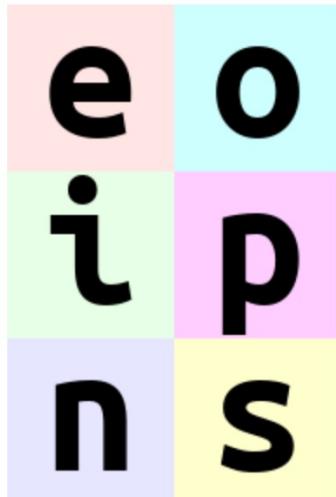
```
In [12]: # finally, combine composition and decomposition:
rearrange(ims, '(b1 b2) h w c -> (b1 h) (b2 w) c ', b1=2)
```

```
Out[12]:
```



```
In [13]: # slightly different composition: b1 is merged with width, b2 with height
# ... so letters are ordered by w then by h
rearrange(ims, '(b1 b2) h w c -> (b2 h) (b1 w) c ', b1=2)
```

```
Out[13]:
```



Meet einops.reduce

In einops-land you don't need to guess what happened

```
x.mean(-1)
```

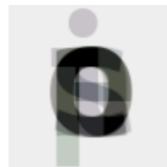
Because you write what the operation does

```
reduce(x, 'b h w c -> b h w', 'mean')
```

if axis is not present in the output – you guessed it – axis was reduced.

```
In [18]: # average over batch  
reduce(ims, 'b h w c -> h w c', 'mean')
```

Out[18]:



```
In [19]: # the previous is identical to familiar:  
ims.mean(axis=0)  
# but is so much more readable
```

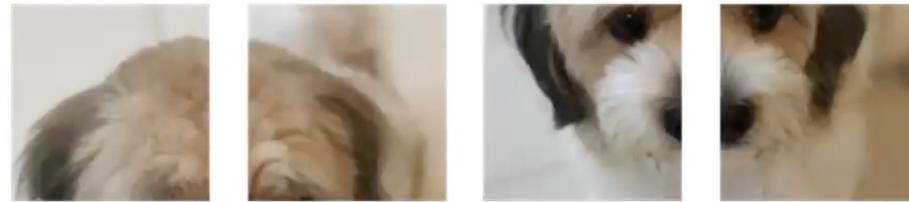
Out[19]:



```
In [20]: # Example of reducing of several axes  
# besides mean, there are also min, max, sum, prod  
reduce(ims, 'b h w c -> h w', 'min')
```

Out[20]:





