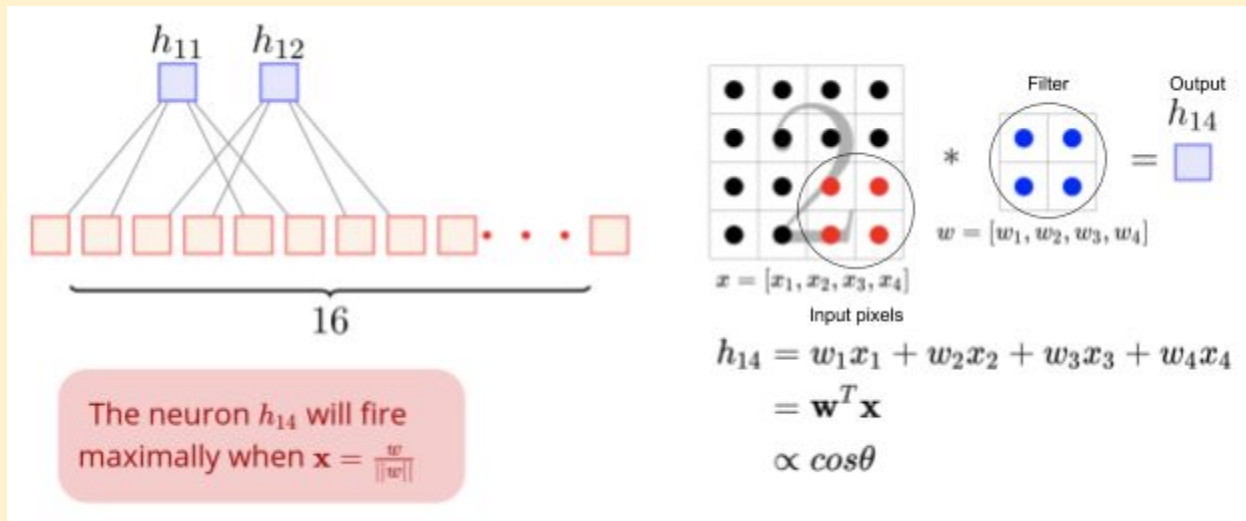


Visualising filters

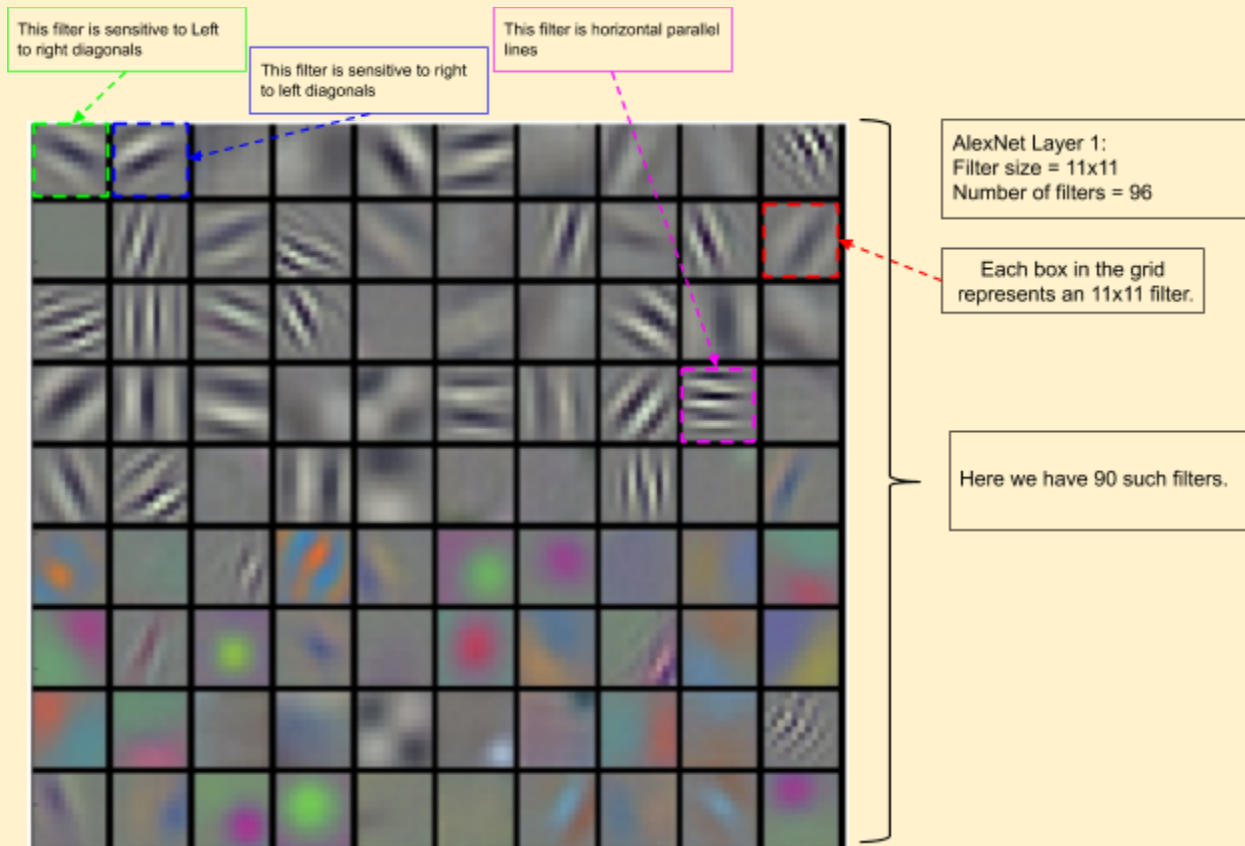
What does a filter capture?

1. We have dealt with filters in all our CNN models so far. Now the question is, what exactly does a filter capture?
2. Let's look at the working of a 2x2 filter on a 4x4 input image



- a. Here, the input image is 4x4 while the filter is 2x2
- b. The red input pixel vector $\mathbf{x} = [x_1, x_2, x_3, x_4]$
- c. The weight vector $\mathbf{w} = [w_1, w_2, w_3, w_4]$
- d. By convolving the input pixels with the filter, we get the output
- e. Output $h_{14} = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4$
- f. $h_{14} = \mathbf{w}^T \mathbf{x}$ (This is the same as the dot product between the two vectors)
- g. $h_{14} \propto \cos(\theta)$ [where θ is the angle between the two vectors] $\cos(\theta) = \frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\| \|\mathbf{x}\|}$
- h. Now for certain inputs, we want the filter to fire (give a high value).
- i. Now, h_{14} will be high when $\cos(\theta)$ is high ($\cos(\theta) = 1$), i.e. when θ is 0. This implies the two vectors \mathbf{w} and \mathbf{x} are in the same direction.
- j. So, we can say that an input vector which aligns with a filter vector yields maximum output.
- k. The neuron h_{14} will fire maximally when $\mathbf{x} = \frac{\mathbf{w}}{\|\mathbf{w}\|}$ (\mathbf{x} is a unit vector in the direction of \mathbf{w})
- l. Thus, when we **slide the 2x2 filter \mathbf{w}** across the 4x4 input region, whenever we **reach a 2x2 region \mathbf{x}** that looks exactly like the filter, we get a high output. For all other regions which do not align with the filter, the output is low.

3. Now, let us visualize the filters in AlexNet



- a. The above image shows us how different patterns are identified by different filters.