

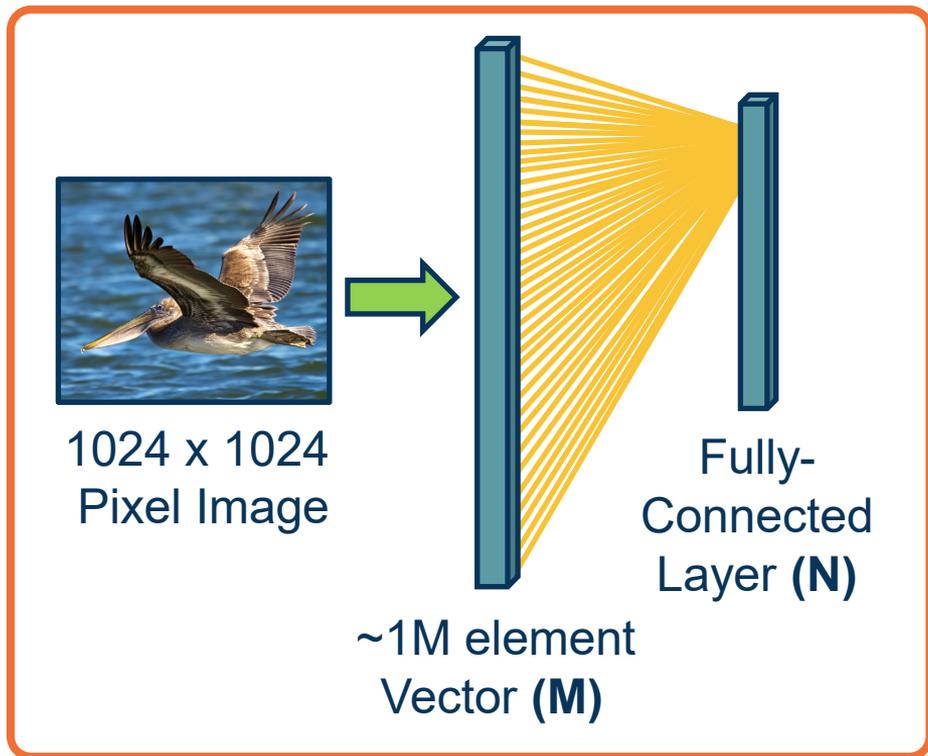
Topics:

- Convolutional Neural Networks

CS 4644-DL / 7643-A

ZSOLT KIRA

The connectivity in linear layers **doesn't** always make sense



How many parameters?

● $M \cdot N$ (weights) + N (bias)

Hundreds of millions of
parameters **for just one layer**

**More parameters => More
data needed**

Is this necessary?

Limitation of Linear Layers

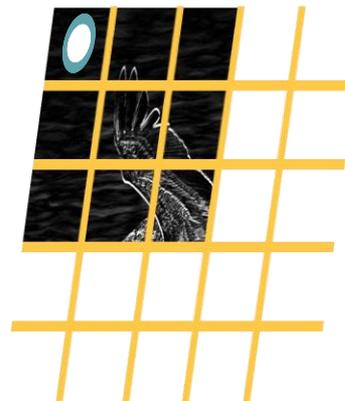
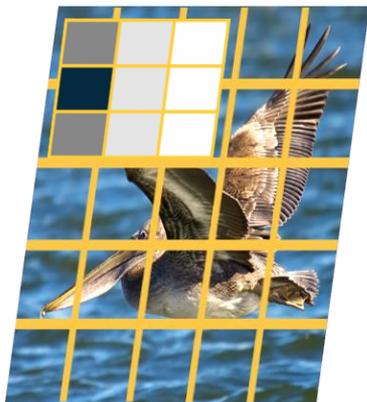
$$X(0:2,0:2) = \begin{bmatrix} 200 & 150 & 150 \\ 100 & 50 & 100 \\ 25 & 25 & 10 \end{bmatrix}$$

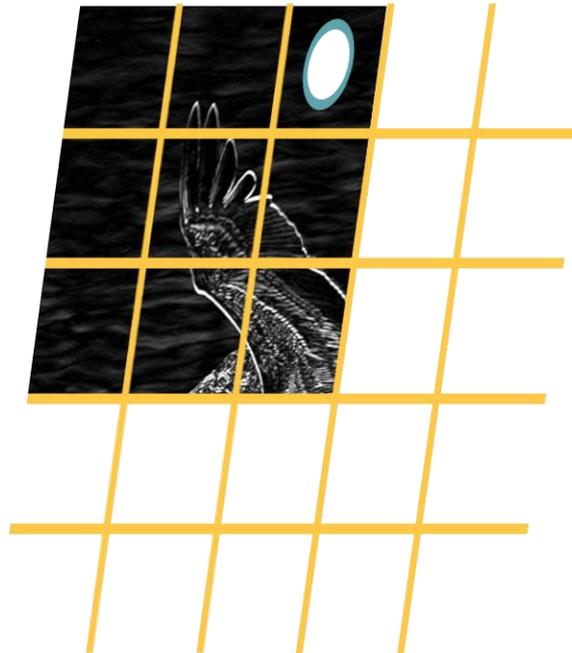
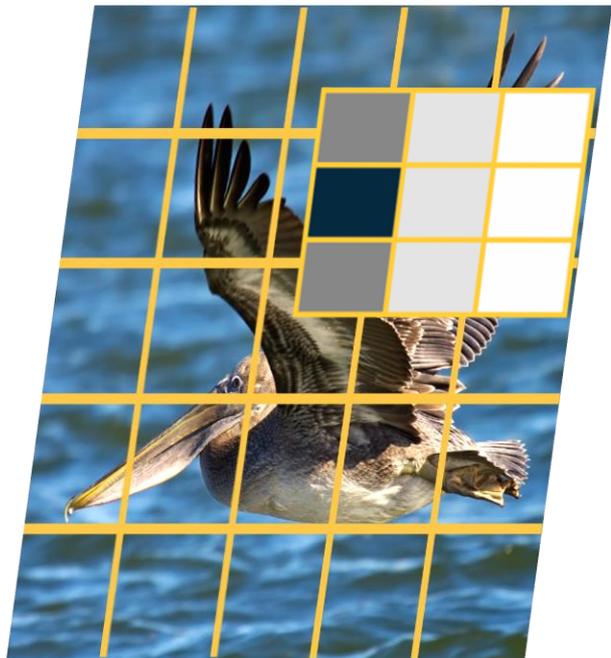
$$K' = \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$



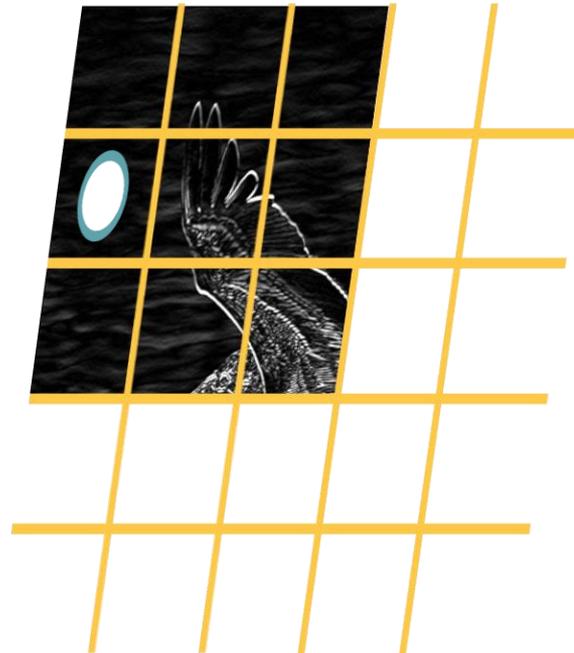
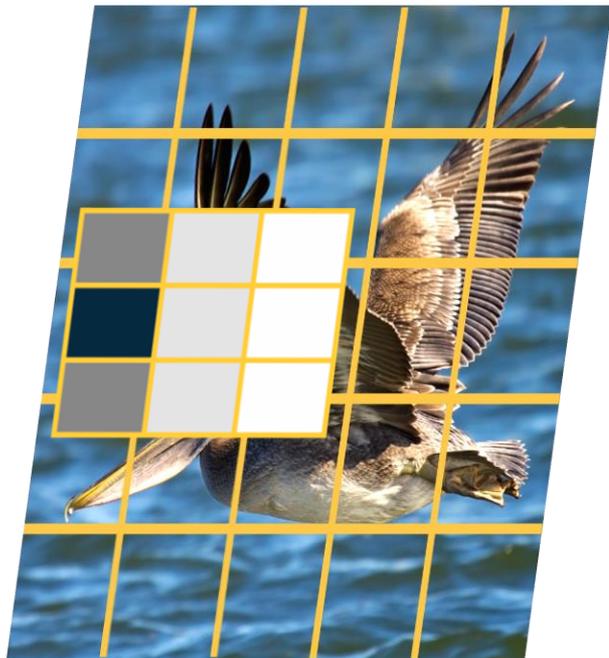
$$X(0:2,0:2) \cdot K' = 65 + \text{bias}$$

Dot product
(element-wise multiply and sum)

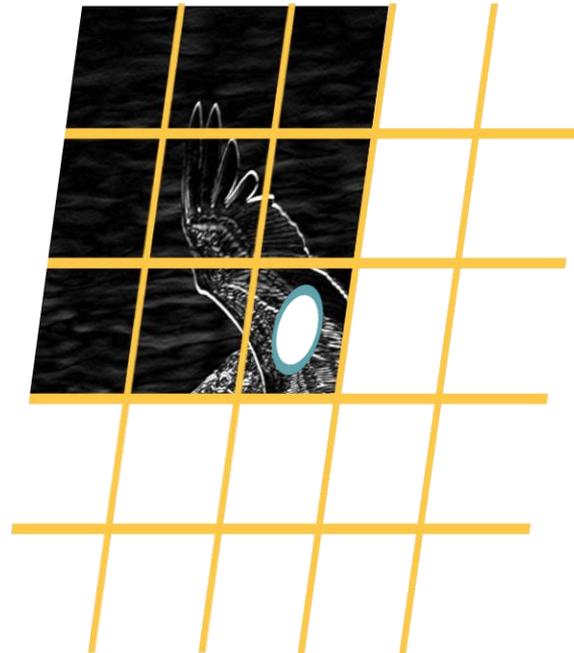
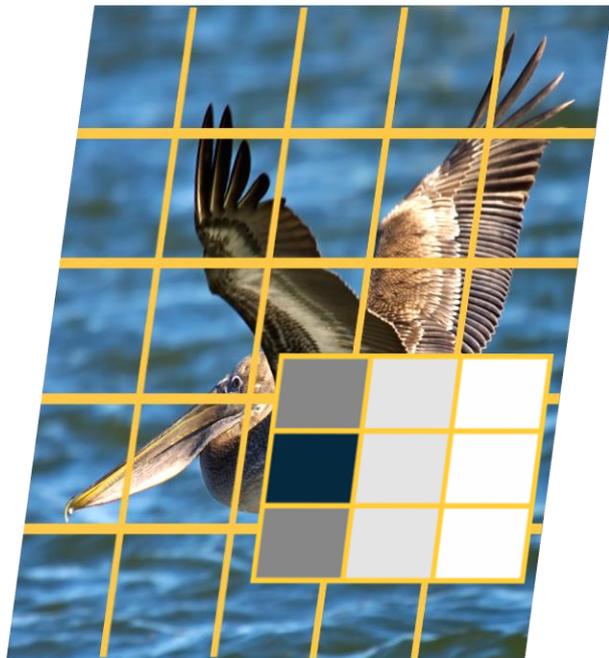




Convolution and Cross-Correlation



Convolution and Cross-Correlation



Convolution and Cross-Correlation

Why Bother with Convolutions?

