

Dimension of the Hidden layer

Let's take the size of the input image to be $n_w^{[0]} \times n_H^{[0]}$.

Typically the filter size f_s , set to be 3, 5 or sometimes 7.

This determines how local we want the model to be.

The size of the hidden layer would be:

$$[n_w^{[0]} - f_s + 1] \times [n_H^{[0]} - f_s + 1]$$

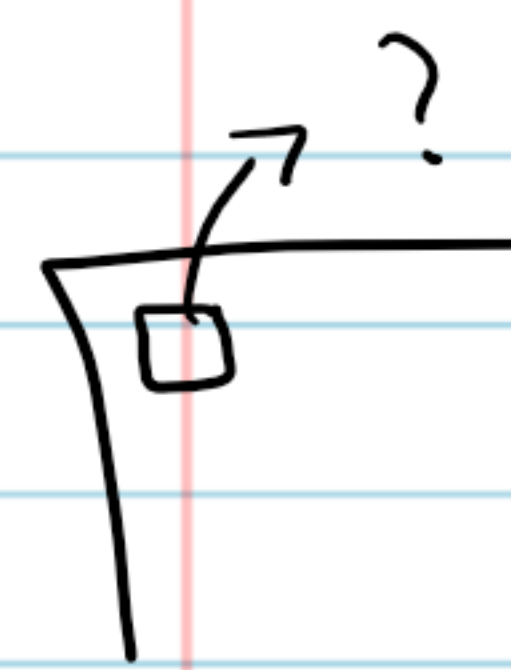
eg. $f_s = 3$

$$\hookrightarrow n_w^{[1]} = n_w^{[0]} - 2, \quad n_H^{[1]} = n_H^{[0]} - 2.$$

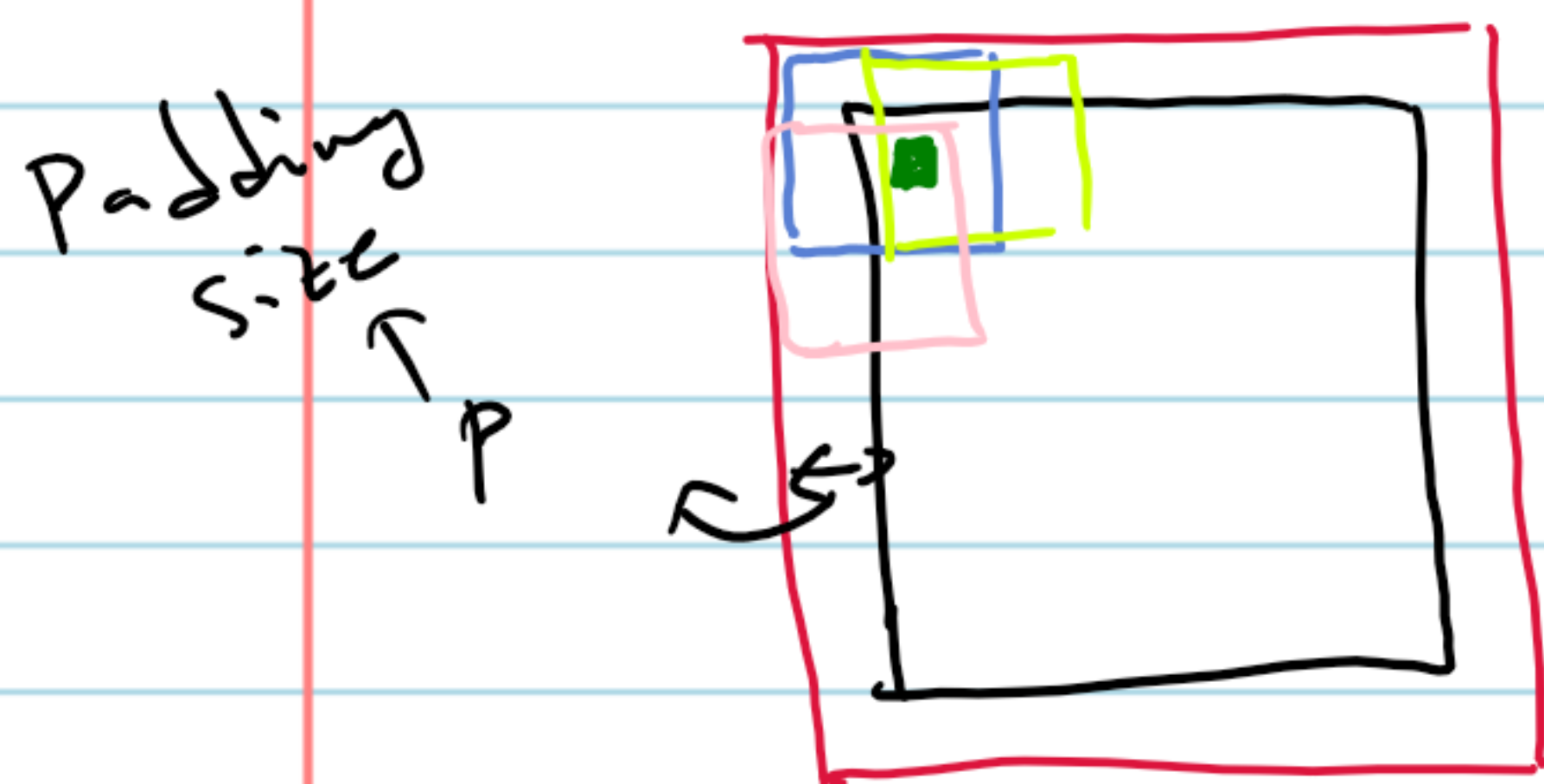
This is basically the number of positions that we can put the filter inside the image.