`rearrange(img, "(p1 h) (p2 w) c -> (p1 p2) h w c", p1=2, p2=2)`

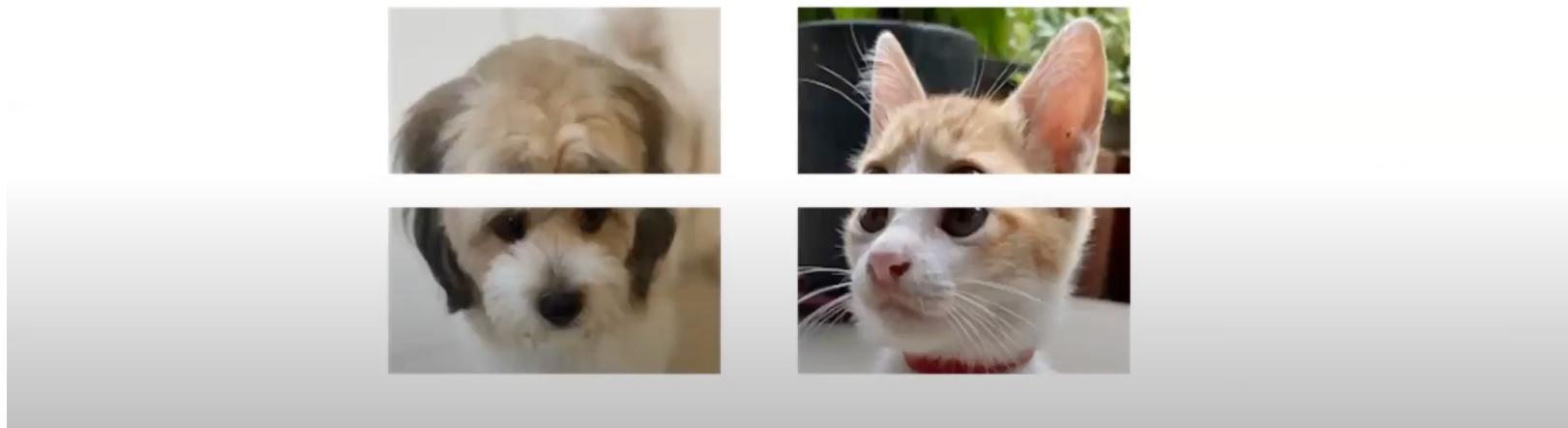
Mixing 2 images



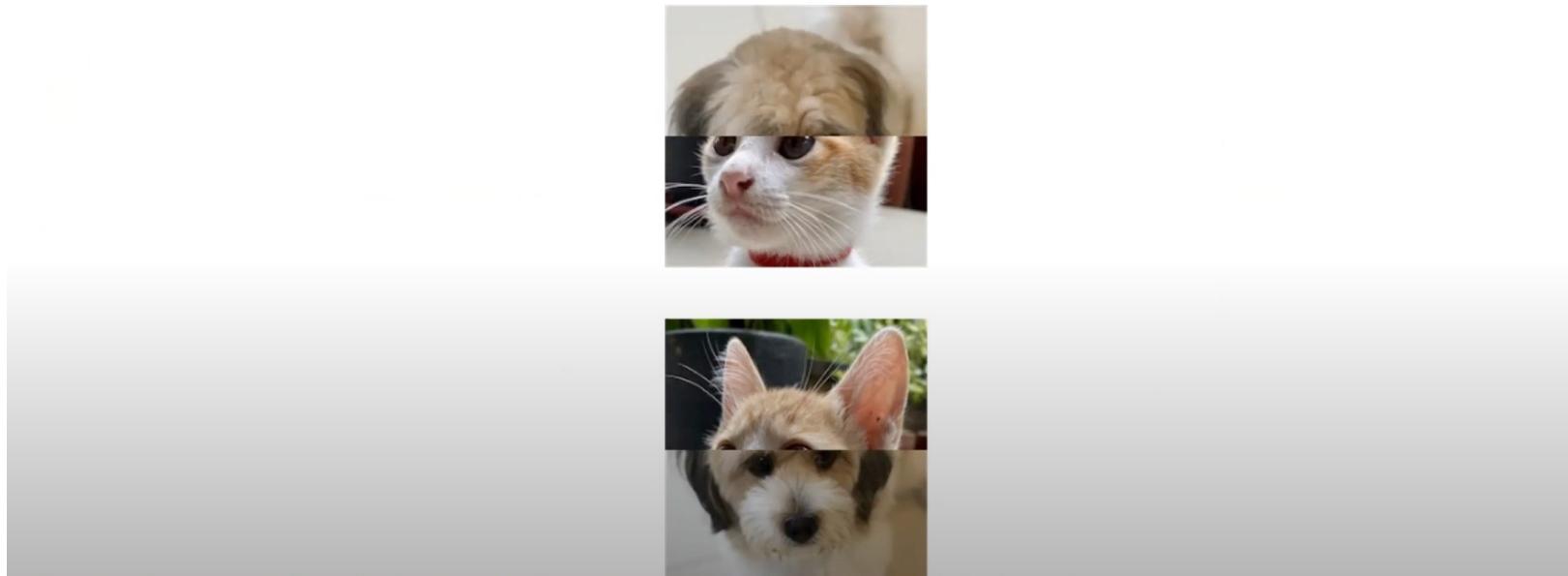
```
img1 = image.imread("aki_dog.jpg")
img2 = image.imread("wonder_cat.jpg")
