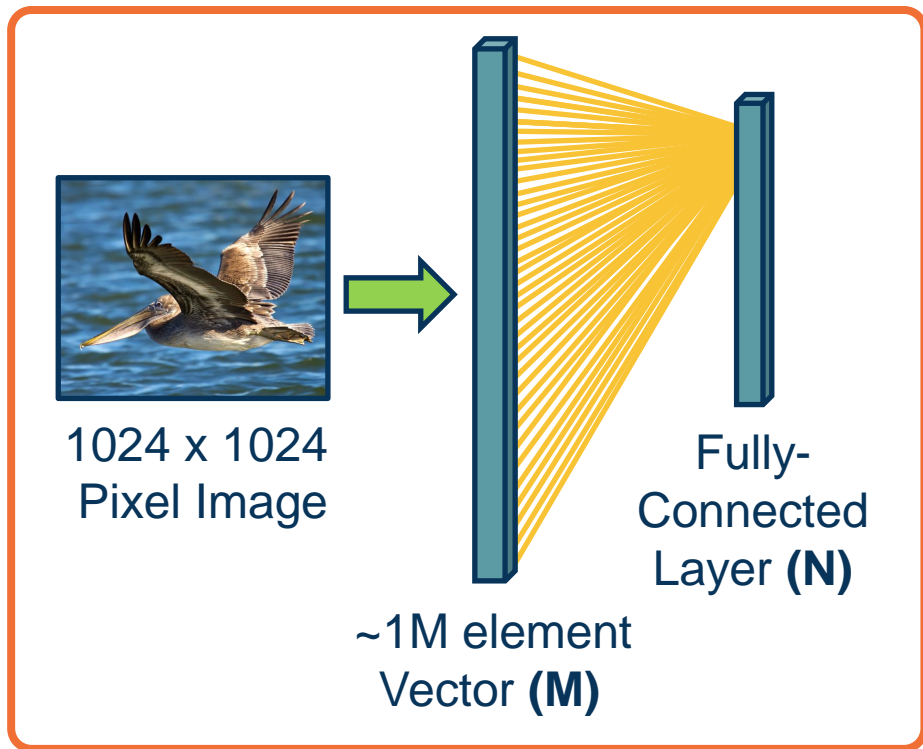Topics:

- Convolutional Neural Networks

# CS 4644-DL / 7643-A
# ZSOLT KIRA

- **Assignment 2 Due Feb 19th**
  - Implement convolutional neural networks
  - Resources (in addition to lectures):
    - [DL book: Convolutional Networks](#)
    - CNN notes https://www.cc.gatech.edu/classes/AY2022/cs7643_spring/assets/L10_cnns_notes.pdf
    - Backprop notes https://www.cc.gatech.edu/classes/AY2022/cs7643_spring/assets/L10_cnns_backprop_notes.pdf
    - HW2 Tutorial (piazza @113)
    - Slower OMSCS lectures on dropbox: Module 2 Lessons 5-6 (M2L5/M2L6) (https://www.dropbox.com/sh/iviro188gq0b4vs/AADdHxX_Uy1TkpF_yvIzX0nPa?dl=0)
- **GPU resources**
  - **For assignments, can use CPU or Google Colab**
  - **Projects:**
    - **Google Cloud Credits**
    - **PACE-ICE**

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



1024 x 1024
Pixel Image

~1M element
Vector **(M)**

Fully-
Connected
Layer **(N)**

**How many parameters?**

⬡ $M*N$ (weights) $+ N$ (bias)

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

**More parameters => More data needed**

**Is this necessary?**

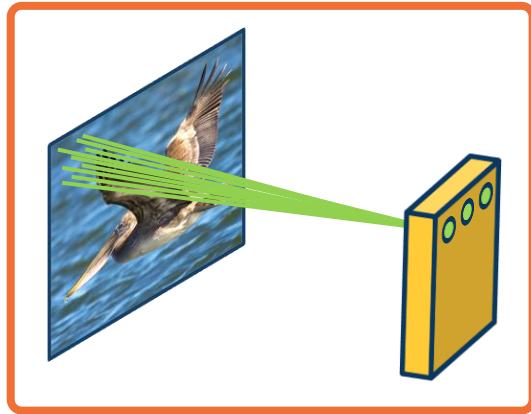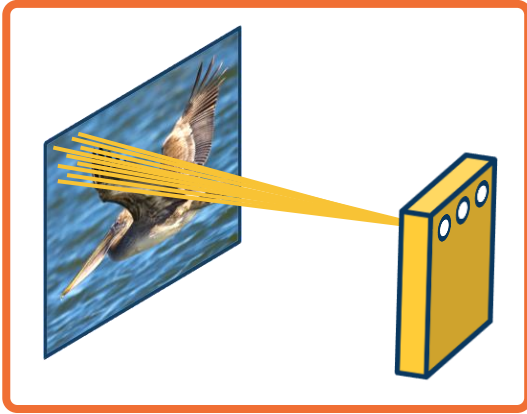# Image features are spatially localized!

- Smaller features repeated across the image
  - Edges
  - Color
  - Motifs (corners, etc.)
- No reason to believe one feature tends to appear in one location vs. another (stationarity)



Can we induce a *bias* in the design of a neural network layer to reflect this?