Convolutions are just **simple linear operations**

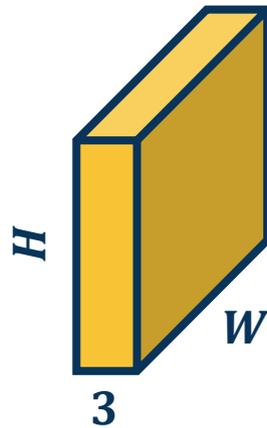
Why bother with this and not just say it's a linear layer with small receptive field?

- ◆ There is a **duality** between them during backpropagation
- ◆ Convolutions have **various mathematical properties** people care about
- ◆ This is **historically** how it was inspired

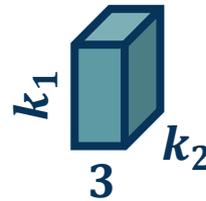


We have shown inputs as a **one-channel image** but in reality they have three channels (red, green, blue)

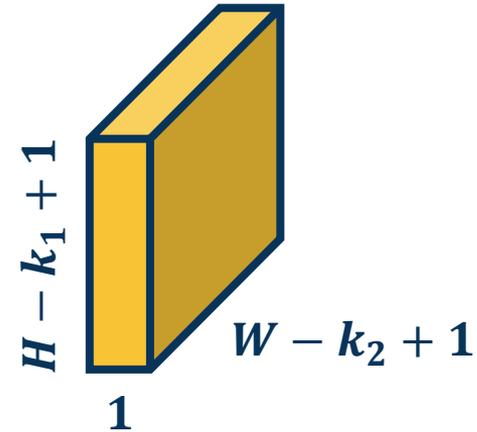
🟡 In such cases, we have **3-channel kernels!**



Image



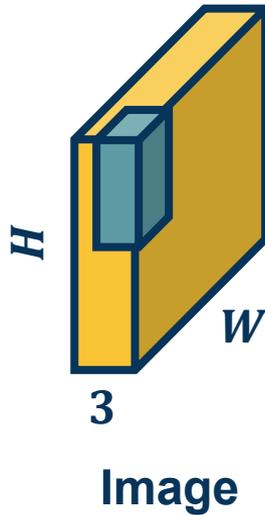
Kernel



Feature Map

We have shown inputs as a **one-channel image** but in reality they have three channels (red, green, blue)

- ◆ In such cases, we have **3-channel kernels!**



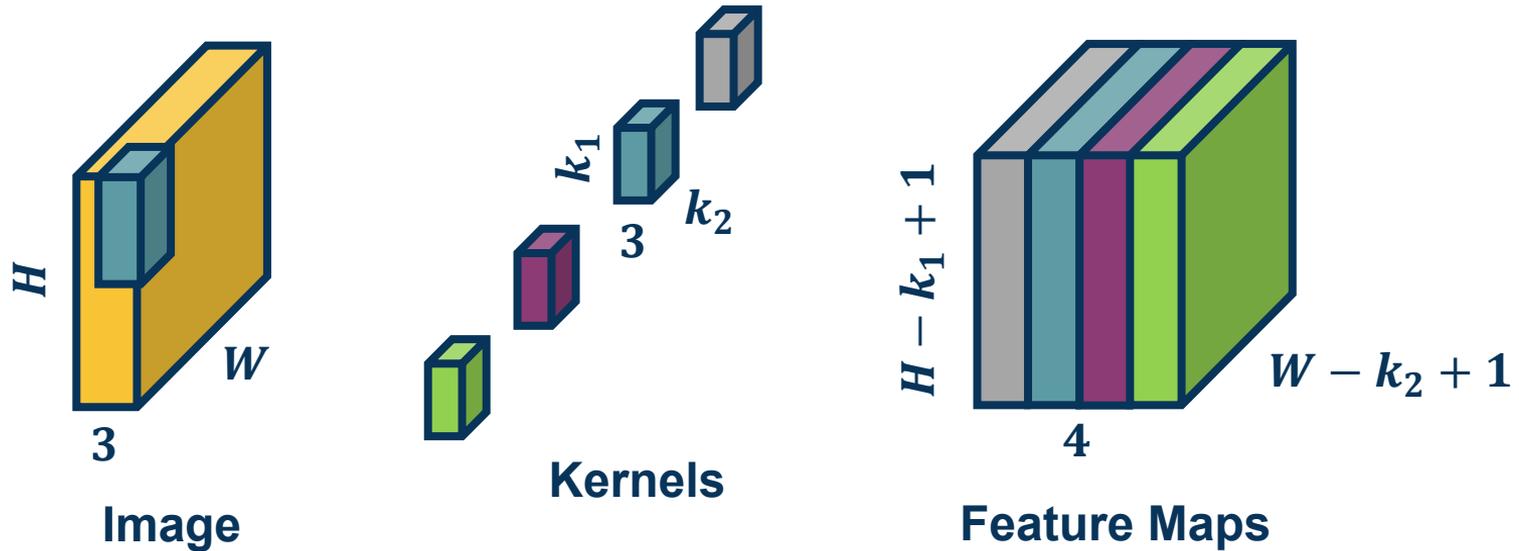
Similar to before, we perform **element-wise multiplication** between kernel and image patch, summing them up (**dot product**)

- ◆ Except with $k_1 * k_2 * 3$ values

We can have **multiple kernels per layer**

- ◆ We stack the feature maps together at the output

Number of channels in output is equal to *number of kernels*

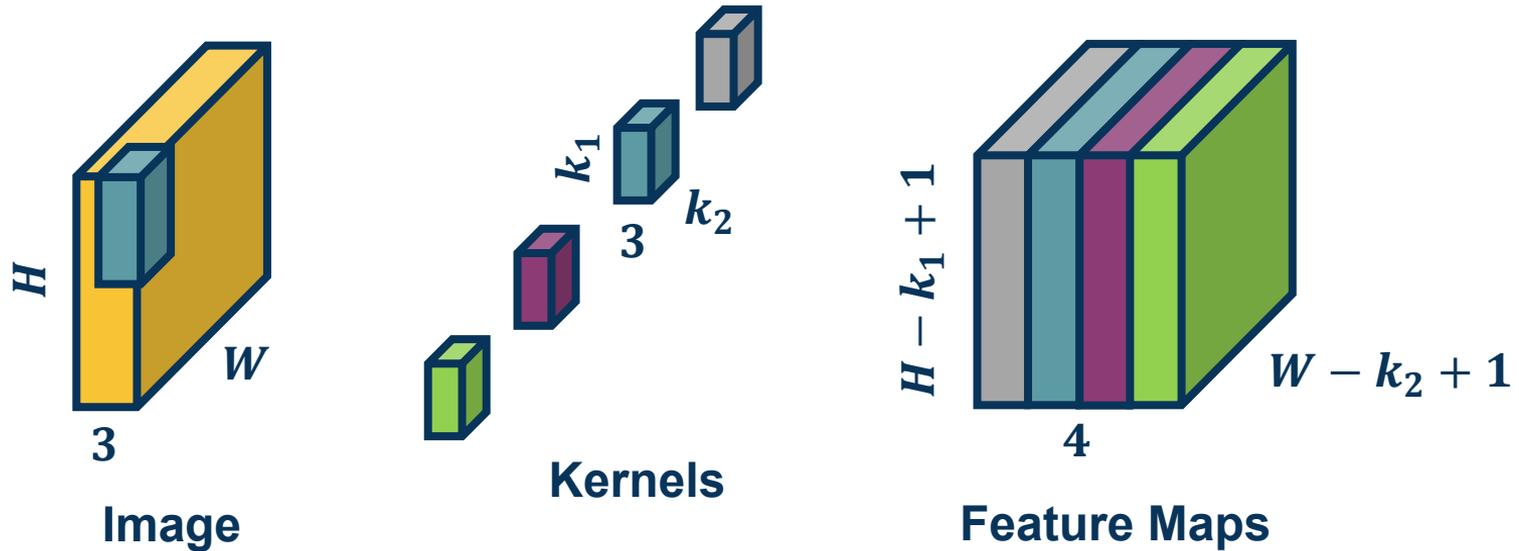


Multiple Kernels

Number of parameters with N filters is: $N * (k_1 * k_2 * 3 + 1)$

Example:

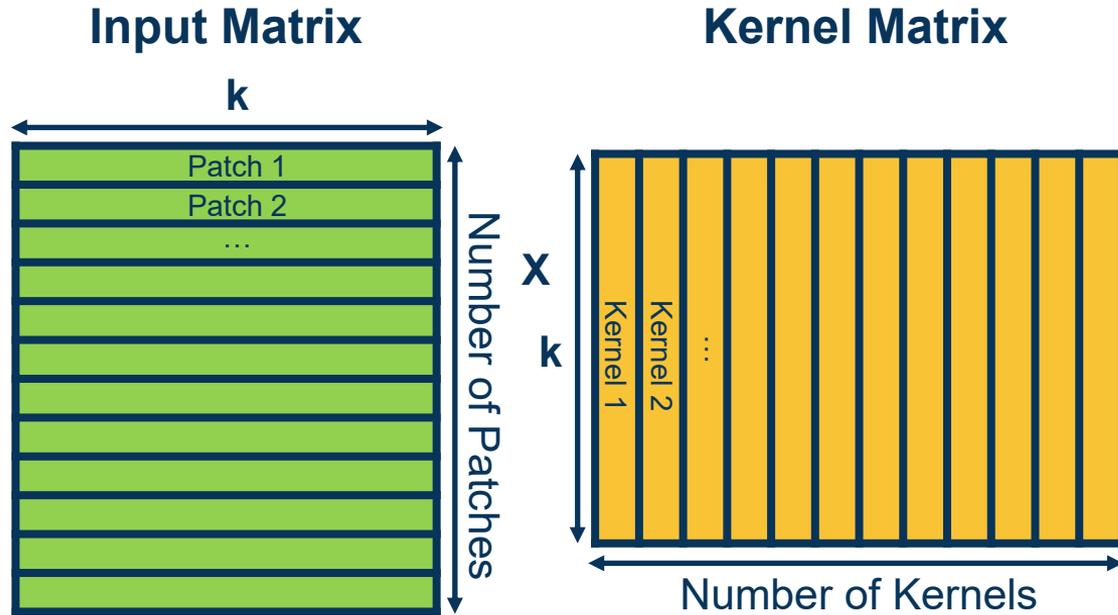
$k_1 = 3, k_2 = 3, N = 4$ input channels = 3, then $(3 * 3 * 3 + 1) * 4 = 112$



Number of Parameters

Just as before, in practice we can **vectorize** this operation

- Step 2: Multiple patches by kernels



Adapted from: <https://petewarden.com/2015/04/20/why-gemm-is-at-the-heart-of-deep-learning/>

Backwards Pass for Convolution Layer

It is instructive to calculate **the backwards pass** of a convolution layer

- Similar to fully connected layer, will be **simple vectorized linear algebra operation!**
- We will see a **duality** between cross-correlation and convolution

$$K = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$$



$$K' = \begin{bmatrix} 9 & 8 & 7 \\ 6 & 5 & 4 \\ 3 & 2 & 1 \end{bmatrix}$$

$$y(r, c) = (x * k)(r, c) = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x(r+a, c+b) k(a, b)$$

(0, 0)

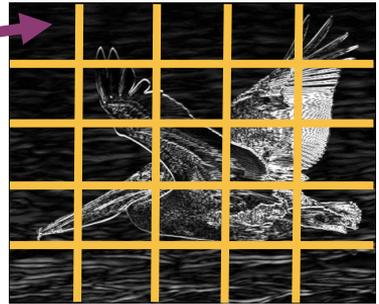


$W = 5$ ($H - 1, W - 1$)

(0, 0)



$k_2 = 3$ ($k_1 - 1, k_2 - 1$)



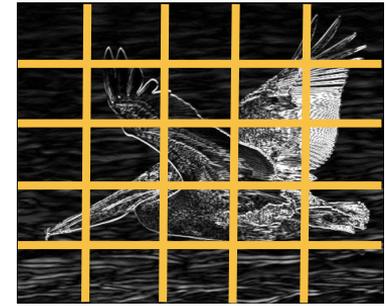
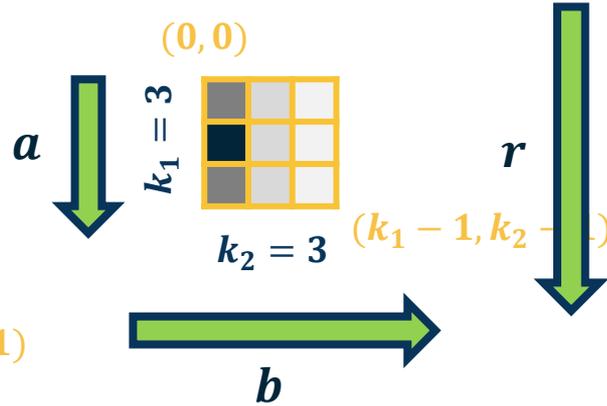
Recap: Cross-Correlation

$$y(r, c) = (x * k)(r, c) = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x(r+a, c+b) k(a, b)$$

(0, 0)



$W = 5$ $(H - 1, W - 1)$



c

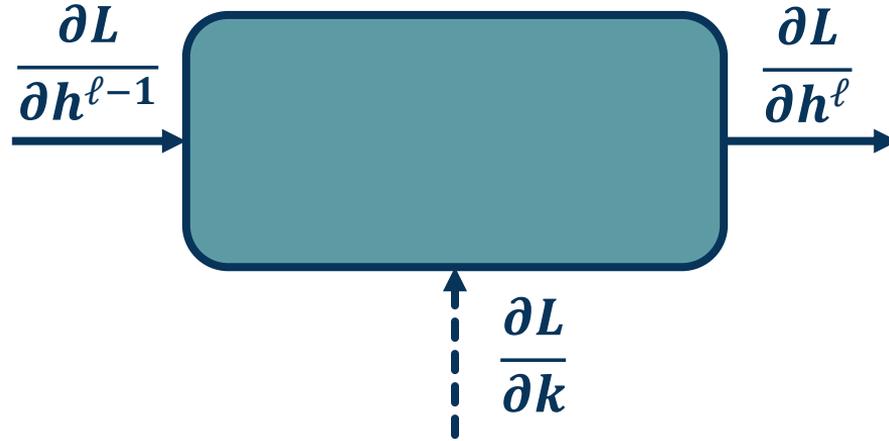
Some simplification: 1 channel input, 1 kernel (channel output), padding (here 2 pixels on right/bottom) to make output the same size

$$y(r, c) = (x * k)(r, c) = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x(r+a, c+b) k(a, b)$$

$$|y| = H \times W$$

$$\frac{\partial L}{\partial y} ? \quad \text{Assume size } H \times W \text{ (add padding)}$$

$$\frac{\partial L}{\partial y(r, c)} \quad \text{to access element}$$



$$\frac{\partial L}{\partial h^{\ell-1}} = \frac{\partial L}{\partial h^{\ell}} \frac{\partial h^{\ell}}{\partial h^{\ell-1}}$$

Gradient for passing back

$$\frac{\partial L}{\partial k} = \frac{\partial L}{\partial h^{\ell}} \frac{\partial h^{\ell}}{\partial k}$$

Gradient for weight update

(weights = k, i.e. kernel values)

Gradient for Convolution Layer

$$\frac{\partial L}{\partial k} = \frac{\partial L}{\partial h^\ell} \frac{\partial h^\ell}{\partial k}$$

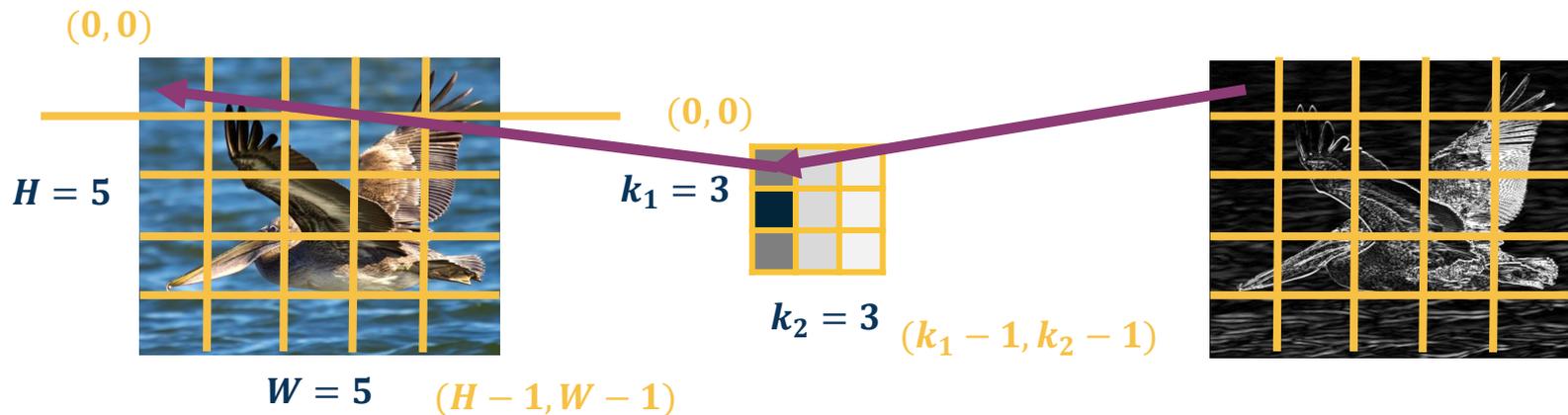
What does this weight affect at the output?

Gradient for weight update

Calculate one pixel at a time

$$\frac{\partial L}{\partial k(a,b)} = \frac{\partial L}{\partial h^\ell} \frac{\partial h^\ell}{\partial k(a,b)}$$

Everything!