Padding:

A point in the middle of the image would be involved in 9 nodes of the Hidden layer. But points in the border would have less contributions.



To balance this out, we could add padding to the image.



$$n_w^{[1]} = n_w^{[0]} - f_s + 1 + 2ps$$

$$n_H^{[1]} = n_H^{[0]} - f_s + 1 + 2ps$$

Padding is also used to keep the hidden layer size the same as the input image.

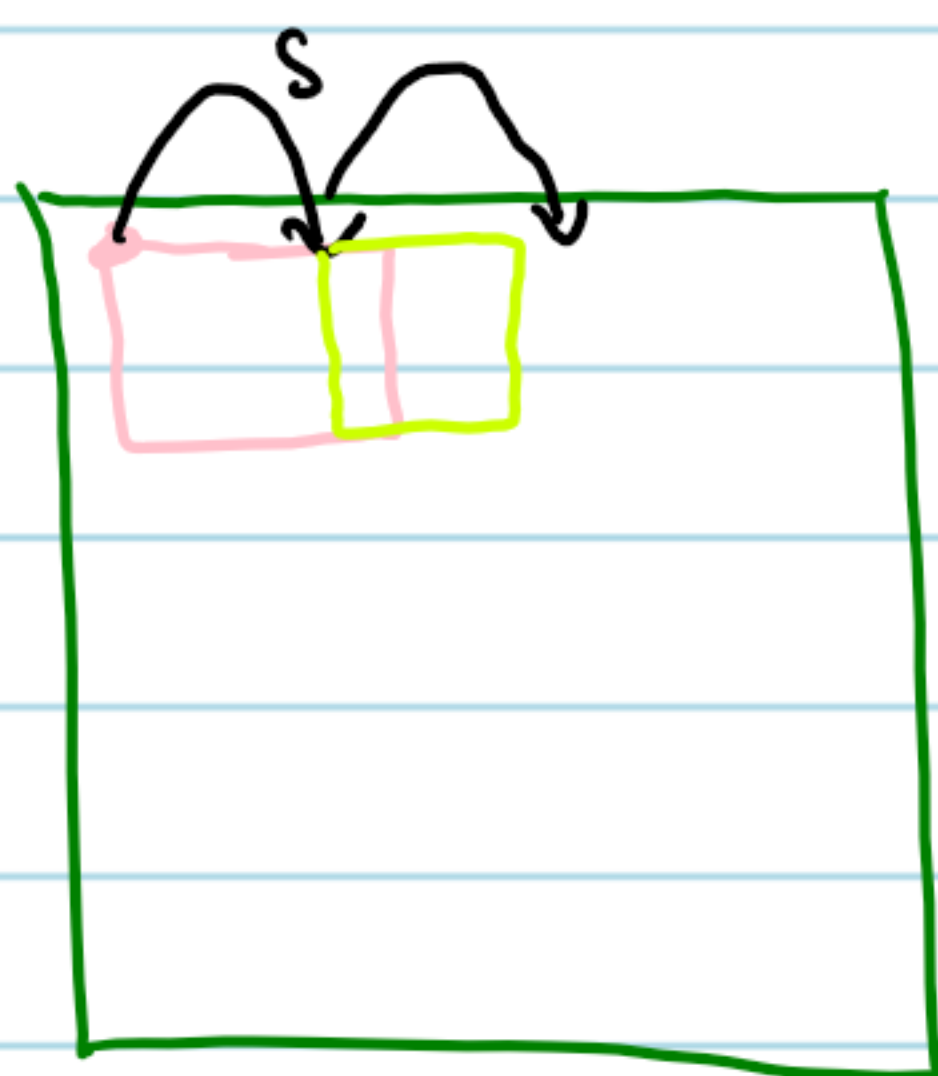
$$n_w = n_w - f_s + 1 + 2ps \Rightarrow ps = \frac{f_s - 1}{2}$$