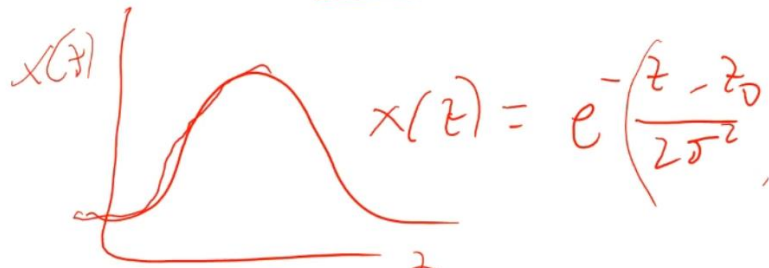Georgia Tech

We can learn **many** such features for this one layer

- Weights are **not** shared across different feature extractors

- **Parameters:** $(K_1 \times K_2 + 1) * M$ where $M$ is number of features we want to learn

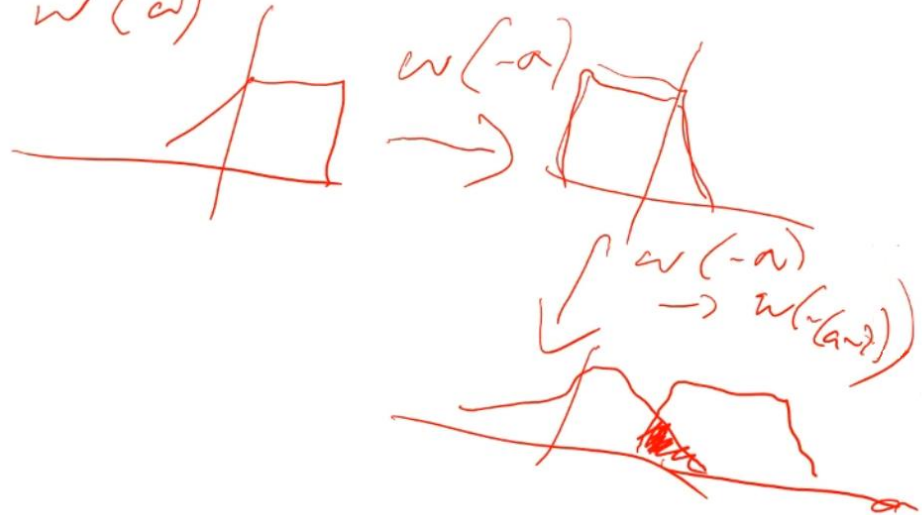# This operation is **extremely common** in electrical/computer engineering!
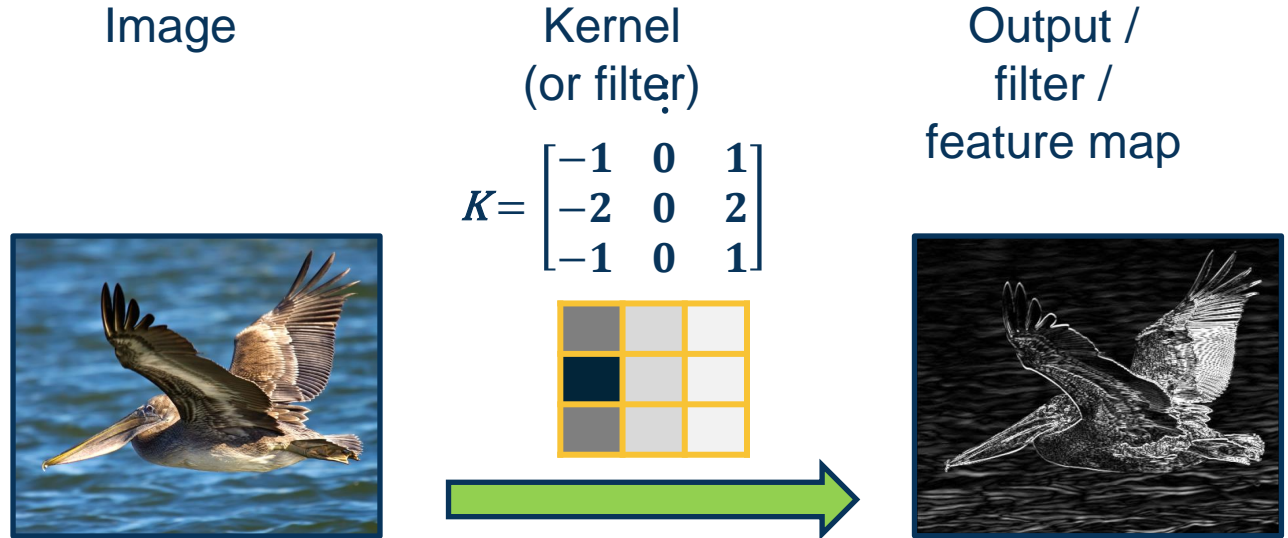


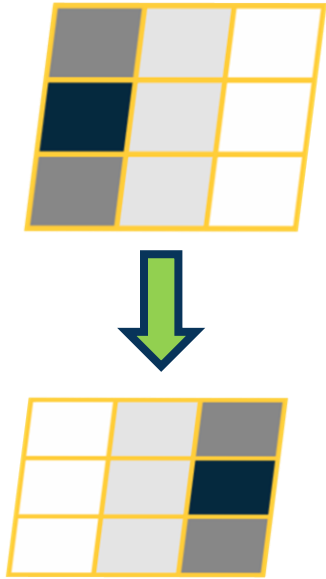*From https://en.wikipedia.org/wiki/Convolution*

We will make this convolution operation **a layer** in the neural network

- Initialize kernel values randomly and optimize them!