What a Kernel Pixel Affects at Output

Need to incorporate all upstream gradients:

$$\left\{ \frac{\partial L}{\partial y(0,0)}, \frac{\partial L}{\partial y(0,1)}, \dots, \frac{\partial L}{\partial y(H,W)} \right\}$$

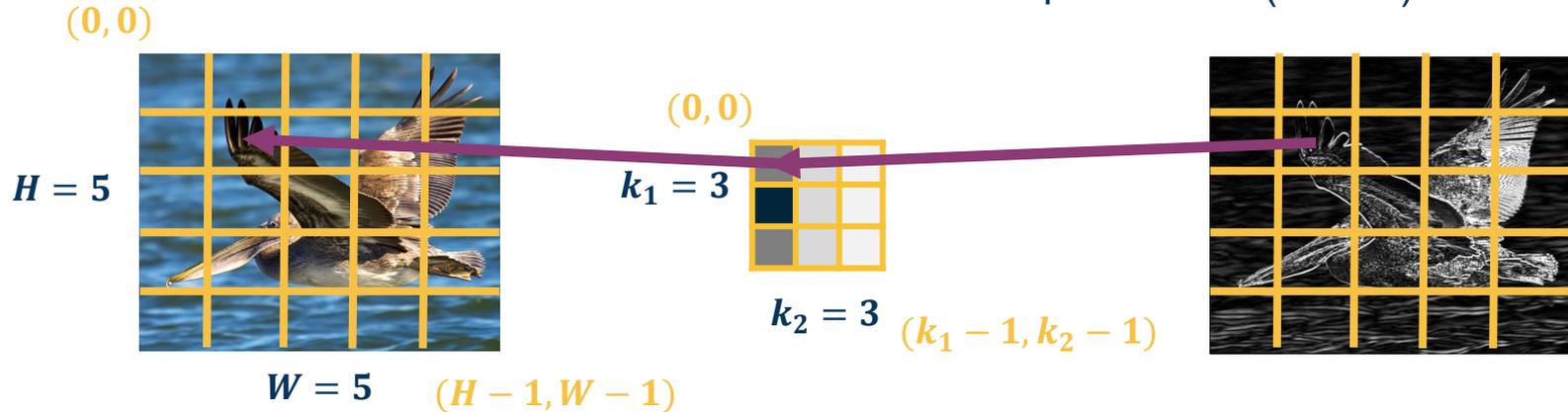
Chain Rule:

$$\frac{\partial L}{\partial k(a,b)} = \sum_{r=0}^{H-1} \sum_{c=0}^{W-1} \frac{\partial L}{\partial y(r,c)} \frac{\partial y(r,c)}{\partial k(a,b)}$$

Sum over
all output
pixels

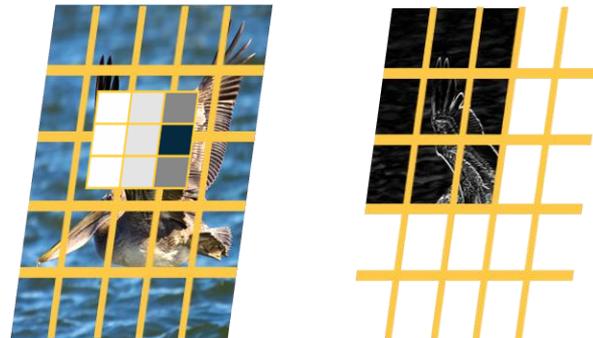
Upstream
gradient
(known)

We will
compute



Chain Rule over all Output Pixels

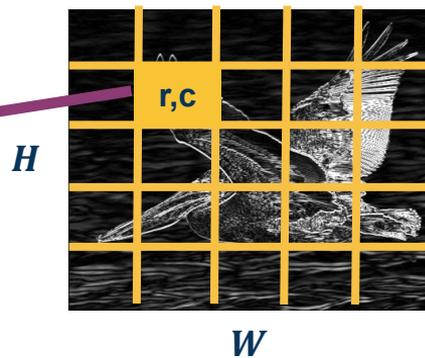
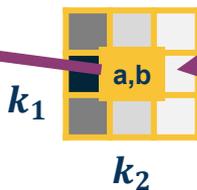
$$\frac{\partial y(r, c)}{\partial k(a, b)} = ?$$



For output at $y(r,c)$, where is placement of kernel?

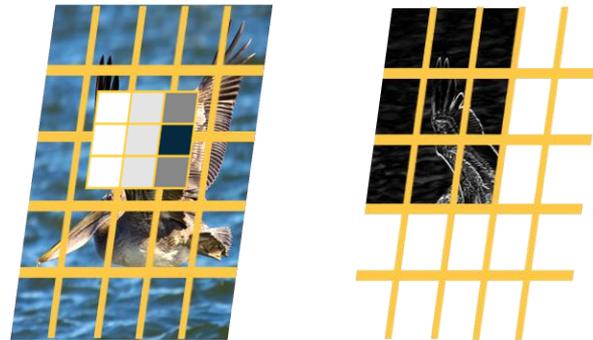


?

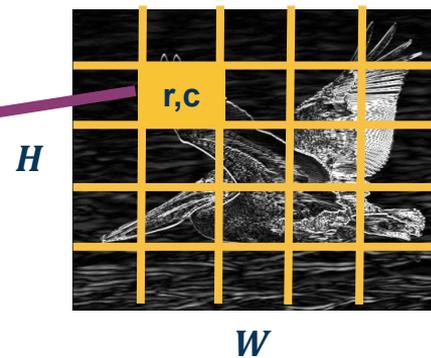
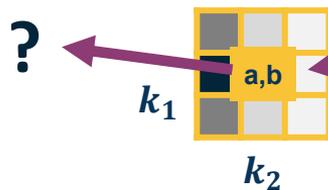


Chain Rule over all Output Pixels

$$\frac{\partial y(r, c)}{\partial k(a, b)} = ?$$



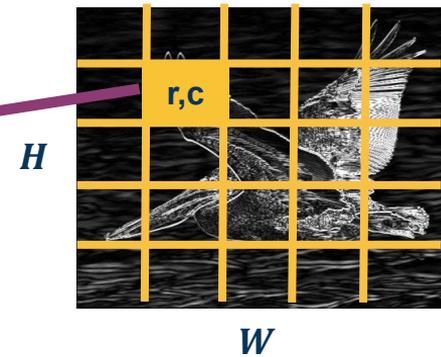
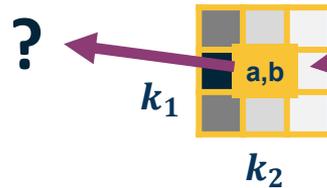
Which x pixel is multiplied by $k(a, b)$



$$\frac{\partial y(r, c)}{\partial k(a, b)} = ?$$

Reasoning:

- Cross-correlation is just “dot product” of kernel and input patch (weighted sum)
- When at pixel $y(r, c)$, kernel is on input x such that $k(0, 0)$ is multiplied by $x(r, c)$
- But we want derivative w.r.t. $k(a, b)$
 - $k(0, 0) * x(r, c)$, $k(1, 1) * x(r + 1, c + 1)$, $k(2, 2) * x(r + 2, c + 2)$
 - => in general $k(a, b) * x(r + a, c + b)$
 - Just like before in fully connected layer, partial derivative w.r.t. $k(a, b)$ *only* has this term (other x terms go away because not multiplied by $k(a, b)$).



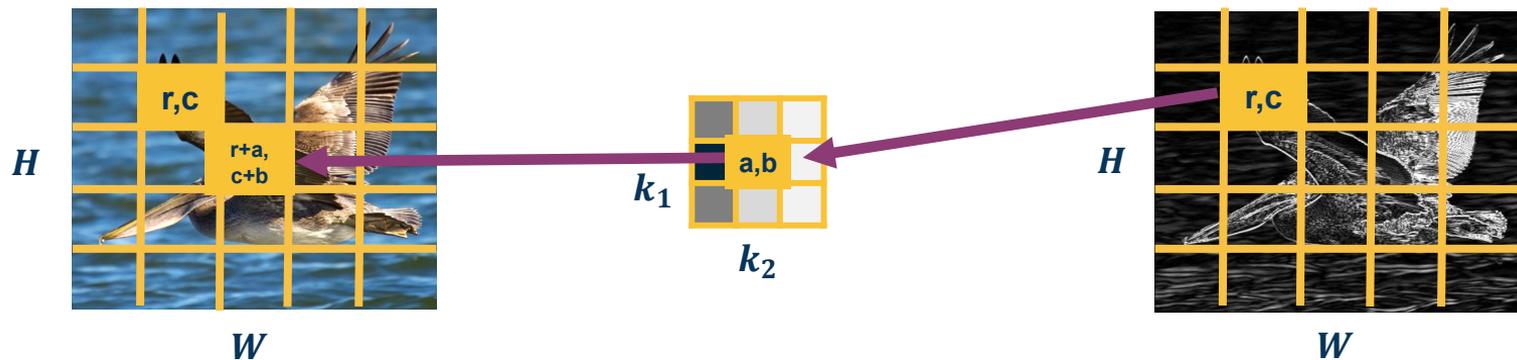
Chain Rule over all Output Pixels

$$\frac{\partial y(r, c)}{\partial k(a, b)} = x(r + a, c + b)$$

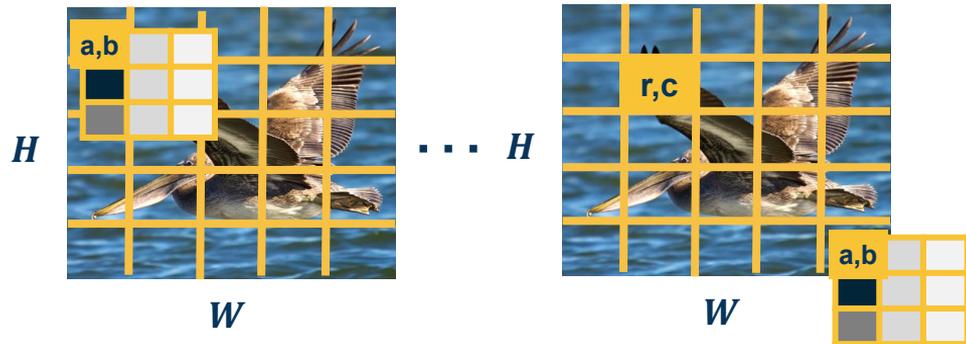
$$\frac{\partial L}{\partial k(a, b)} = \sum_{r=0}^{H-1} \sum_{c=0}^{W-1} \frac{\partial L}{\partial y(r, c)} x(r + a, c + b)$$

Does this look familiar?

Cross-correlation
between upstream
gradient and input!
(until $k_1 \times k_2$ output)



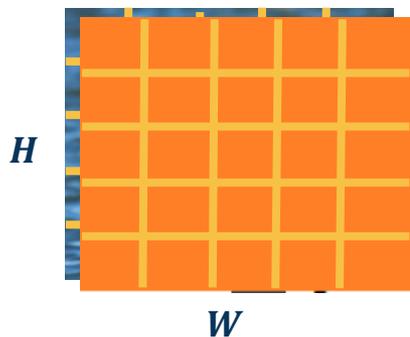
Forward Pass



Does this look familiar?