imgs = np.array([img1, img2])
```



```
imgs = rearrange(imgs, "b (k h) w c -> k b h w c", k =2)
```

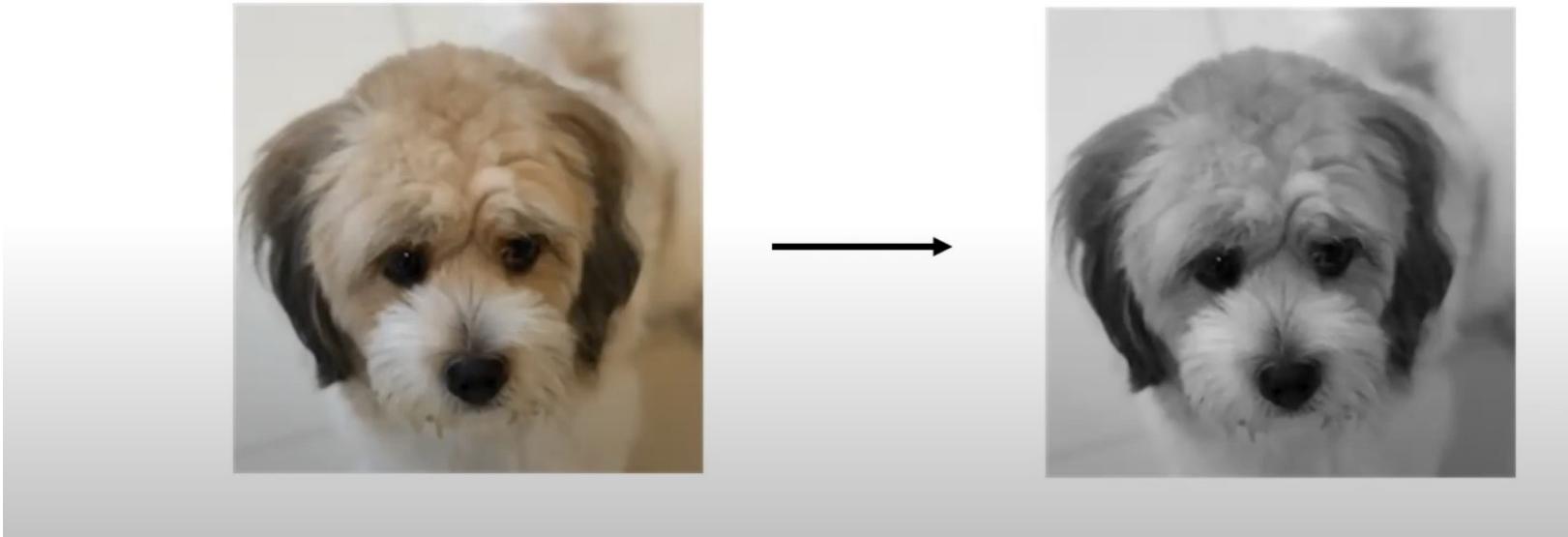


```
imgs = rearrange(imgs, "i j h w c -> j (i h) w c")
```