e.g. for $f_s = 3 \rightarrow ps = 1$

$f_s = 5 \rightarrow ps = 2$

Stride

We could also set the stepsize for moving the filter over the image.



Typically, we take $s=1$ but it can be treated as another hyperparameter.

With stride of s ,

$$n_w^{[l]} = \left\lfloor \frac{n_w^{[0]} - f_s + 2ps}{s} \right\rfloor + 1$$

$$n_h^{[l]} = \left\lfloor \frac{n_h^{[0]} - f_s + 2ps}{s} \right\rfloor + 1$$

So the hyperparameters for the convolution would be

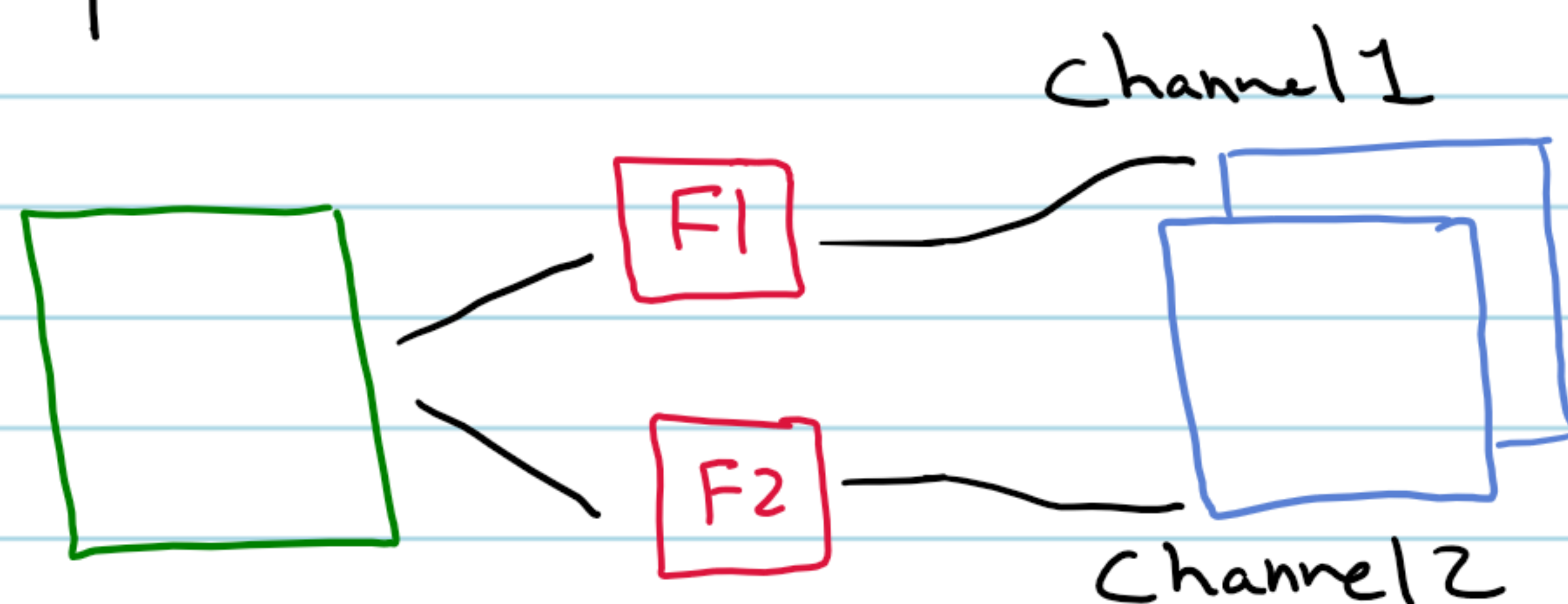
Filter-size: f_s typically 3, 5, 7

Padding size: ps Typically: "Same", "Valid" \rightarrow no padding

Strides, s " 1.