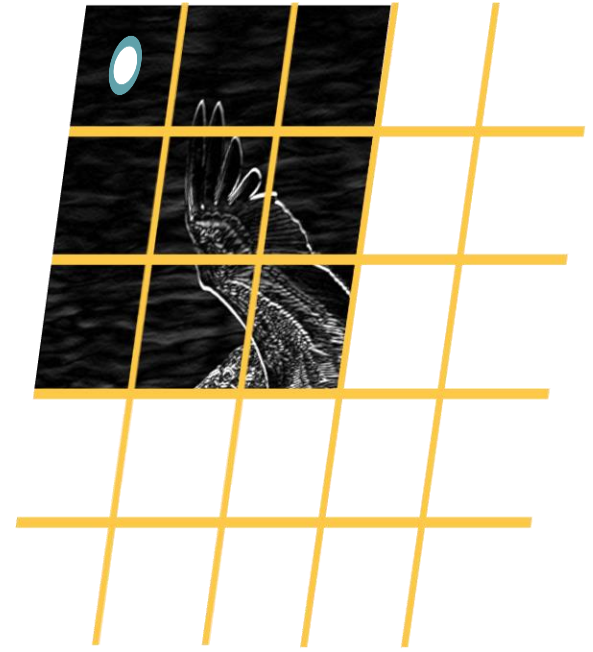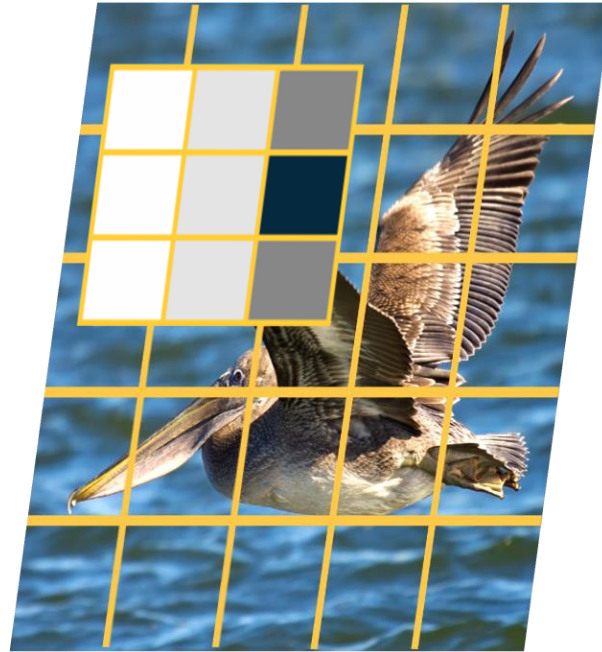- These are our parameters (plus a bias term per filter)

Image

Kernel (or filter)

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

Output / filter / feature map



**2D Convolution**

**2D Discrete Convolution**

**1. Flip kernel (rotate 180 degrees)**

**2. Stride along image**

$$y(r,c) = (x * k)(r,c) = \sum_{a=-\frac{H-1}{2}}^{\frac{H-1}{2}} \sum_{b=-\frac{W-1}{2}}^{\frac{W-1}{2}} x(a,b)\, k(r-a, c-b)$$

$$\left(-\frac{H-1}{2}, -\frac{W-1}{2}\right)$$

$H = 5$

$(0,0)$

$k_1 = 3$

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

$W = 5$

$$\left(\frac{H-1}{2}, \frac{W-1}{2}\right)$$

$$y(0,0) = x(-2,-2)k(2,2) + x(-2,-1)k(2,1) + x(-2,0)k(2,0) +$$
$$x(-2,1)k(2,-1) + x(-2,2)k(2,-2) + \ldots$$

**Mathematics of Discrete 2D Convolution**

$$y(r, c) = (x * k)(r, c) = \sum_{a=-\frac{K_1-1}{2}}^{\frac{k_1-1}{2},} \sum_{b=-\frac{k_2-1}{2}}^{\frac{k_2-1}{2},} x(r-a, c-b)\, k(a, b)$$
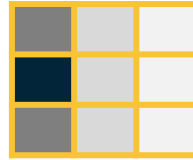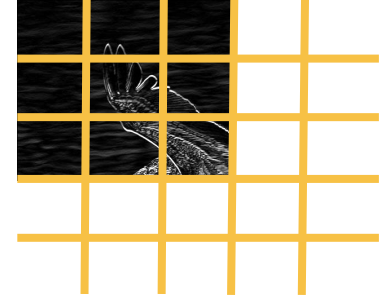
$(0, 0)$

$H = 5$

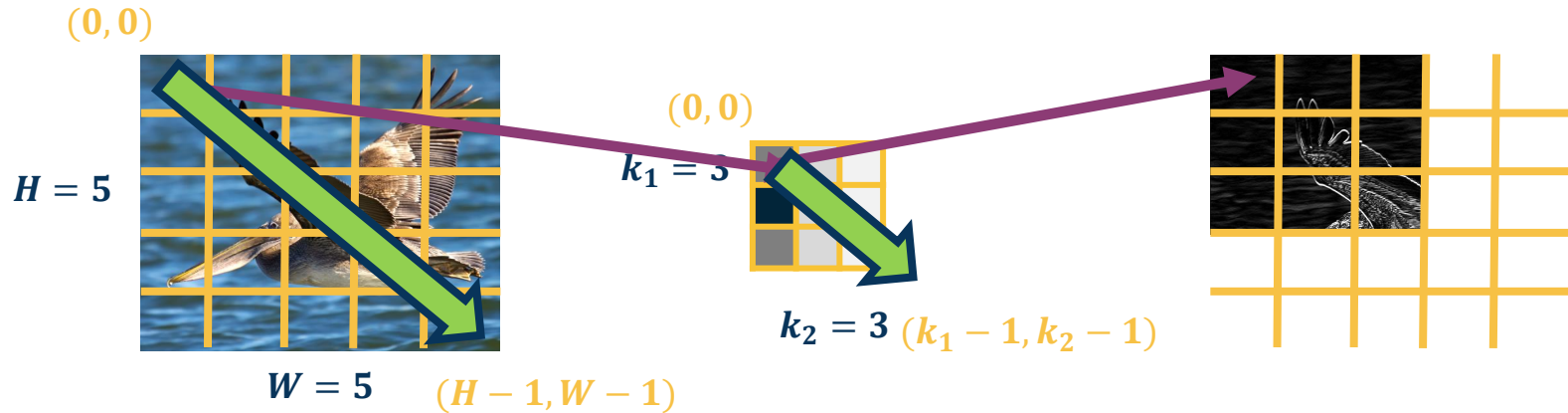$\left(-\dfrac{k_1 - 1}{2}, -\dfrac{k_2 - 1}{2}\right)$

