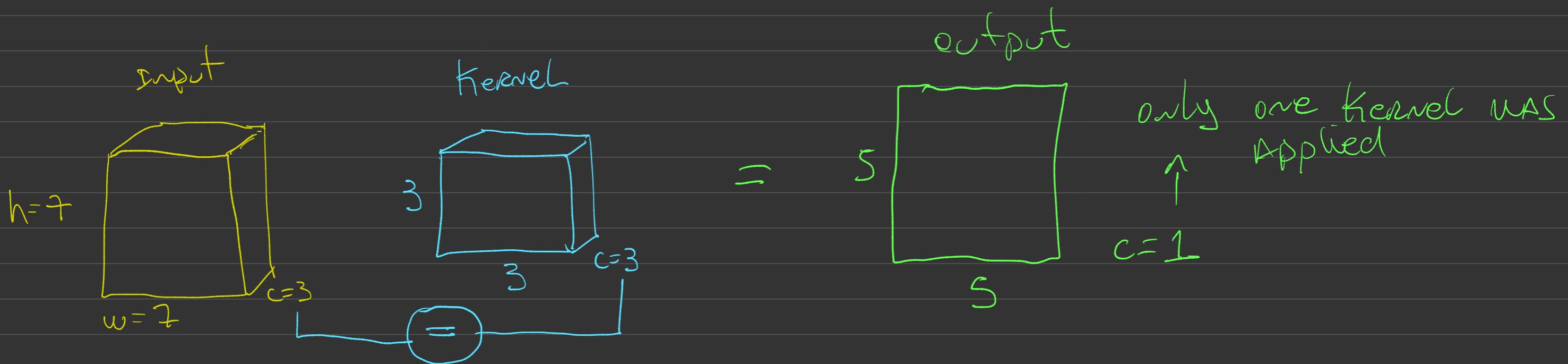
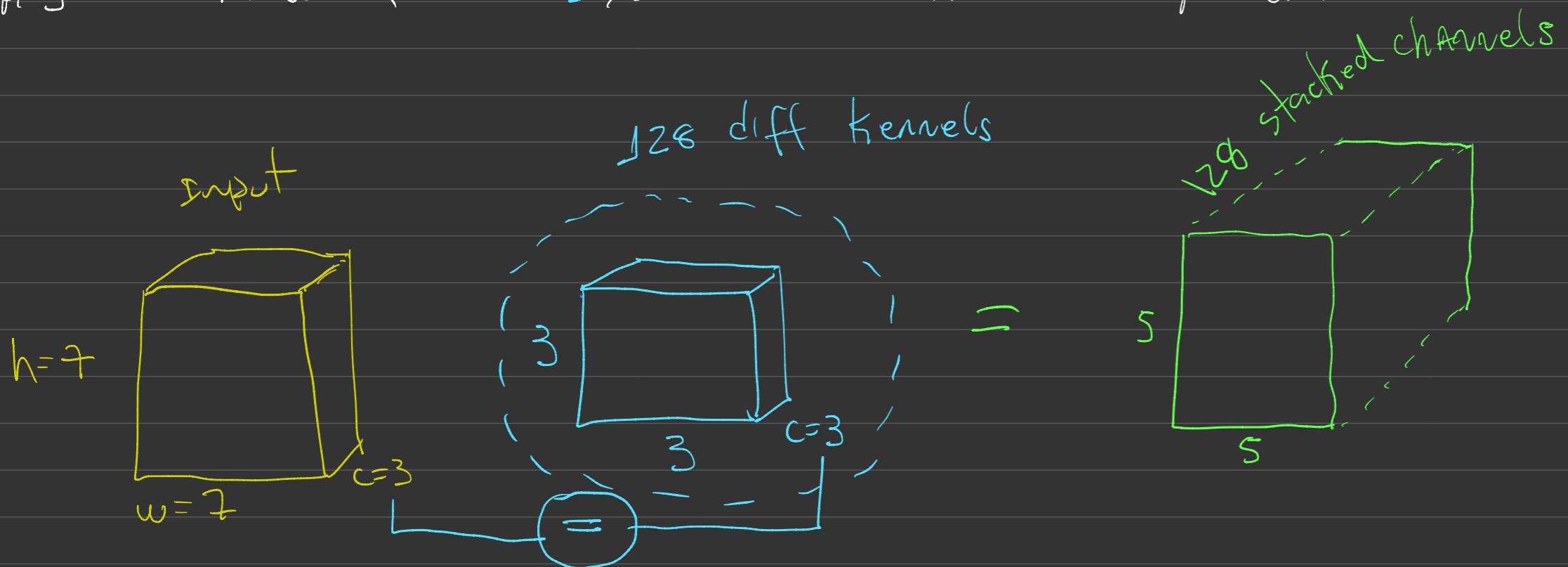