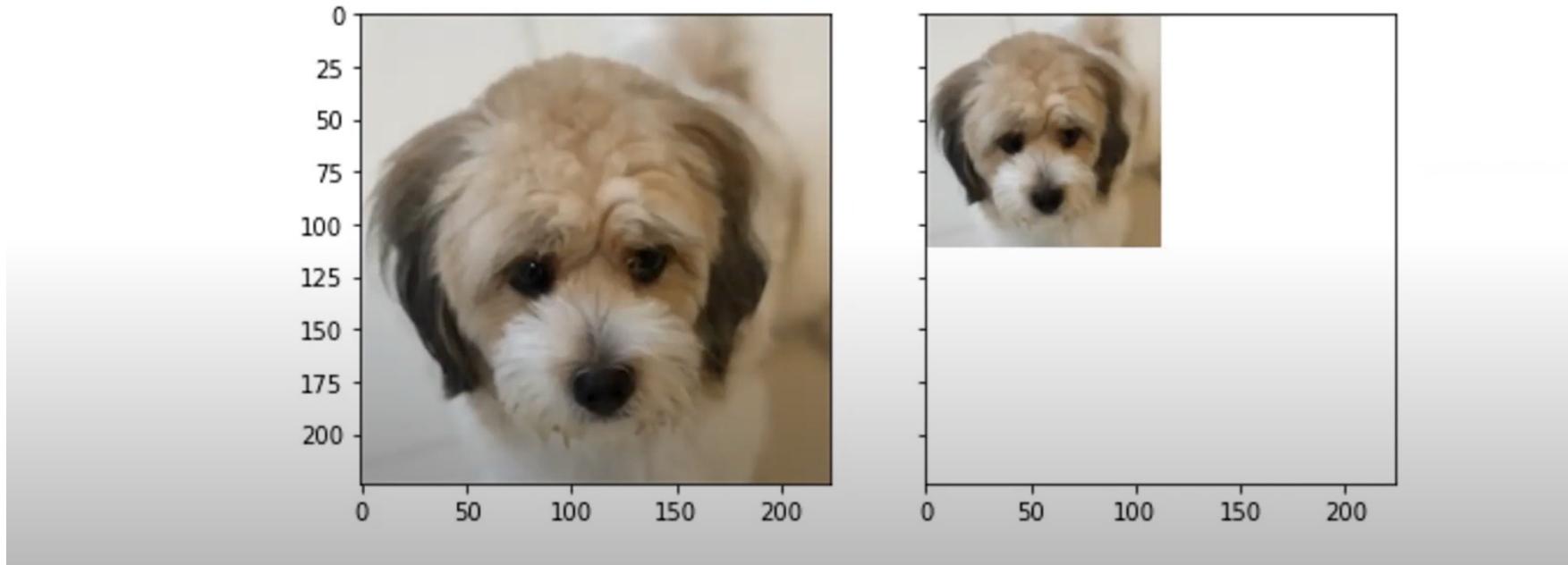
RGB to Grayscale

```
img = reduce(img, "h w c -> h w", 'mean')
```