$k_1 = 3$

$k_2 = 3$

$\left(\dfrac{k_1 - 1}{2}, \dfrac{k_2 - 1}{2}\right)$

$W = 5$

$(H - 1, W - 1)$

**Centering Around the Kernel**

Georgia Tech

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



$(0,0)$

$H = 5$

$k_1 = 3$

$(0,0)$

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

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

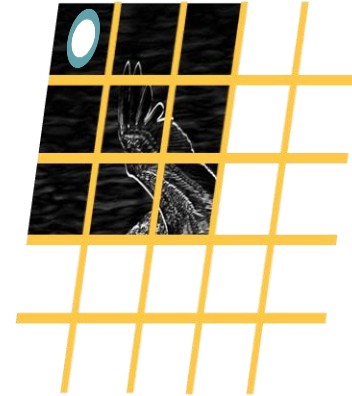**Since we will be learning these kernels, this change does not matter!**
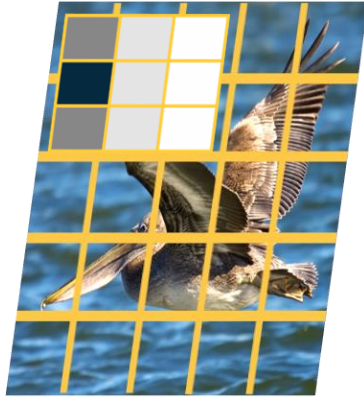
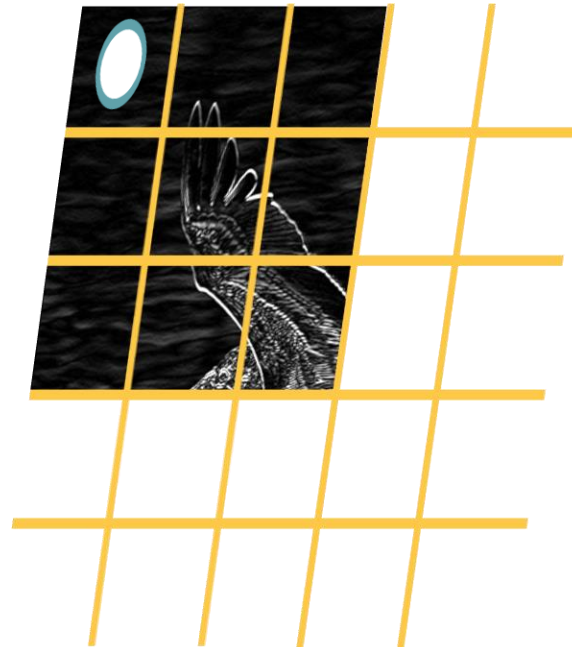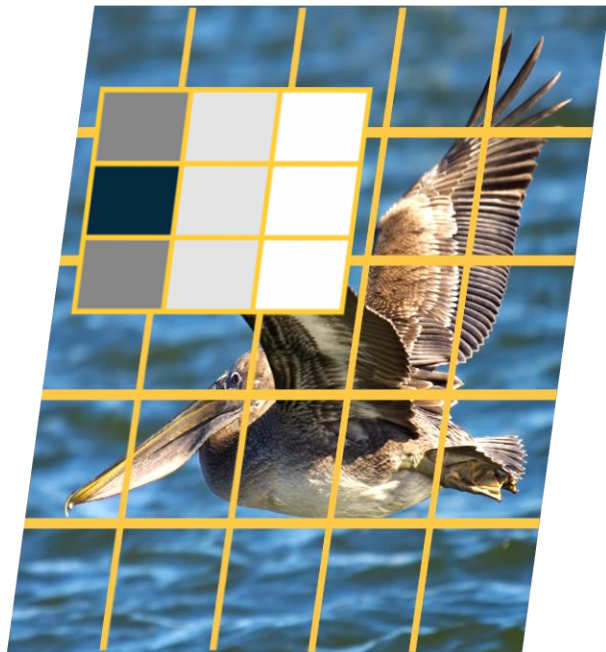$$X(0:2, 0:2) = \begin{bmatrix} 200 & 150 & 150 \\ 100 & 50 & 100 \\ 25 & 25 & 10 \end{bmatrix} \qquad K' = \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix} \qquad X(0:2,0:2) \cdot K' = 65 + \text{bias}$$
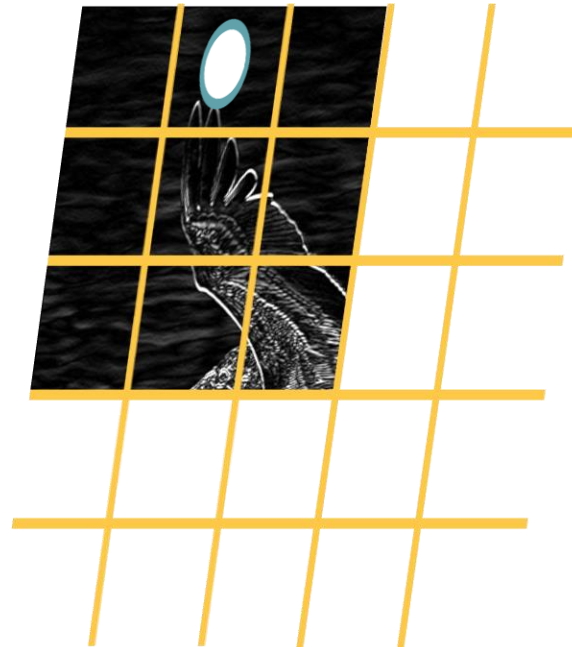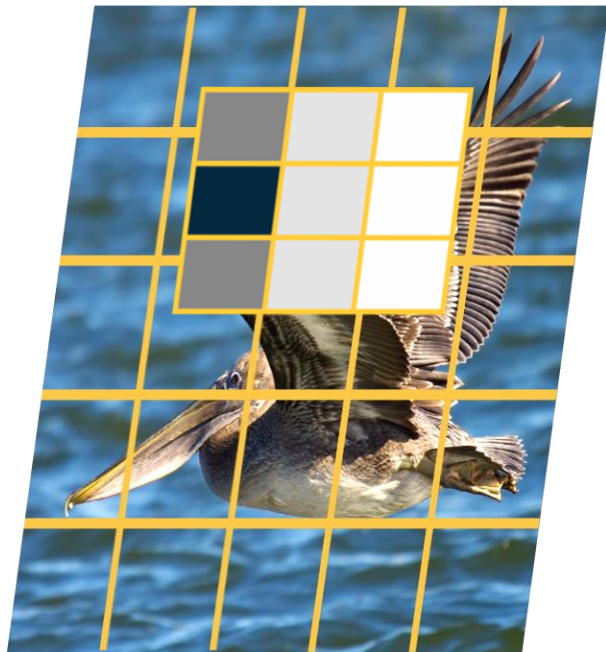
Dot product
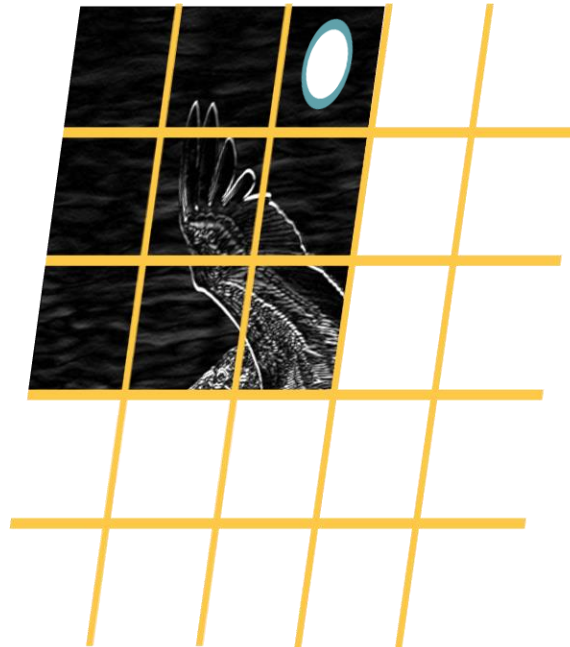(element-wise multiply and sum)