# Standard 2D Convolutions (1)



To generate an output of 128 channels we must apply 128 kernels of  $3 \times 3 \times 3$ , so we will have 128 output channels.



# Standard 2D Convolutions (2)

number of operations

$$\begin{aligned} C_{in} &= 3 \\ H_{in} &= 7 \\ W_{in} &= 7 \end{aligned}$$

$$\begin{aligned} C_{out} &= 128 \\ H_{out} &= 5 \\ W_{out} &= 5 \end{aligned}$$

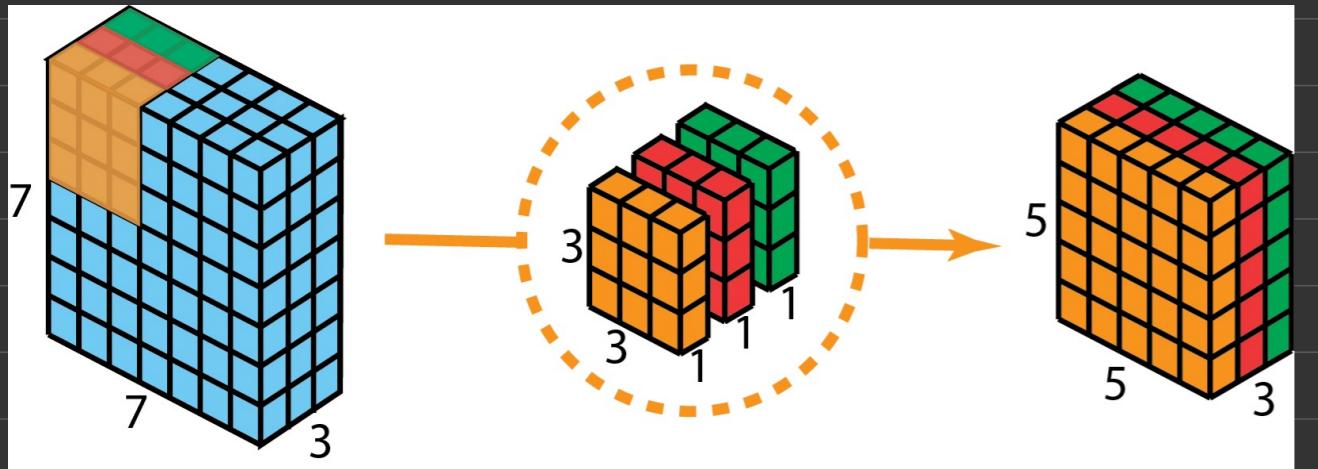
Filter

$$\begin{aligned} C_{in} &= 3 \\ H &= 3 \\ W &= 3 \end{aligned} \quad ] * C_{out}$$

There are 128  $3 \times 3 \times 3$  kernels that makes  $5 \times 5$  times  $C_{out}$

$$= 128 \times 3 \times 3 \times 3 \times 5 \times 5 = 86400 \text{ multiplications}$$