Channels

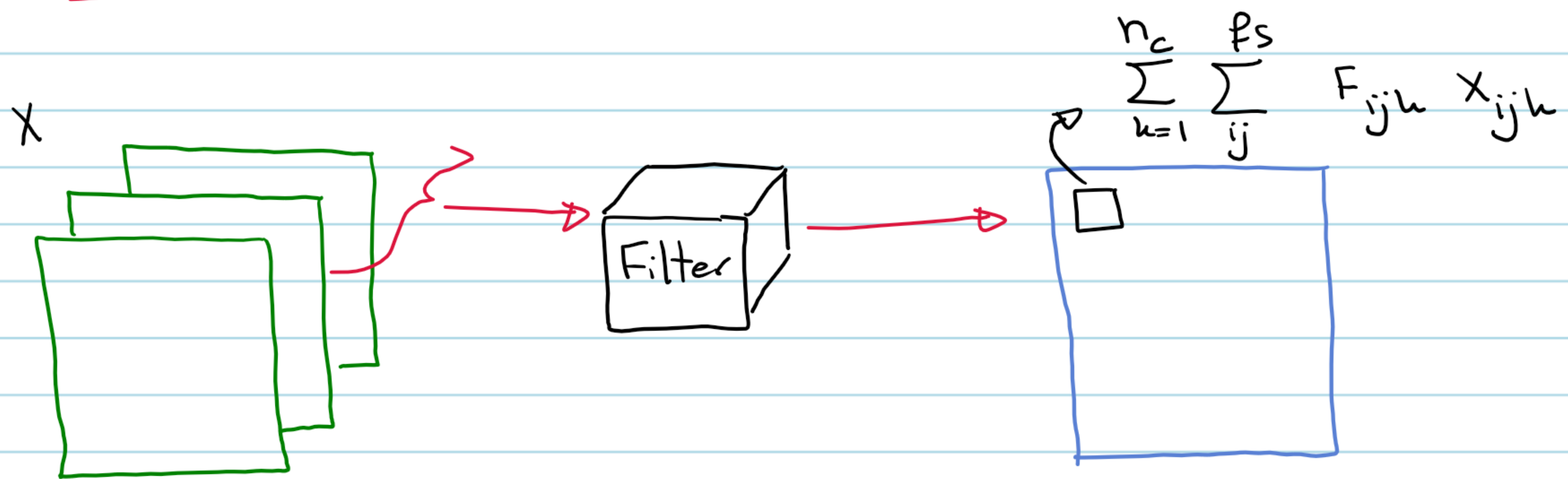
A filter is optimized to extract a certain feature through the training process. But it may not be enough to rely on only one single feature. For instance, for edge detection, we probably want to have both the vertical & horizontal edges. This cannot be done with a single filter. So we use multiple filters.



We refer to the # of filters as n_c , the number of channels. This, sort of, represents the number of features that the conv-layer extracts.

Even the input could have (as usually does) multiple channels. For a typical image, these could be RGB channels.

Applying Filter to layers with channels:

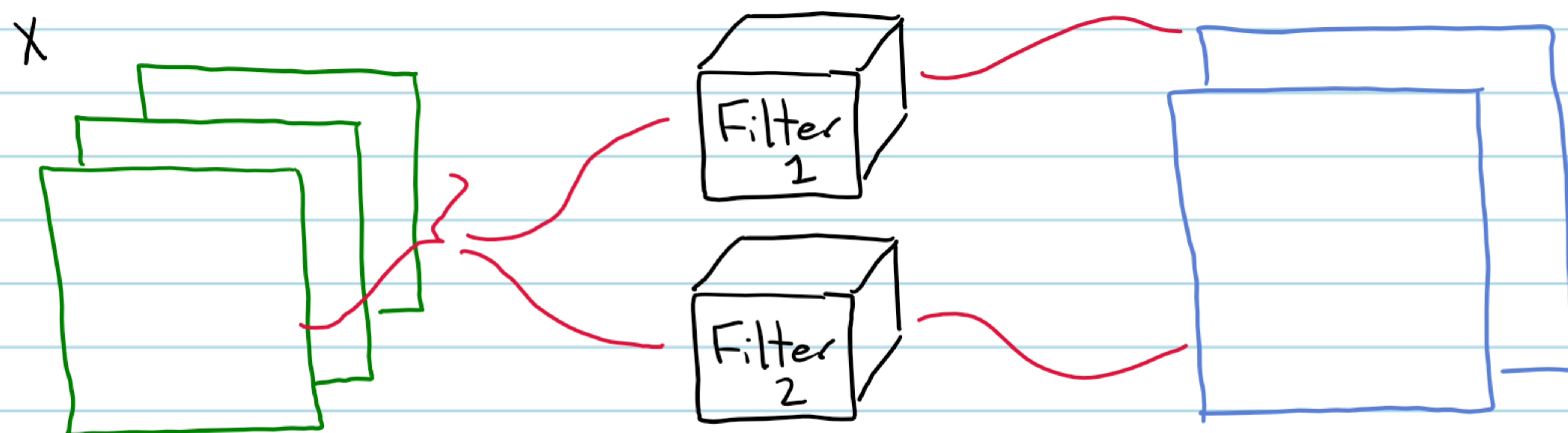


* Filter needs to have the same dimension as the input layer. That means it should be

$$f_s \times f_s \times n_c^{[l-1]}.$$

* But the outcome would be one channel, i.e. each filter makes one channel.

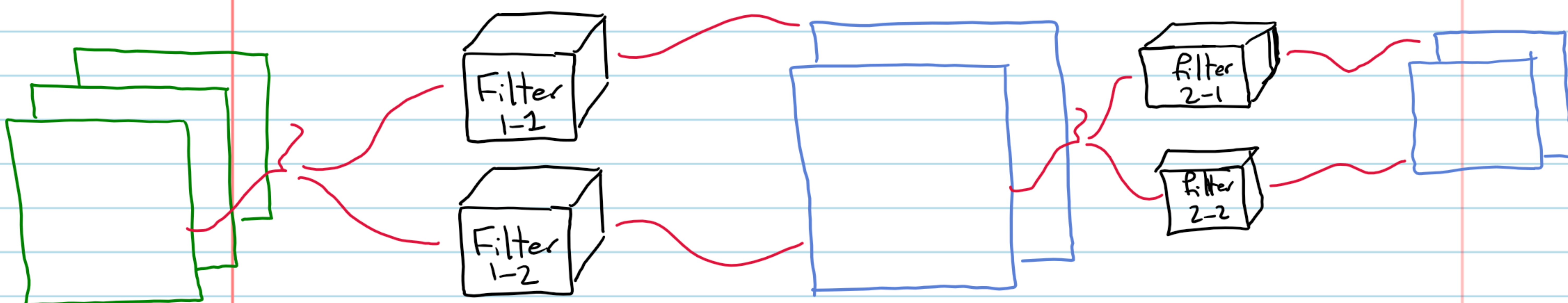
* To increase the # channels, we would need to add more filters.



The number of channels of the hidden layer is given by the # of filters.

Deep conv. net:

We can apply more than just one layer of convolution.



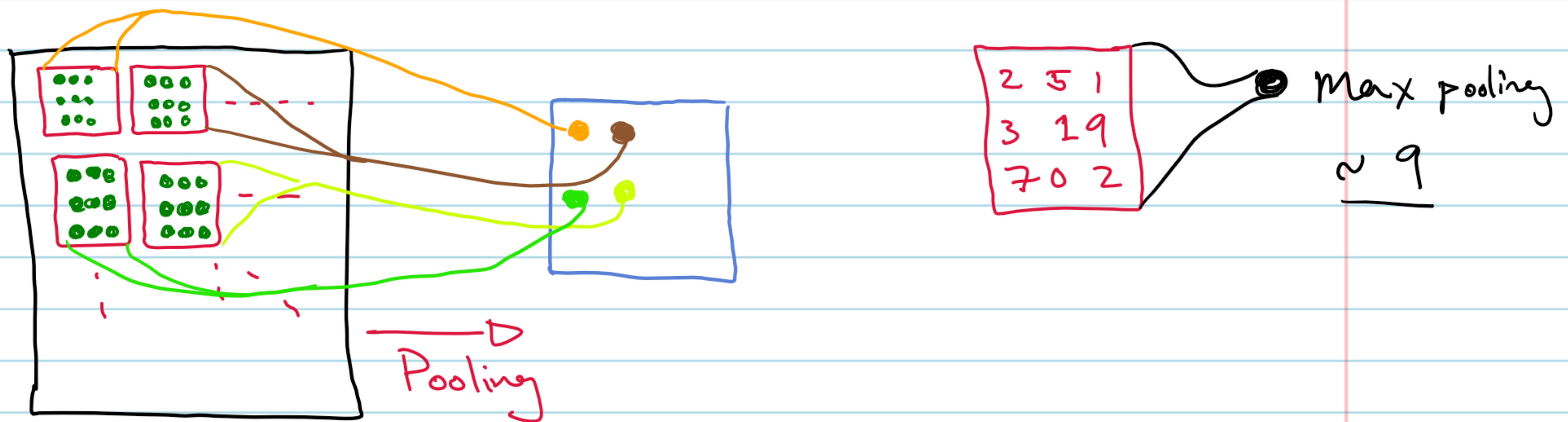
This would mean that the first layer extracts certain features and the following layer builds one top of that to extract more complex features.