## 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

# Input & Output Sizes

# Convolution Layer Hyper-Parameters

**Parameters**

- **in_channels** (*int*) – Number of channels in the input image
- **out_channels** (*int*) – Number of channels produced by the convolution
- **kernel_size** (*int* or *tuple*) – Size of the convolving kernel
- **stride** (*int* or *tuple*, *optional*) – Stride of the convolution. Default: 1
- **padding** (*int* or *tuple*, *optional*) – Zero-padding added to both sides of the input. Default: 0
- **padding_mode** (*string*, *optional*) – `'zeros'`, `'reflect'`, `'replicate'` or `'circular'`. Default: `'zeros'`

Convolution operations have several hyper-parameters

**Output size** of vanilla convolution operation is $(H - k_1 + 1) \times (W - k_2 + 1)$

⬡ This is called a **"valid" convolution** and only applies kernel within image



$(0, 0)$
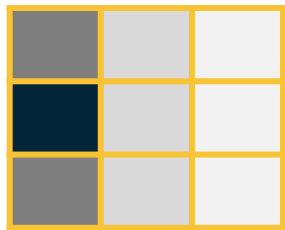
$H = 5$

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

$(0, 0)$

$k_1 = 3$

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

$H - k_1 + 1$
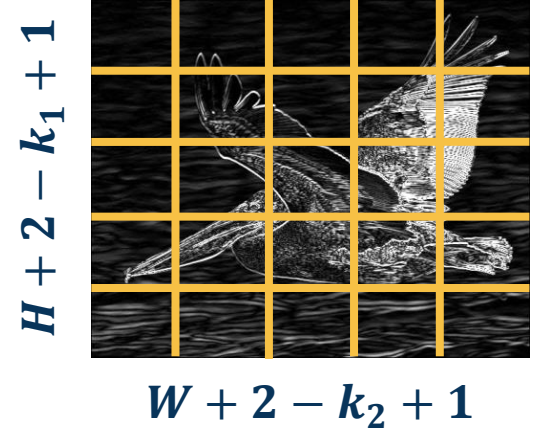
$W - k_2 + 1$

Georgia Tech

We can **pad the images** to make the output the same size:

- Zeros, mirrored image, etc.