Downsize

```
reduce(img, "(h 2) (w 2) c -> h w c", 'mean')
```

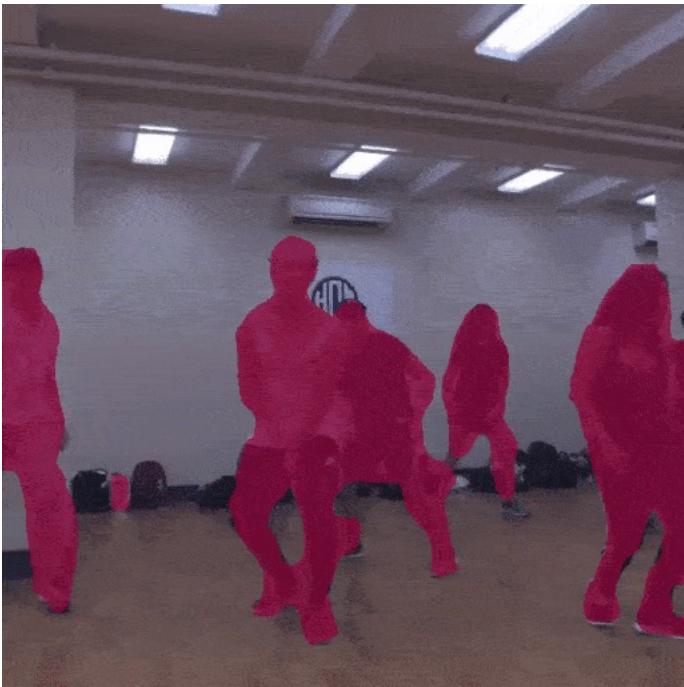


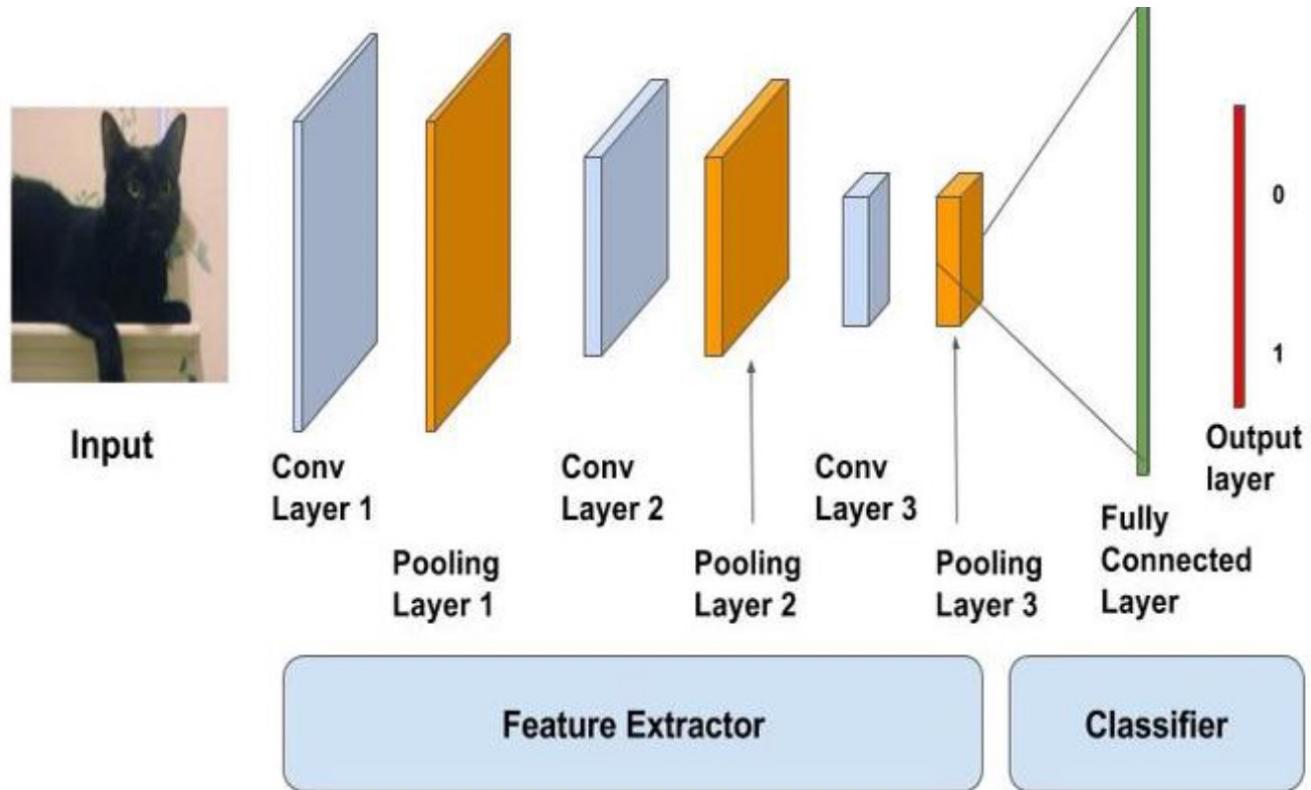
Implementing **ViT** in PyTorch