Pooling:

This is fairly similar to coarse-graining.

The idea is that we want to reduce the size of the hidden layers as we go deeper in the neural net, i.e. n_H & n_W should decrease. This is b/c we are looking for a few key features that help with the classification task.

With max pooling, we take a window and move it over the layer and pick the max value in that window.



In some sense, the max pooling checks if a certain property exists anywhere in the window of the layer.

We can also take the average which is known as "average pooling".

Hyperparameters of Pooling

Windows size : w_s typically 2 3 max vs avg
stride 2 2

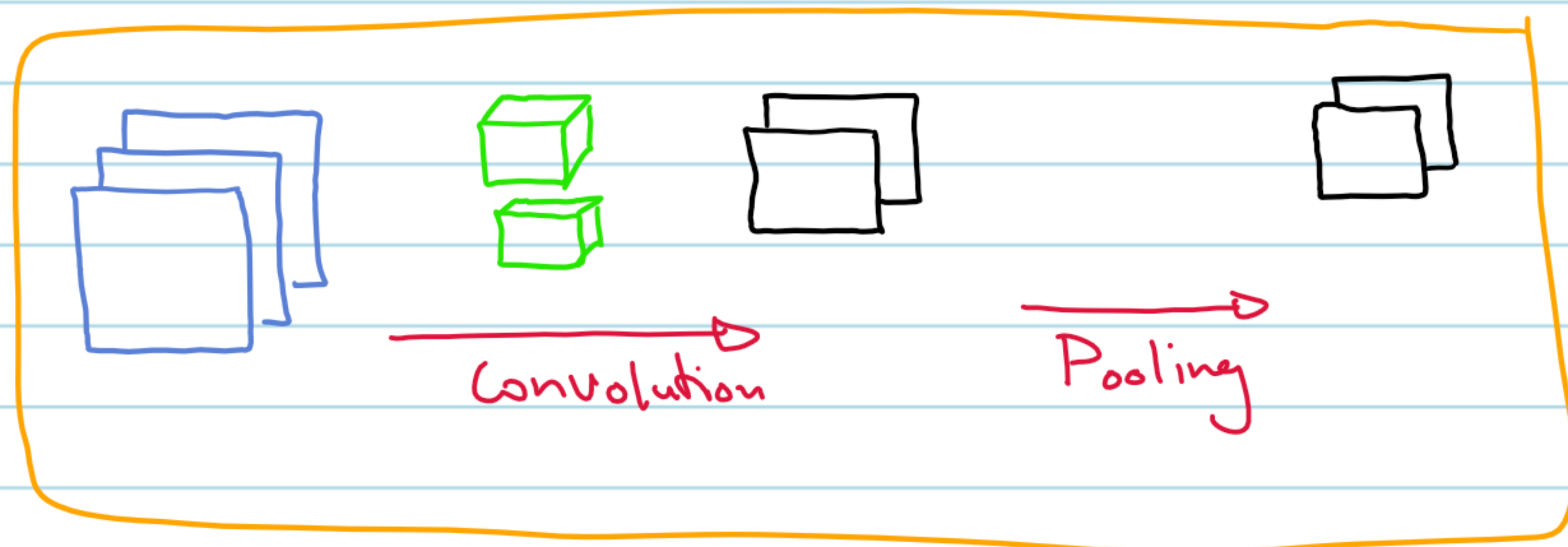
But no learnable parameter.

*Note Pooling is applied to channels separately and does not change n_c .

Size $n_w \times n_H \times n_c$

$$\left\lfloor \frac{n_w - w_s}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_H - w_s}{s} + 1 \right\rfloor \times n_c$$

One layer



One unit of convolution

Fully Connected Layer

Often a full conv_net is composed of several convolution units followed by a couple of fully connected layers.

One can think of the convolution units as feature extractors and the FC layers as the classifier that works/uses the extracted features for classification.

Full Network