- Note padding often refers to pixels added to **one size** ($P = 1$ here)



$H + 2$

$W + 2$

$k_1$

$k_2$

$H + 2 - k_1 + 1$

$W + 2 - k_2 + 1$

Georgia Tech

We can move the filter along the image using larger steps **(stride)**

⬡ This can potentially result in **loss of information**

⬡ Can be used for **dimensionality reduction** (not recommended)

**Stride = 2 (every other pixel)**

$H$

$W$

$(H - k_1)/2 + 1$

$(W - k_2)/2 + 1$

**Stride**

# Stride can result in **skipped pixels**, e.g. stride of 3 for 5x5 input



$H$

$W$

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**!



$H$ $W$ 3

**Image**

$k_1$ $k_2$ 3

**Kernel**

$H - k_1 + 1$ $W - k_2 + 1$ 1

**Feature Map**

**Multi-Channel Inputs**

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**

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

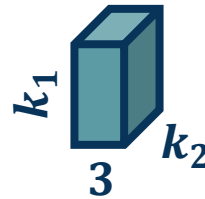**Number of channels in output is equal to _number of kernels_**



$H$

$W$

3

**Image**

$k_1$

$k_2$

3

**Kernels**

$H - k_1 + 1$

$W - k_2 + 1$

4

**Feature Maps**

Georgia Tech

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



**Image**

**Kernels**

**Feature Maps**

**Number of Parameters**

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

⬡ **Step 1:** Lay out image patches in vector form (note can overlap!)

**Input Image**



Im2col
=>

| Patch 1 |
|---|
| Patch 2 |
| ... |

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

**Vectorization**

Georgia Tech

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/*

**Vectorization**

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

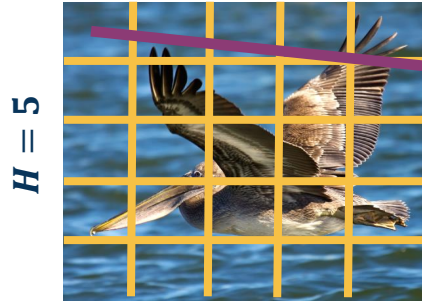$$K = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$$

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

Georgia Tech

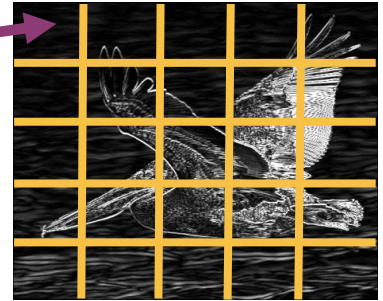$$y(r,c) = (x * k)(r,c) = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x(r+a, c+b)\, k(a,b)$$

$(0,0)$

$H = 5$

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

$(0,0)$

$k_1 = 3$

$k_2 = 3$

$(k_1 - 1, k_2 - 1)$

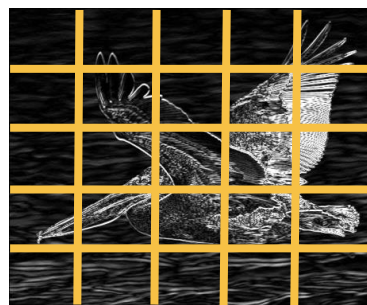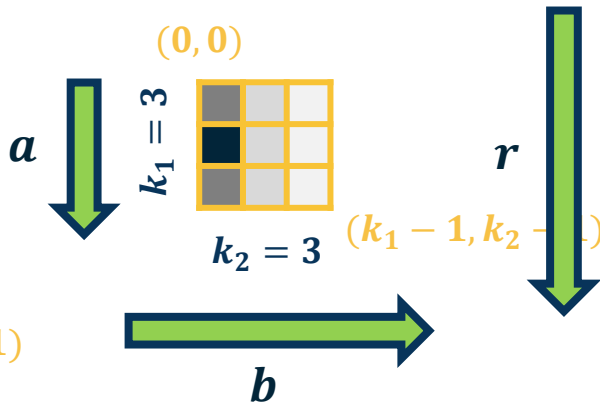**Recap: Cross-Correlation**

Georgia Tech

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



$(0,0)$

$H = 5$

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

$a$

$(0,0)$

$k_1 = 3$

$k_2 = 3$

$(k_1 - 1, k_2 - 1)$

$b$

$r$

$c$

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

Georgia Tech

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

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

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

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

Georgia Tech

$$\frac{\partial L}{\partial h^{\ell-1}} \qquad \boxed{\phantom{XXXXXXXX}} \qquad \frac{\partial L}{\partial h^{\ell}}$$