# Depthwise Separable Convolutions (1)



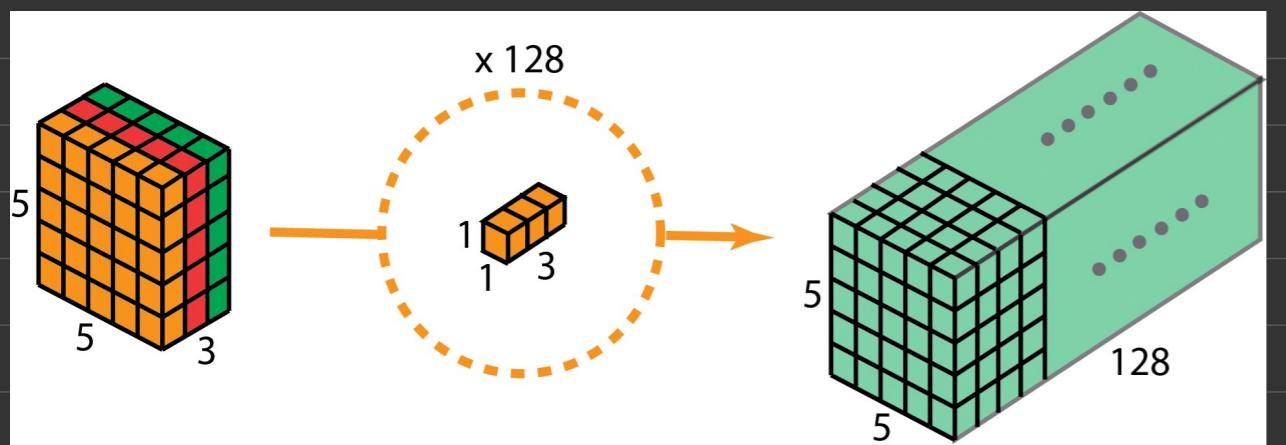
1. Split the Input image in channels  $I_1 I_2 I_3$
2. Split the Filter in channels  $k_1 k_2 k_3$
3. Apply each kernel in the respective  $I_n$  channel

$$I_1 \otimes k_1 = \text{Out}_1 [S, S, 1]$$

$$I_2 \otimes k_2 = \text{Out}_2 [S, S, 1]$$

$$I_3 \otimes k_3 = \text{Out}_3 [S, S, 1]$$

4. Stack the output. This produces an output  $[S, S, 3]$



5. Apply  $C_{out}$  filters of  $1 \times 1 \times C_{in}$   
Each output will be of  $S \times S \times 1$
6. Stack the 128 outputs will produce  $(S, S, 128)$

- Number of Operations

$$\begin{aligned} C_{in} &= 1 \\ H_{in} &= 7 \\ W_{in} &= 7 \end{aligned} \Bigg] \times 3$$

$$\begin{aligned} C_{out} &= 128 \\ H_{out} &= 5 \\ W_{out} &= 5 \end{aligned}$$

Filter\_1

$$\begin{aligned} C_{in} &= 1 \\ H &= 3 \\ W &= 3 \end{aligned} \Bigg] \times 3$$

Filter\_2

$$\begin{aligned} C_{in} &= 3 \\ H &= 1 \\ W &= 1 \end{aligned} \Bigg] \times 128$$

$$\begin{aligned} 1) \text{Filter\_1: } 3 \times 3 \times 3 \times 1 \times 5 \times 5 &= 675 \text{ mul} \\ 2) \text{Filter\_2: } 128 \times 1 \times 1 \times 3 \times 5 \times 5 &= 9600 \text{ mul} \end{aligned} \Bigg] = \boxed{10.275 \text{ mult}}$$

So, for an image with arbitrary size, how much time can we save if we apply depthwise separable convolution. Let's generalize the above examples a little bit. Now, for an input image of size  $H \times W \times D$ , we want to do 2D convolution (stride=1, padding=0) with  $N_c$  kernels of size  $h \times h \times D$ , where  $h$  is even. This transform the input layer ( $H \times W \times D$ ) into the output layer ( $H-h+1 \times W-h+1 \times N_c$ ). The overall multiplications needed is

$$N_c \times h \times h \times D \times (H-h+1) \times (W-h+1)$$

On the other hand, for the same transformation, the multiplication needed for depthwise separable convolution is

$$\begin{aligned} D \times h \times h \times 1 \times (H-h+1) \times (W-h+1) + N_c \times 1 \times 1 \times D \times (H-h+1) \times (W-h+1) \\ = (h \times h + N_c) \times D \times (H-h+1) \times (W-h+1) \end{aligned}$$