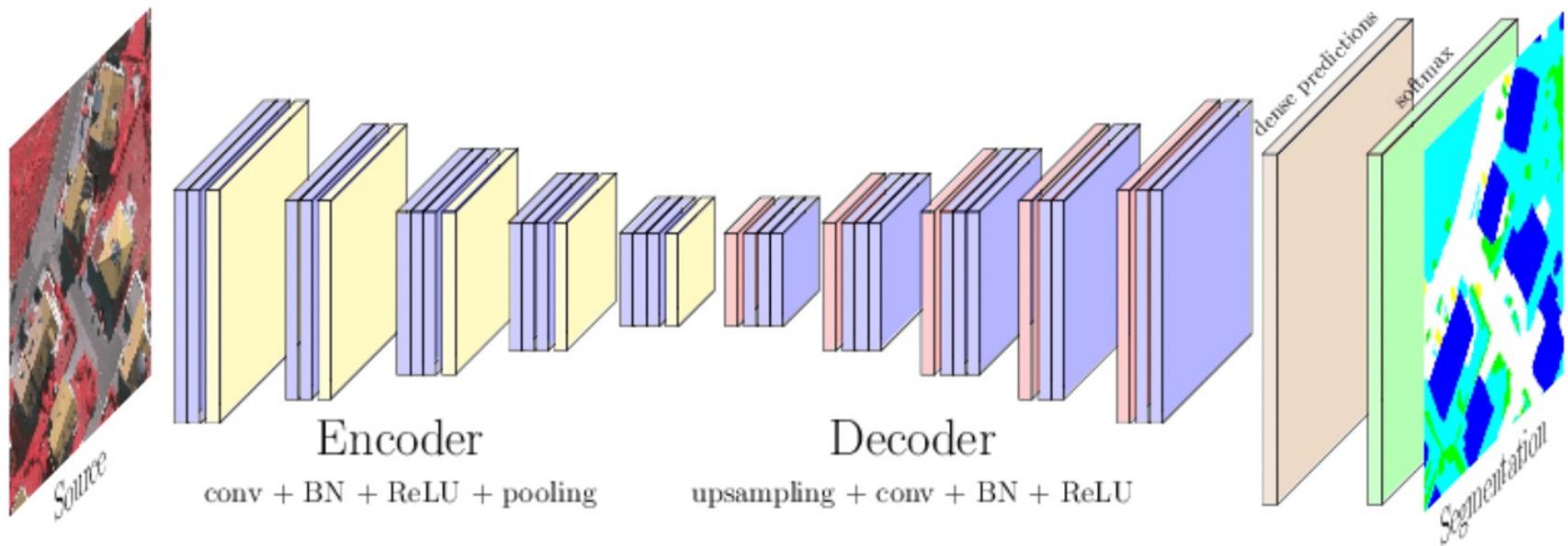


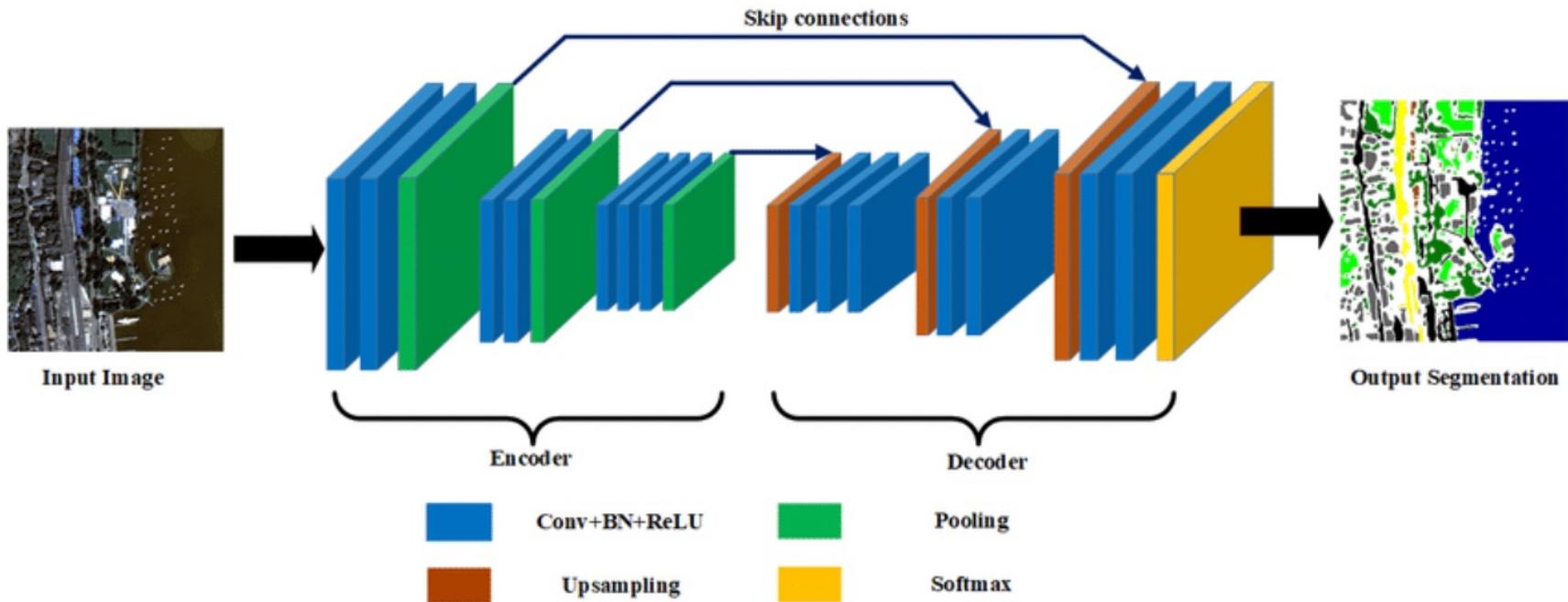
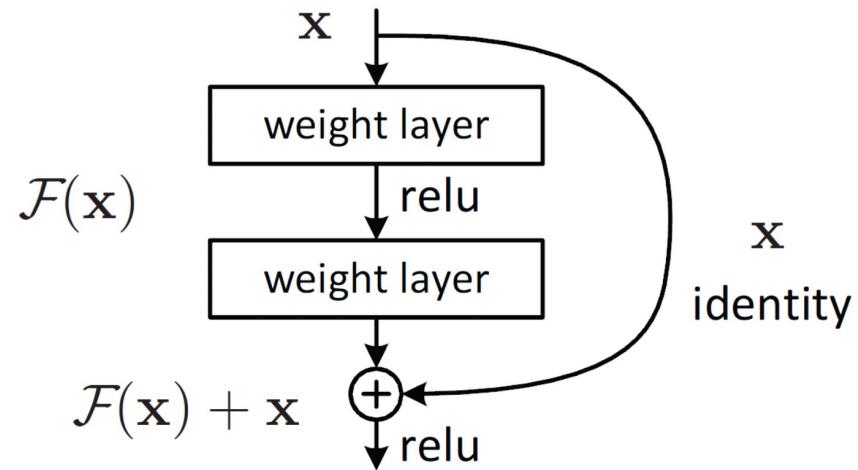
<https://medium.com/towards-data-science/implementing-visualtransformer-in-pytorch-184f9f16f632>