$$\frac{\partial L}{\partial k}$$

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

**Gradient for passing back**

**Gradient for weight update**

(weights = k, i.e. kernel values)
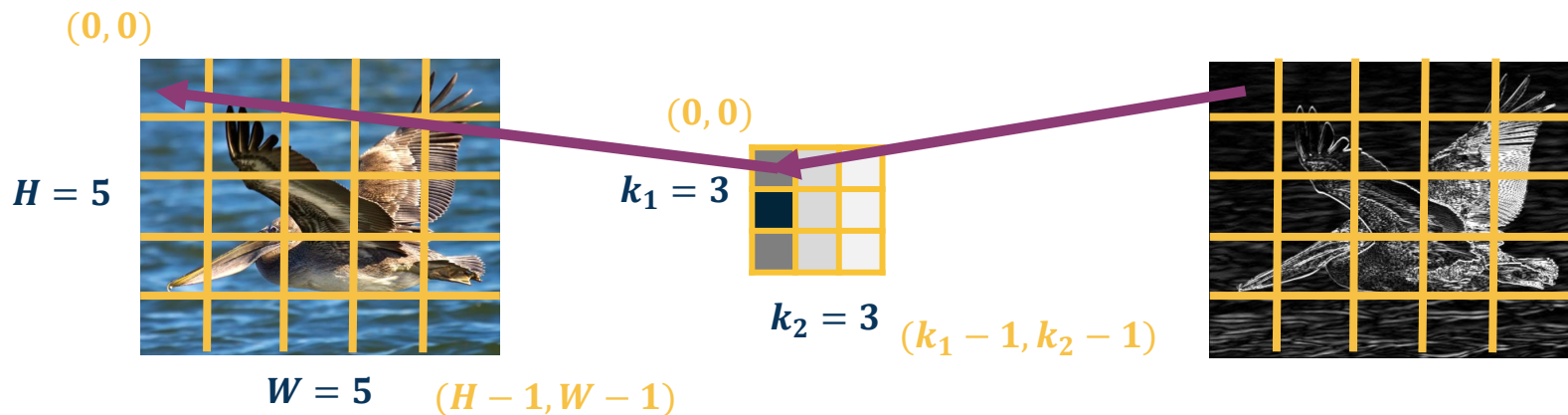
**Backpropagation Chain Rule**

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

**Gradient for weight update**

Calculate one pixel at a time $\dfrac{\partial L}{\partial k(a, b)}$



What does this weight affect at the output?

Everything!

$(0, 0)$

$H = 5$

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

$(0, 0)$

$k_1 = 3$

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

Need to incorporate all upstream gradients:

$$\left\{ \frac{\partial L}{\partial y(0,0)}, \frac{\partial L}{\partial y(0,1)}, \cdots, \frac{\partial L}{\partial y(H,W)} \right\}$$
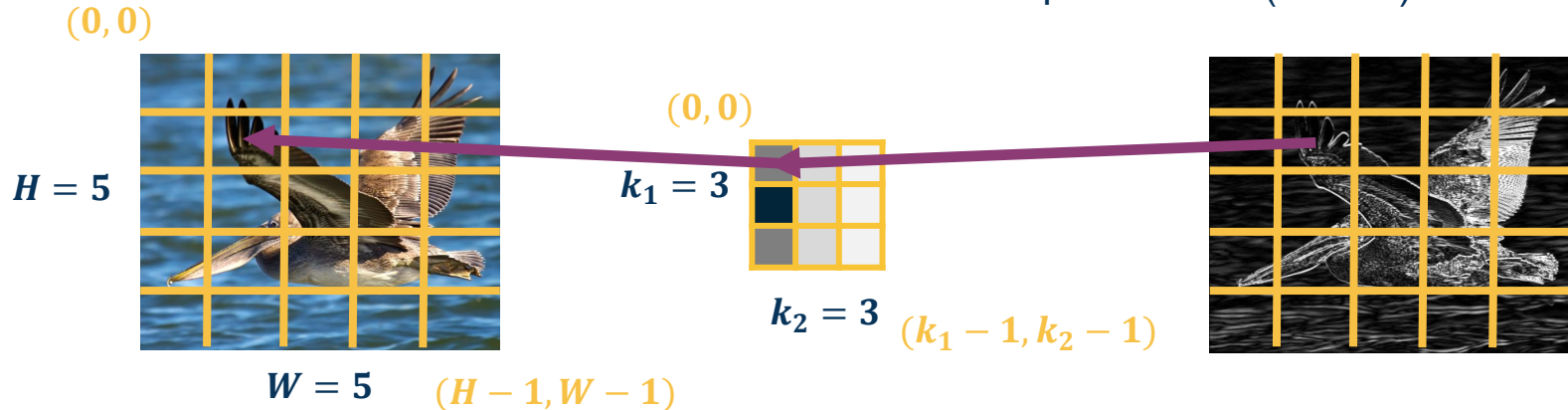
Chain Rule:

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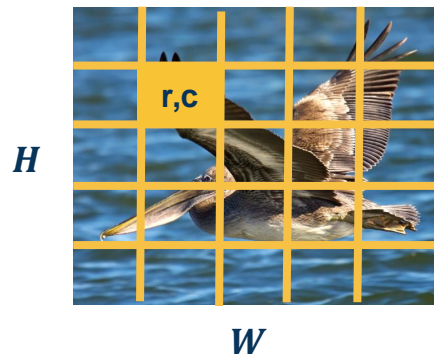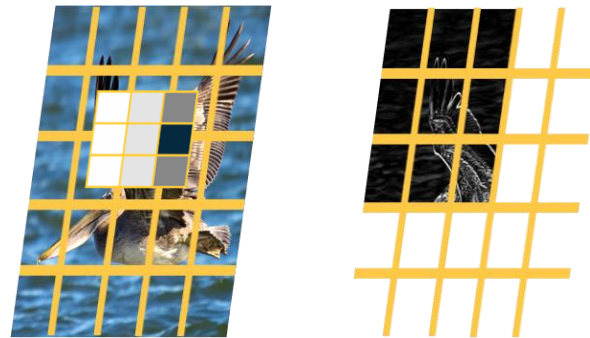
Sum over all output pixels

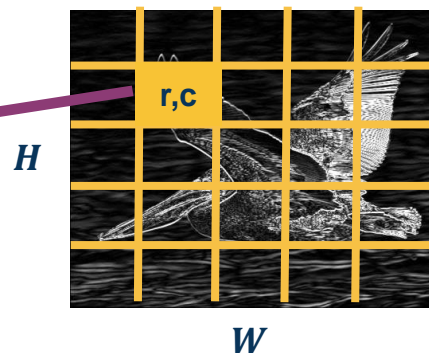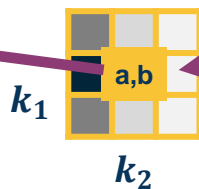Upstream gradient (known)

We will compute

$(0,0)$

$H = 5$

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

$(0,0)$

$k_1 = 3$

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

**Chain Rule over all Output Pixels**

Georgia Tech

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

r,c

H

W

?