Cross-correlation
between upstream
gradient and input!
(until $k_1 \times k_2$ output)

Backward Pass $k(0, 0)$



Backward Pass $k(2, 2)$



Forward and Backward Duality

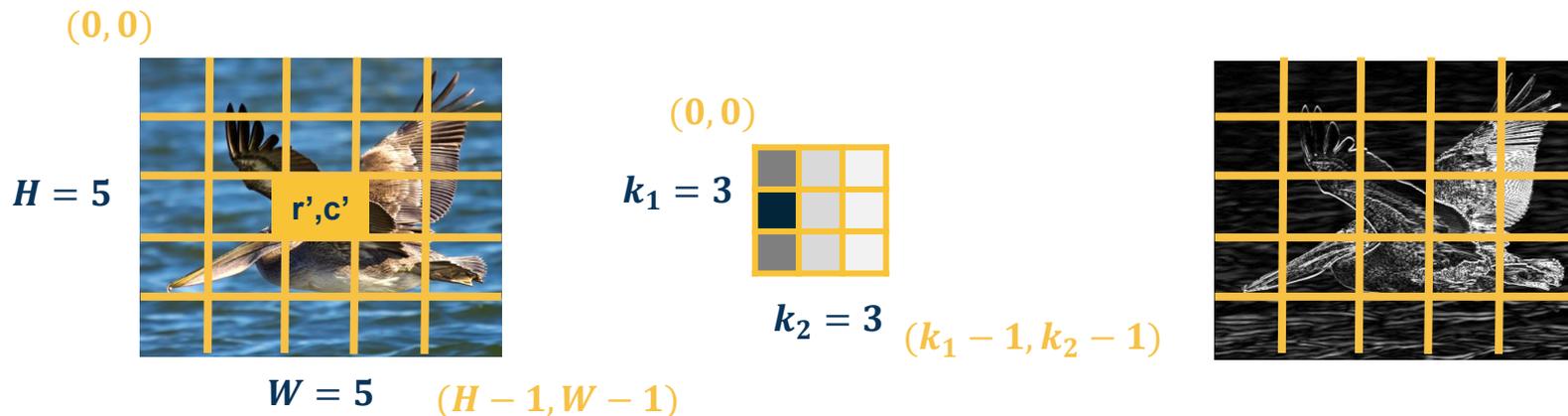
$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x}$$

Gradient for input (to pass to prior layer)

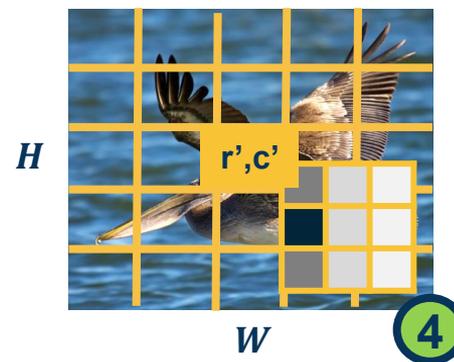
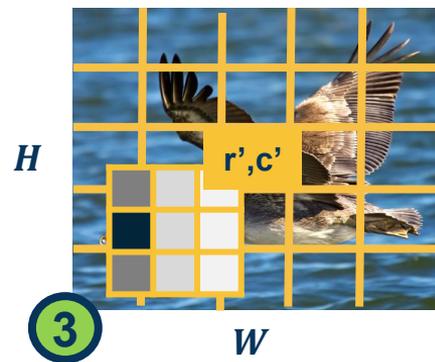
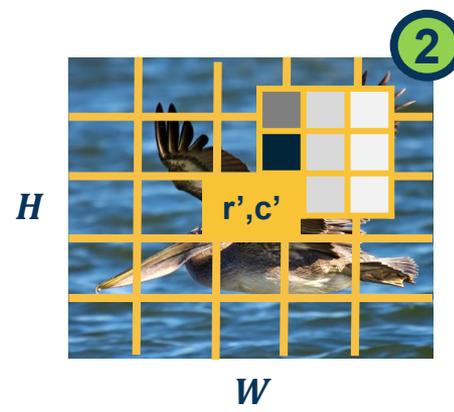
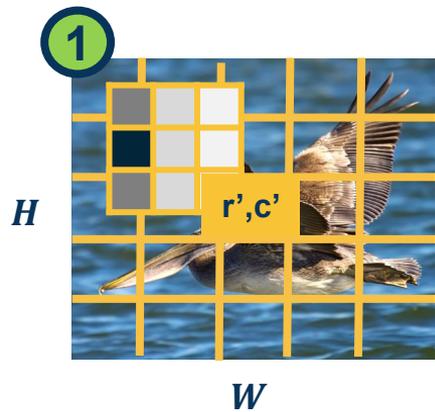
Calculate one pixel at a time $\frac{\partial L}{\partial x(r', c')}$

What does this input pixel affect at the output?

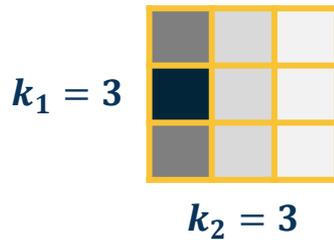
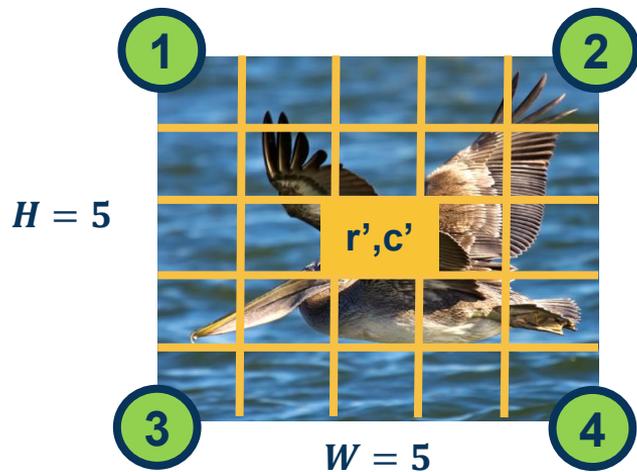
Neighborhood around it (where part of the kernel touches it)



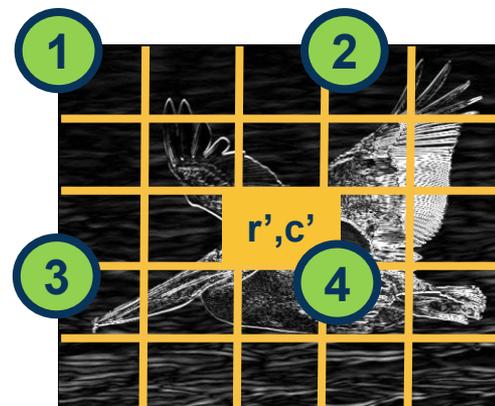
What an Input Pixel Affects at Output



Extents of Kernel Touching the Pixel



$$(r' - k_1 + 1, \\ c' - k_2 + 1)$$



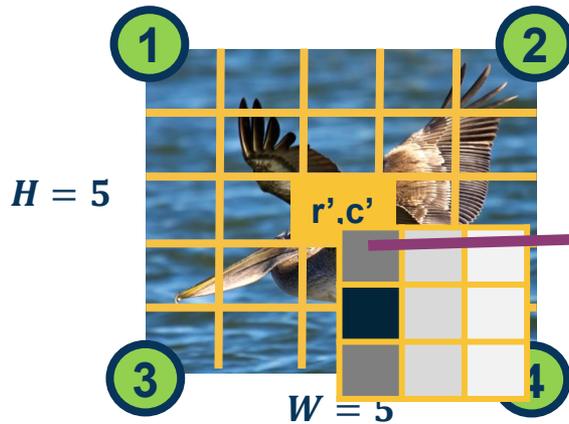
This is where the corresponding locations are for the **output**

Extents at the Output

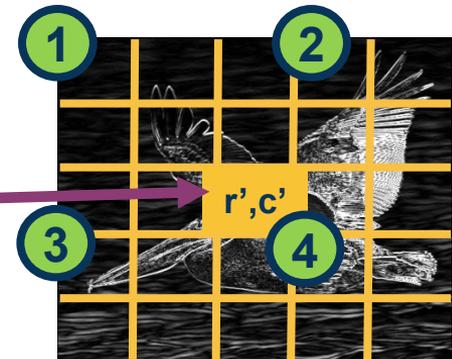
Chain rule for affected pixels (sum gradients):

$$\frac{\partial L}{\partial x(r', c')} = \sum_{\text{Pixels } p} \frac{\partial L}{\partial y(p)} \frac{\partial y(p)}{\partial x(r', c')}$$

$$\frac{\partial L}{\partial x(r', c')} = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} \frac{\partial L}{\partial y(?, ?)} \frac{\partial y(?, ?)}{\partial x(r', c')}$$



$(r' - k_1 + 1, c' - k_2 + 1)$



Summing Gradient Contributions

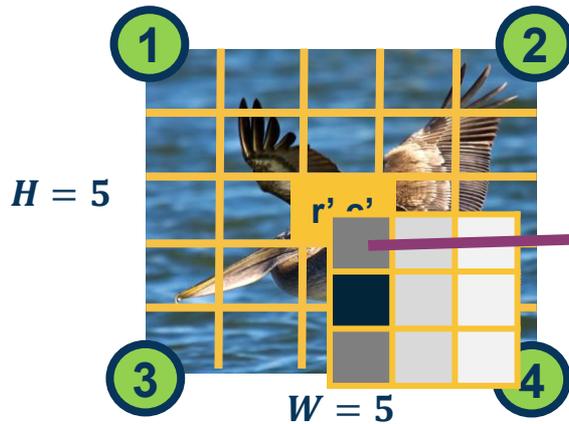
Chain rule for affected pixels (sum gradients):

$$\frac{\partial L}{\partial x(r', c')} = \sum_{\text{Pixels } p} \frac{\partial L}{\partial y(p)} \frac{\partial y(p)}{\partial x(r', c')}$$