image segmentation



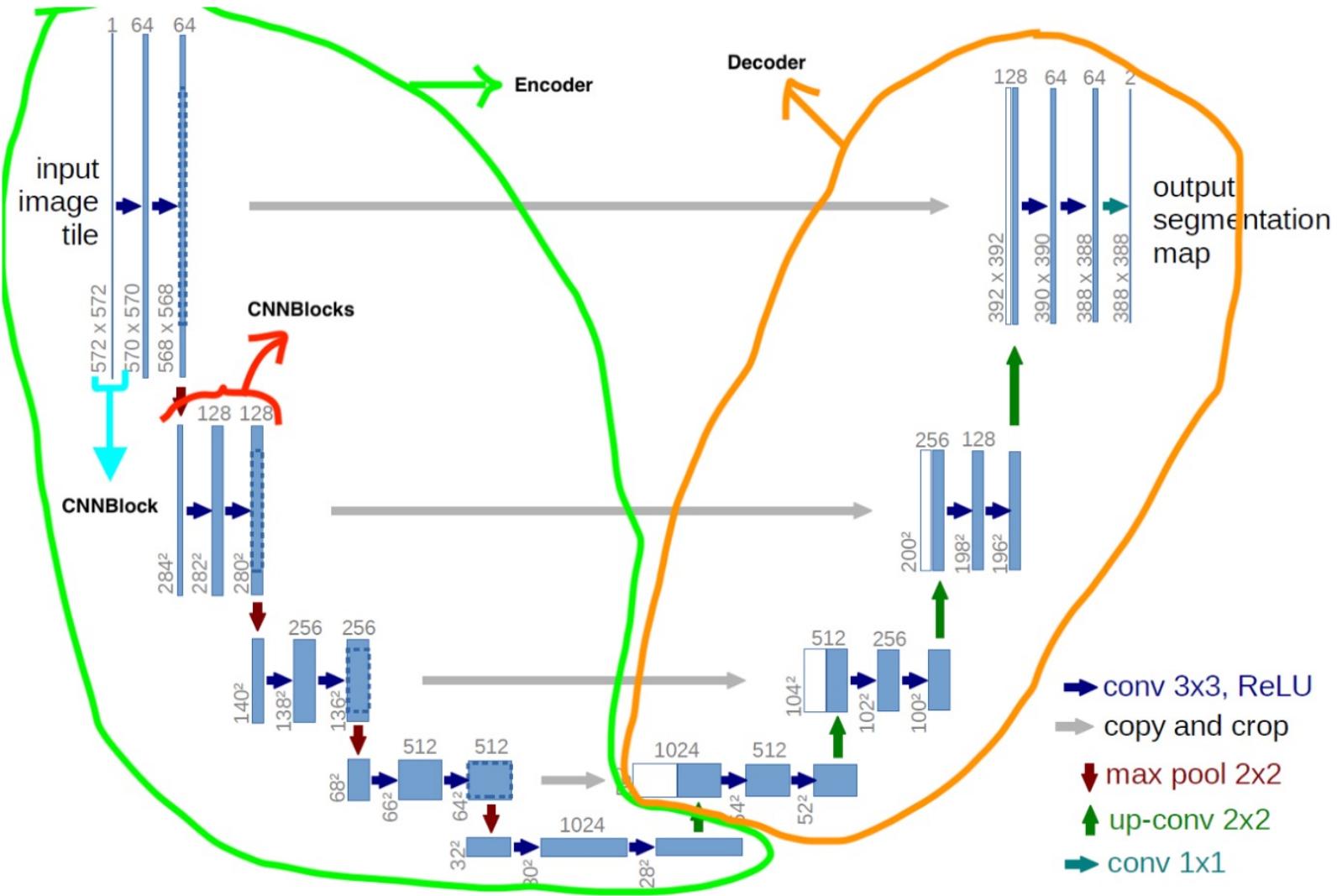




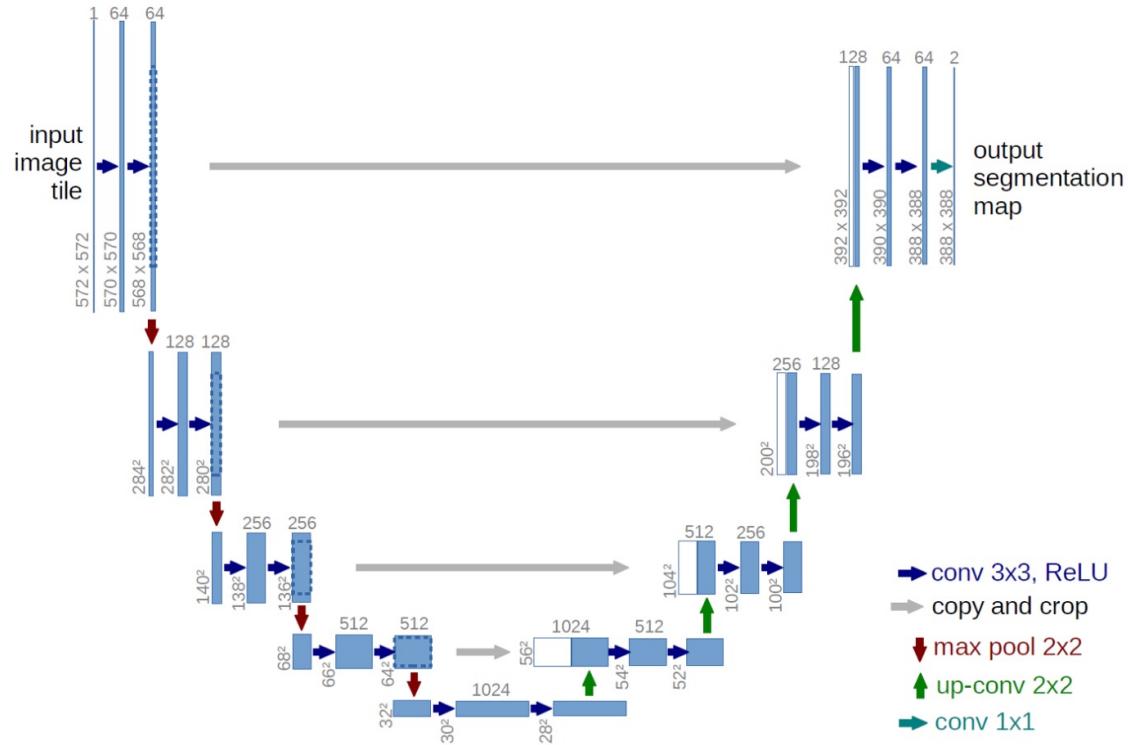


<https://medium.com/mlearning-ai/semantic-segmentation-with-pytorch-u-net-from-scratch-502d6565910a>

https://medium.com/@Chinmay_Paranjape/satellite-imagery-segmentation-using-u-net-4ec7f265ddbe