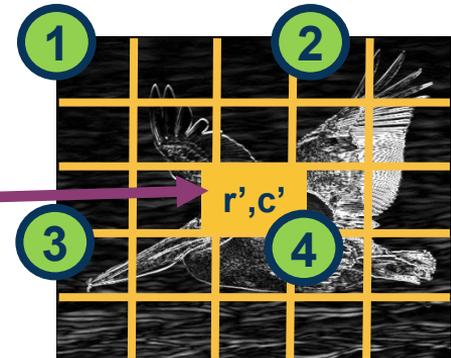
$$\frac{\partial L}{\partial x(r', c')} = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} \frac{\partial L}{\partial y(?, ?)} \frac{\partial y(?, ?)}{\partial x(r', c')}$$

$$x(r', c') * k(0, 0) \Rightarrow y(r', c')$$

$$x(r', c') * k(1, 1) \Rightarrow ?$$



$$(r' - k_1 + 1, c' - k_2 + 1)$$



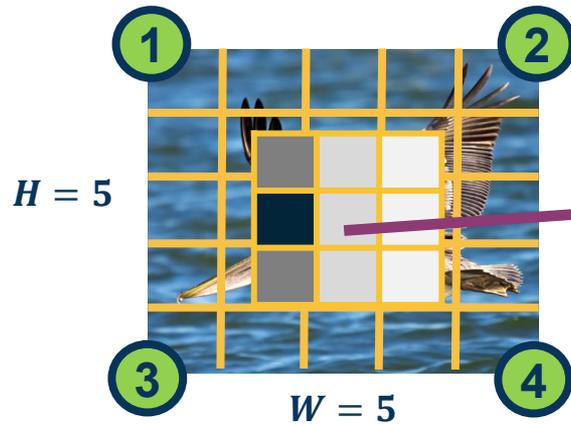
Summing Gradient Contributions

Chain rule for affected pixels (sum gradients):

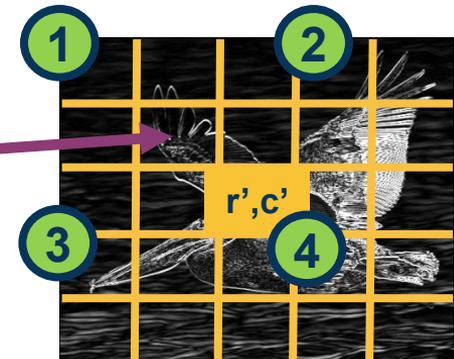
$$\frac{\partial L}{\partial x(r', c')} = \sum_{\text{Pixels } p} \frac{\partial L}{\partial y(p)} \frac{\partial y(p)}{\partial x(r', c')}$$

$$\frac{\partial L}{\partial x(r', c')} = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} \frac{\partial L}{\partial y(?, ?)} \frac{\partial y(?, ?)}{\partial x(r', c')}$$

$$\begin{aligned} x(r', c') * k(0, 0) &\Rightarrow y(r', c') \\ x(r', c') * k(1, 1) &\Rightarrow y(r' - 1, c' - 1) \\ \dots \\ x(r', c') * k(a, b) &\Rightarrow y(r' - a, c' - b) \end{aligned}$$



$(r' - k_1 + 1, c' - k_2 + 1)$



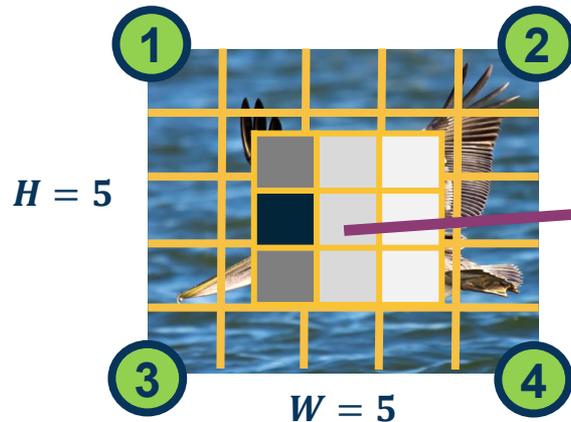
Summing Gradient Contributions

Chain rule for affected pixels (sum gradients):

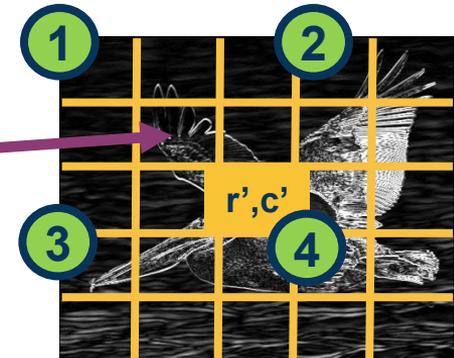
$$\frac{\partial L}{\partial x(r', c')} = \sum_{\text{Pixels } p} \frac{\partial L}{\partial y(p)} \frac{\partial y(p)}{\partial x(r', c')}$$

$$\frac{\partial L}{\partial x(r', c')} = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} \frac{\partial L}{\partial y(r' - a, c' - b)} \frac{\partial y(r' - a, c' - b)}{\partial x(r', c')}$$

Let's derive it analytically this time (as opposed to visually)



$(r' - k_1 + 1, c' - k_2 + 1)$



Summing Gradient Contributions

Definition of cross-correlation (use a', b' to distinguish from prior variables):

$$y(\mathbf{r}', \mathbf{c}') = (\mathbf{x} * \mathbf{k})(\mathbf{r}', \mathbf{c}') = \sum_{a'=0}^{k_1-1} \sum_{b'=0}^{k_2-1} x(\mathbf{r}' + \mathbf{a}', \mathbf{c}' + \mathbf{b}') k(\mathbf{a}', \mathbf{b}')$$

Plug in what we actually wanted :

$$y(\mathbf{r}' - \mathbf{a}, \mathbf{c}' - \mathbf{b}) = (\mathbf{x} * \mathbf{k})(\mathbf{r}', \mathbf{c}') = \sum_{a'=0}^{k_1-1} \sum_{b'=0}^{k_2-1} x(\mathbf{r}' - \mathbf{a} + \mathbf{a}', \mathbf{c}' - \mathbf{b} + \mathbf{b}') k(\mathbf{a}', \mathbf{b}')$$

What is $\frac{\partial y(\mathbf{r}' - \mathbf{a}, \mathbf{c}' - \mathbf{b})}{\partial x(\mathbf{r}', \mathbf{c}')} = \mathbf{k}(\mathbf{a}, \mathbf{b})$

(we want term with $x(\mathbf{r}', \mathbf{c}')$ in it;
this happens when $\mathbf{a}' = \mathbf{a}$ and $\mathbf{b}' = \mathbf{b}$)

Plugging in to earlier equation:

$$\begin{aligned}\frac{\partial L}{\partial x(r', c')} &= \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} \frac{\partial L}{\partial y(r' - a, c' - b)} \frac{\partial y(r' - a, c' - b)}{\partial x(r', c')} \\ &= \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} \frac{\partial L}{\partial y(r' - a, c' - b)} k(a, b)\end{aligned}$$

Again, all operations can be implemented via matrix multiplications (same as FC layer)!

Does this look familiar?

Convolution between upstream gradient and kernel!

(can implement by flipping kernel and cross-correlation)

Backwards is Convolution

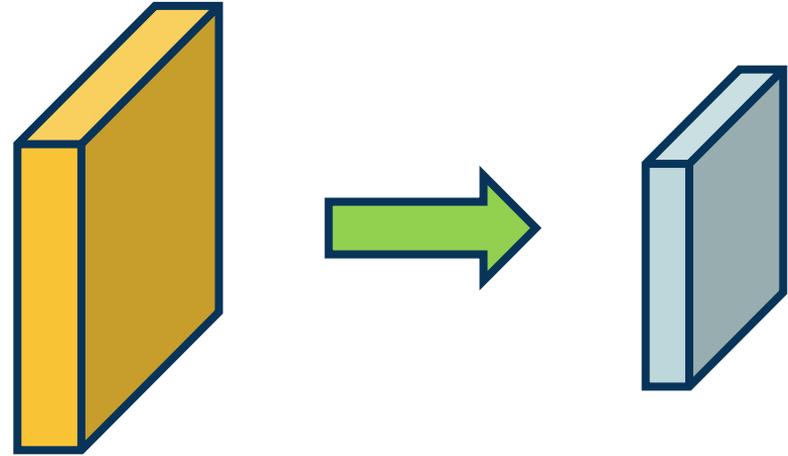
- Convolutions are mathematical descriptions of striding linear operation
- In practice, we implement **cross-correlation neural networks!** (still called convolutional neural networks due to history)
 - Can connect to convolutions via duality (flipping kernel)
 - Convolution formulation has mathematical properties explored in ECE
- Duality for forwards and backwards:
 - **Forward:** Cross-correlation
 - **Backwards w.r.t. K :** Cross-correlation b/w upstream gradient and input
 - **Backwards w.r.t. X :** Convolution b/w upstream gradient and kernel
 - In practice implement via cross-correlation and flipped kernel
- All operations still implemented via **efficient linear algebra** (e.g. matrix-matrix multiplication)

Pooling Layers

➤ **Dimensionality reduction** is an important aspect of machine learning

➤ Can we make a layer to **explicitly down-sample** image or feature maps?

➤ **Yes!** We call one class of these operations **pooling** operations



Parameters

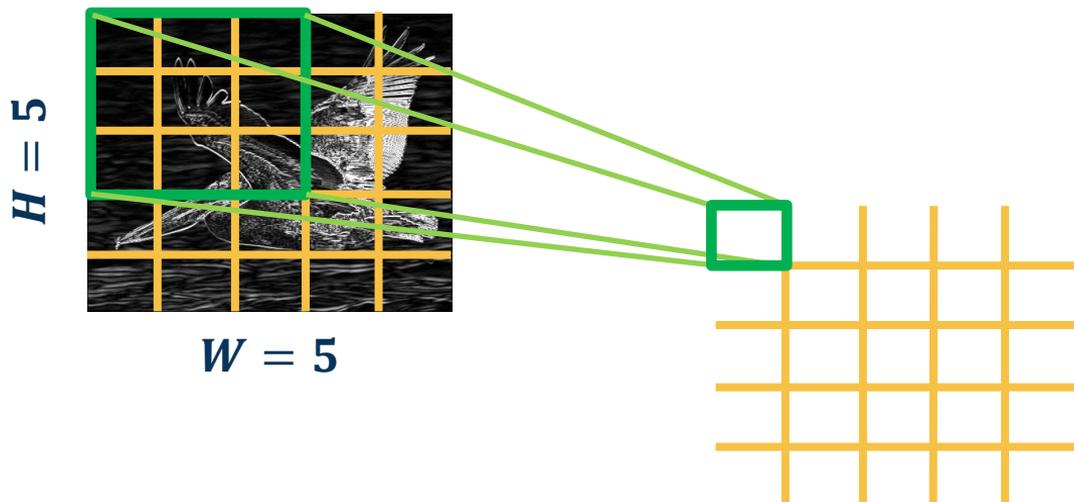
- **kernel_size** – the size of the window to take a max over
- **stride** – the stride of the window. Default value is `kernel_size`
- **padding** – implicit zero padding to be added on both sides

From: <https://pytorch.org/docs/stable/generated/torch.nn.MaxPool2d.html#torch.nn.MaxPool2d>

Example: Max pooling

- Stride window across image but perform per-patch **max operation**

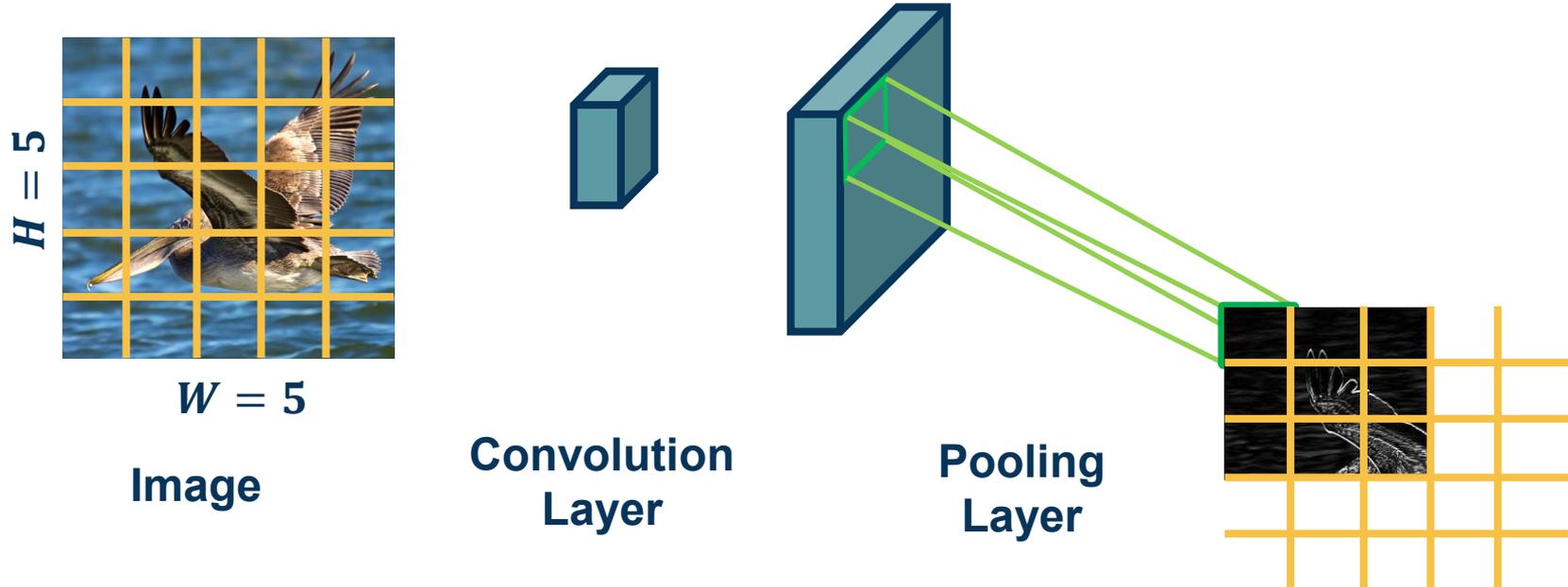
$$X(0:2, 0:2) = \begin{bmatrix} 200 & 150 & 150 \\ 100 & 50 & 100 \\ 25 & 25 & 10 \end{bmatrix} \Rightarrow \max(0:2, 0:2) = 200$$



How many learned parameters does this layer have?

None!

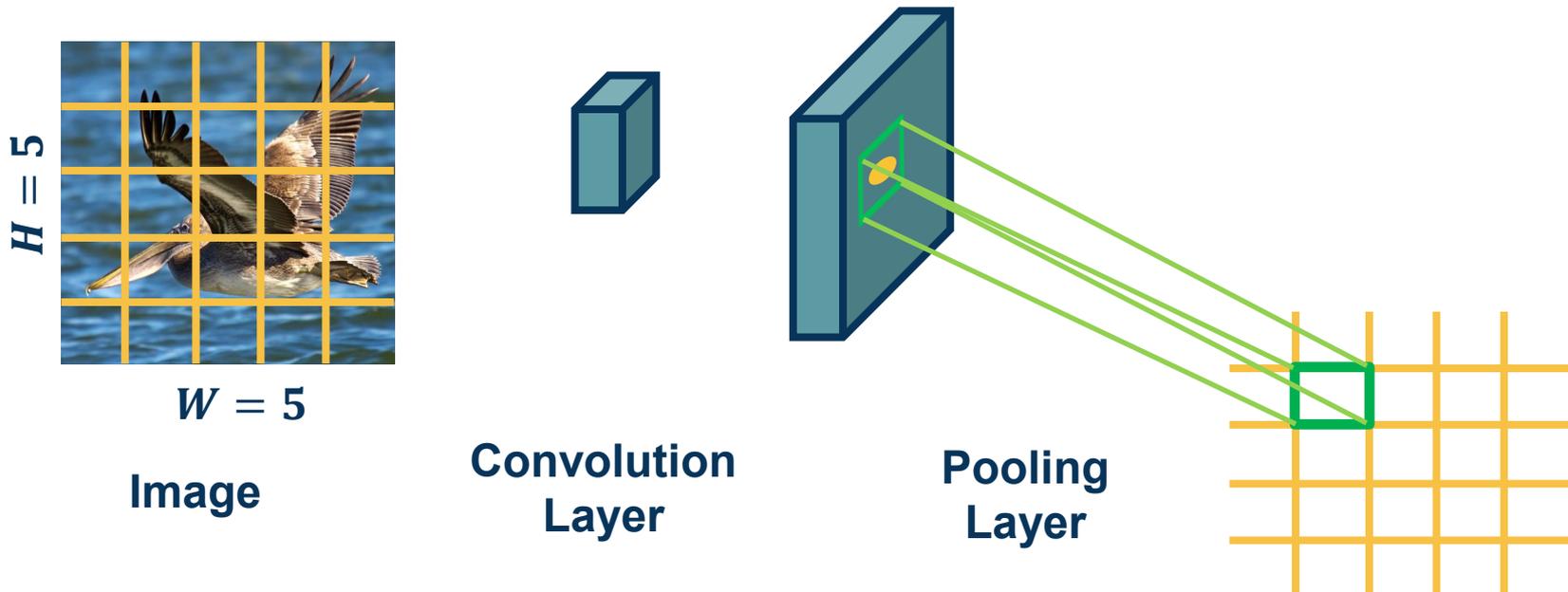
Since the **output** of convolution and pooling layers are **(multi-channel) images**, we can sequence them just as any other layer



Combining Convolution & Pooling Layers

This combination adds some **invariance** to translation of the features

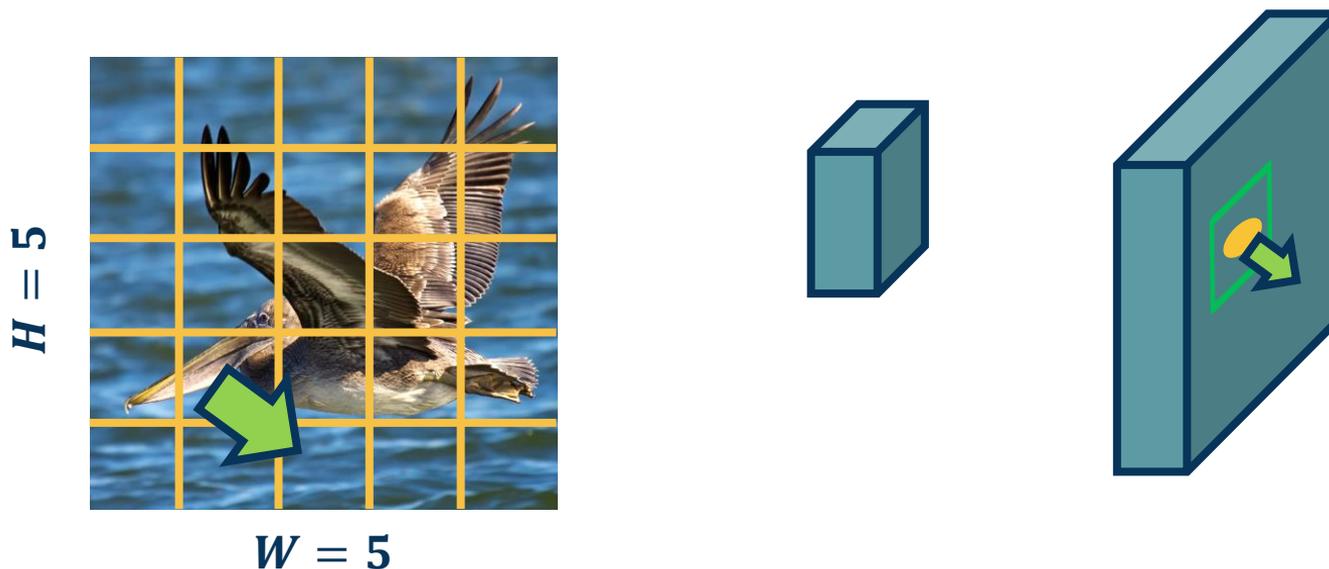
- If feature (such as beak) translated a little bit, output values still **remain the same**