U-net Architecture



- ➡ conv 3x3, ReLU
- ➡ copy and crop
- ⬇ max pool 2x2
- ⬆ up-conv 2x2
- ➡ conv 1x1

```

1 class CNNBlock(nn.Module):
2     def __init__(self,
3                  in_channels,
4                  out_channels,
5                  kernel_size=3,
6                  stride=1,
7                  padding=0):
8         super(CNNBlock, self).__init__()
9
10        self.seq_block = nn.Sequential(
11            nn.Conv2d(in_channels, out_channels, kernel_size, stride, padding, bias=False),
12            nn.BatchNorm2d(out_channels),
13            nn.ReLU(inplace=True)
14        )
15
16    def forward(self, x):
17        x = self.seq_block(x)
18        return x

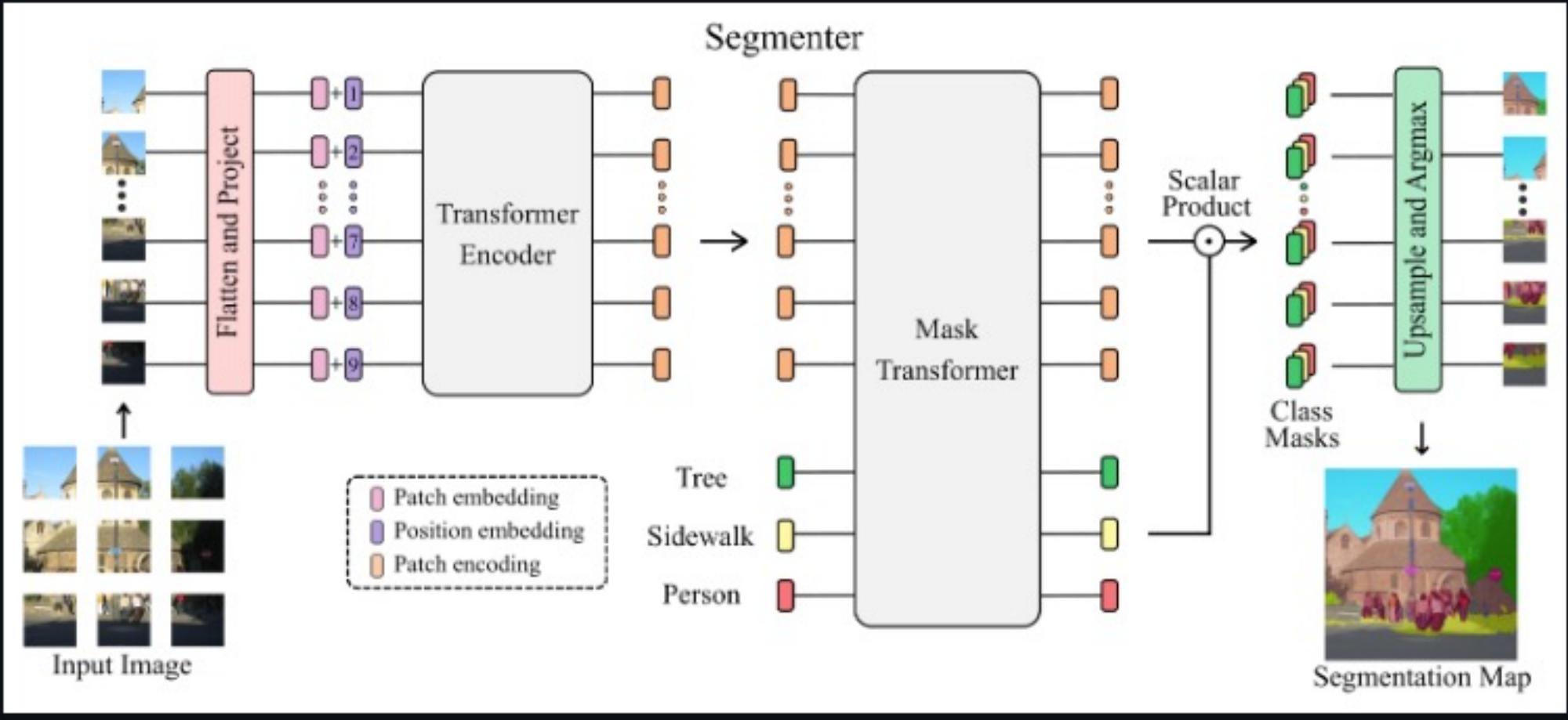
```

```

1 class CNNBlocks(nn.Module):
2     """
3     Parameters:
4     n_conv (int): creates a block of n_conv convolutions
5     in_channels (int): number of in_channels of the first block's convolution
6     out_channels (int): number of out_channels of the first block's convolution
7     expand (bool) : if True after the first convolution of a block the number of channels
8     """
9
10    def __init__(self,
11                 n_conv,
12                 in_channels,
13                 out_channels,
14                 padding):
15        super(CNNBlocks, self).__init__()
16        self.layers = nn.ModuleList()
17        for i in range(n_conv):
18
19            self.layers.append(CNNBlock(in_channels, out_channels, padding=padding))
20            # after each convolution we set (next) in_channel to (previous) out_channels
21            in_channels = out_channels
22
23    def forward(self, x):
24        for layer in self.layers:
25            x = layer(x)
26        return x

```

Segmenter: Transformer for Semantic Segmentation



crack Image Dataset

[Try Pre-Trained Model](#)

VERSIONS

2022-09-29 2:14pm v2 Sep 30, 2022

2022-09-02 11:09pm v1 Sep 2, 2022

2022-09-29 2:14pm
Version 2 Generated Sep 30, 2022

[Download](#)

IMAGES

4029 images [View All Images >](#)

TRAIN / TEST SPLIT

Training Set 3.7K images	Validation Set 200 images	Testing Set 112 images
---------------------------------	----------------------------------	-------------------------------

PREPROCESSING **Resize:** Stretch to 416x416

<https://universe.roboflow.com/chonnam-national-university/crack-bphdr/dataset/2>