Invariance

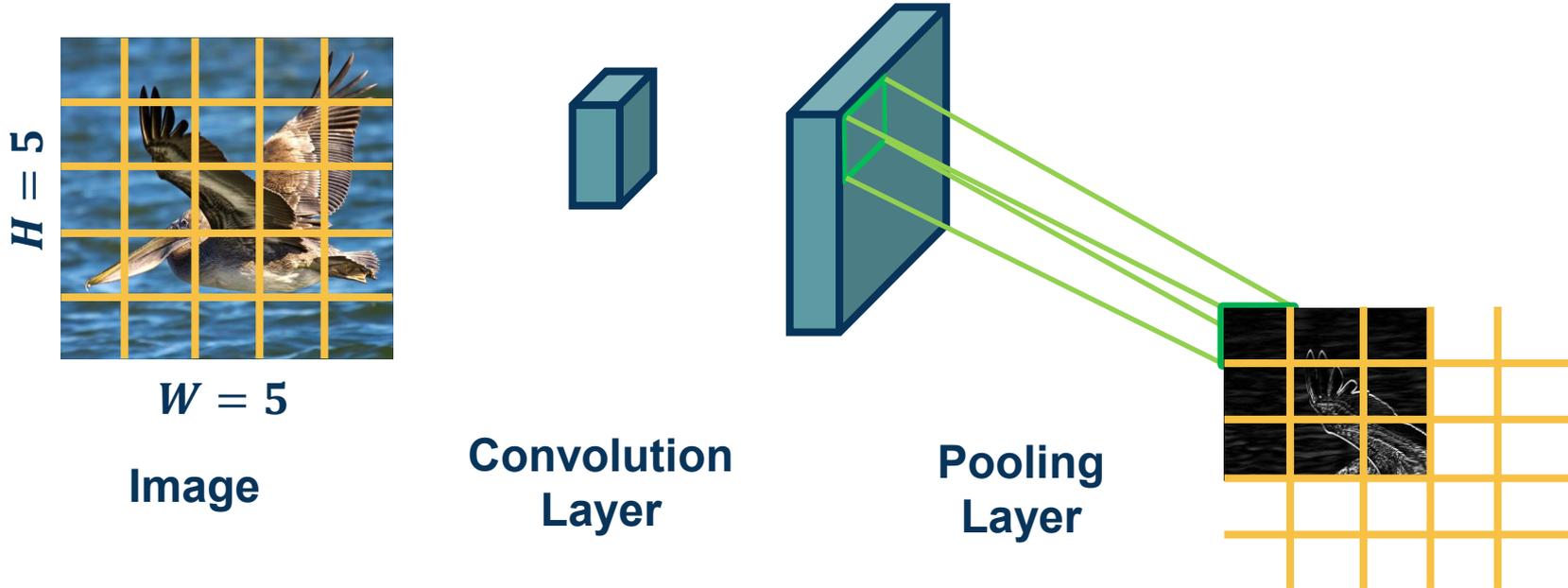
Convolution by itself has the property of **equivariance**

- ◆ If feature (such as beak) translated a little bit, output values **move by the same translation**



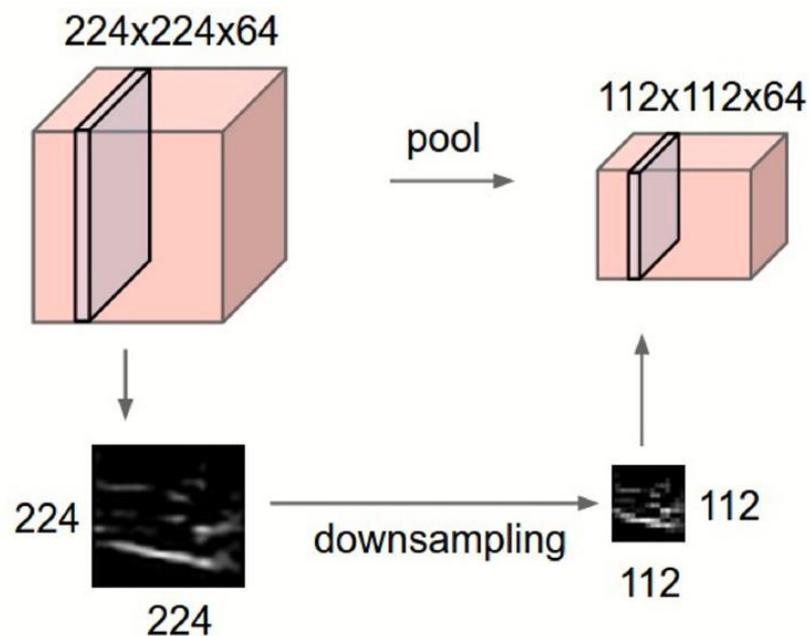
Simple Convolutional Neural Networks

Since the **output** of convolution and pooling layers are **(multi-channel) images**, we can sequence them just as any other layer



Combining Convolution & Pooling Layers

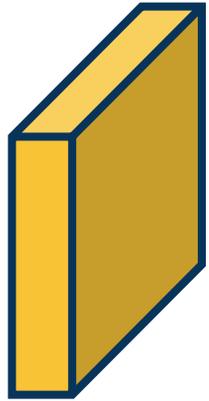
- makes the representations spatially smaller
- saves computation (GPU mem & speed), allows go deeper
- operates over each activation map independently:



From: Slides by CS 231n, Dante Xu

Pooling with Tensors

Convolutional Neural Networks (CNNs)



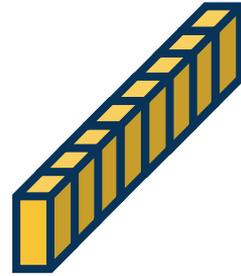
Image



Convolution +
Non-Linear
Layer



Pooling
Layer

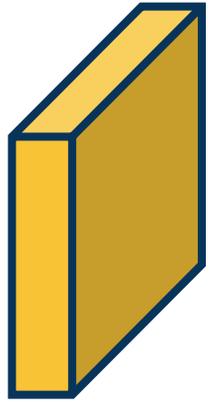


Convolution +
Non-Linear
Layer

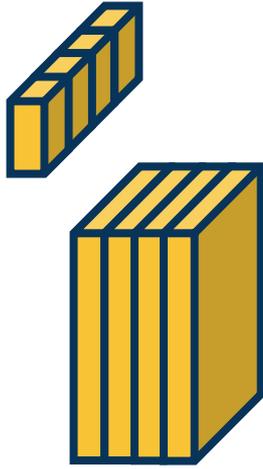


Useful,
lower-
dimensional
features

Alternating Convolution and Pooling



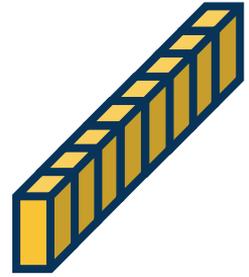
Image



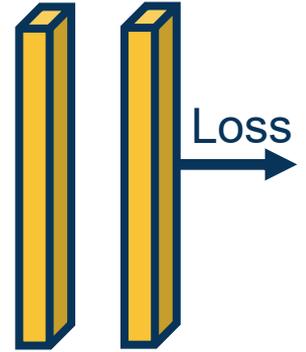
Convolution +
Non-Linear
Layer



Pooling
Layer

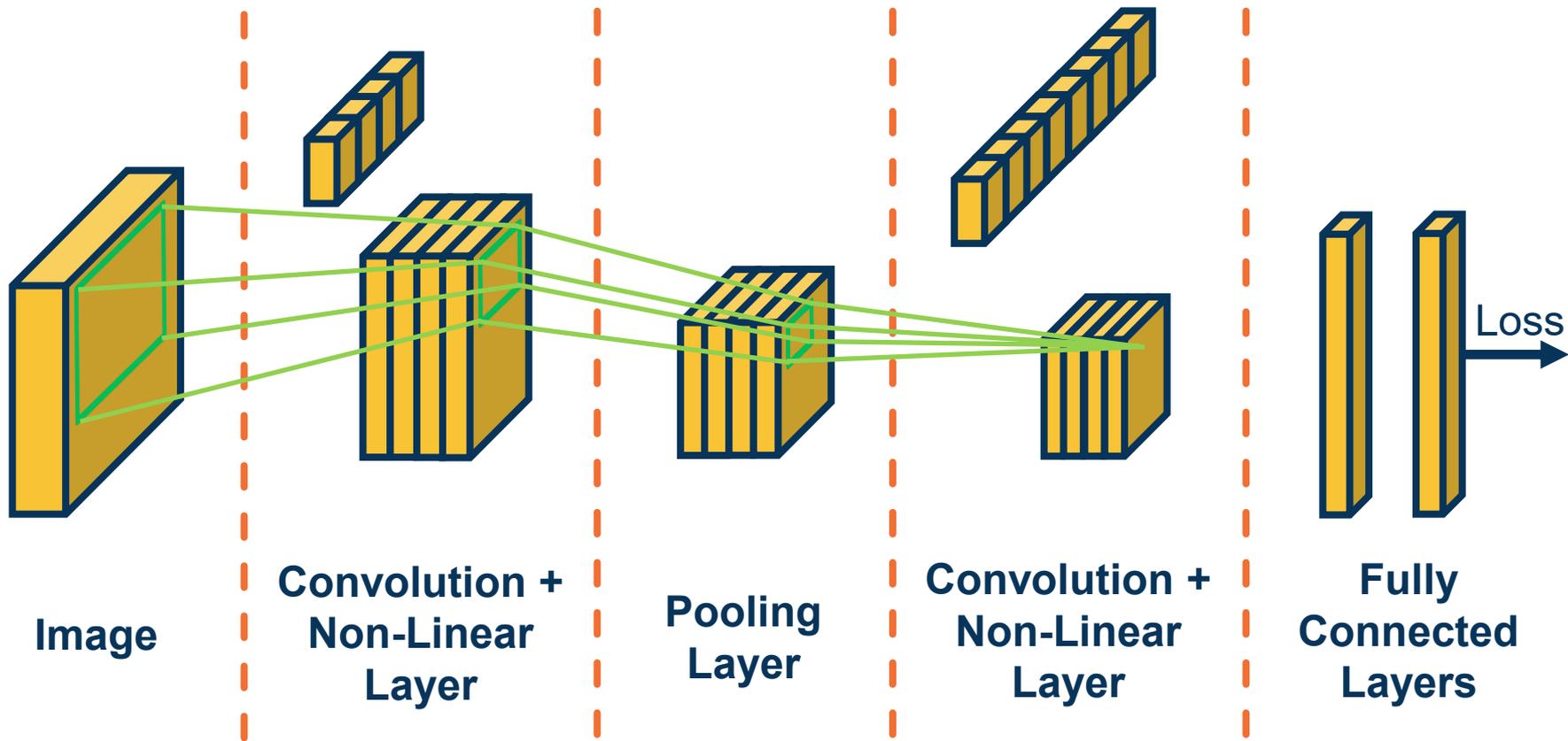


Convolution +
Non-Linear
Layer



Fully
Connected
Layers

Adding a Fully Connected Layer



Receptive Fields