$k_1$

a,b

$k_2$

r,c
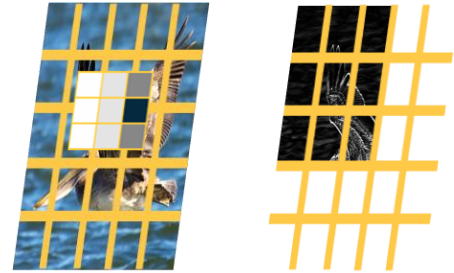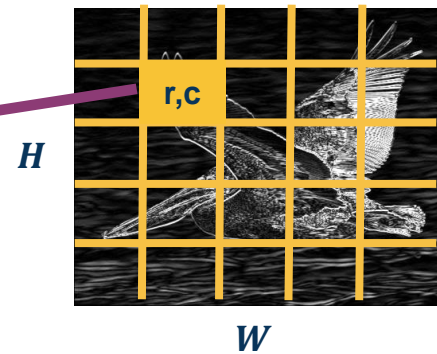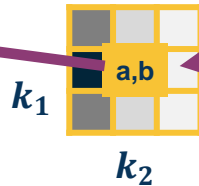
H

W

Georgia Tech

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

**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 $k(a,b)$).
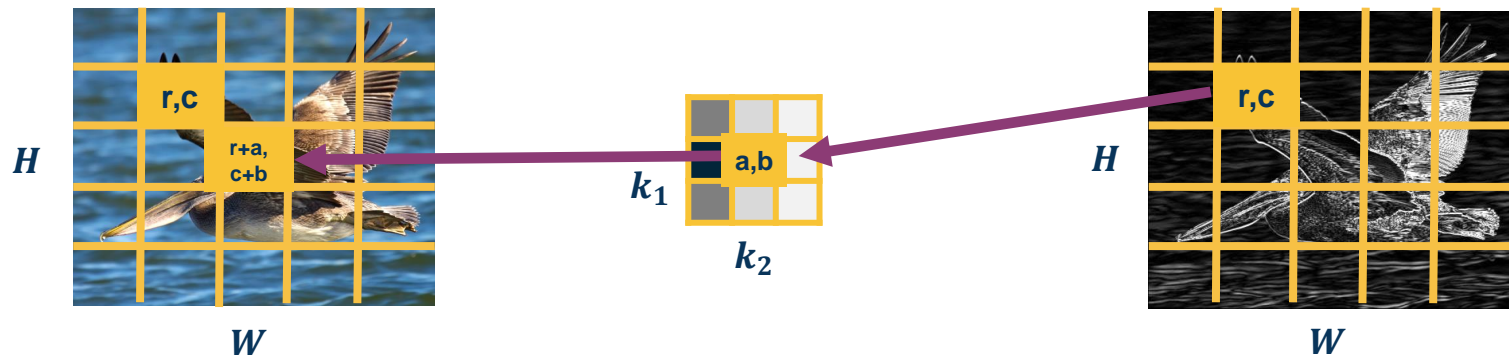


$H$

$W$

**?**

$k_1$

a,b

$k_2$

$H$

$W$

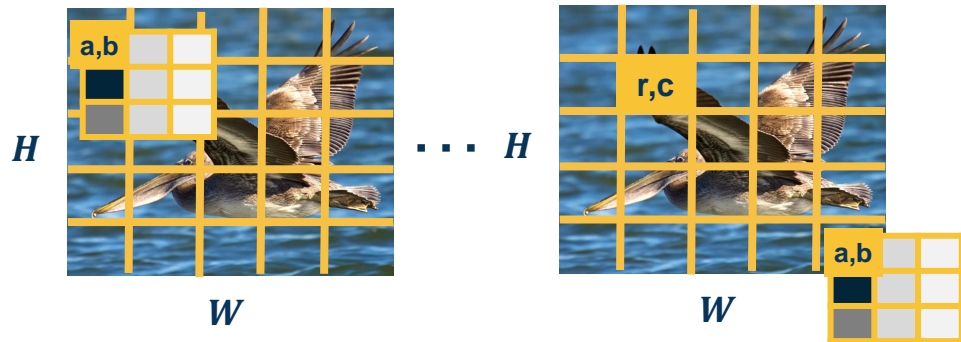$$\frac{\partial y(r,c)}{\partial k(a,b)} = x(r+a, c+b)$$

**Does this look familiar?**

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

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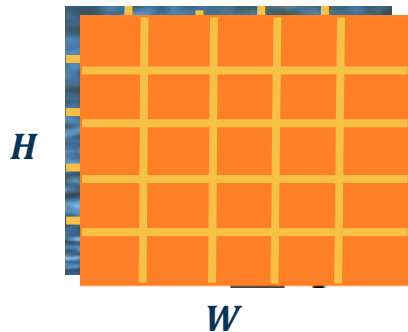
(until $k_1 \times k_2$ output)



**Gradients and Cross-Correlation**

**Forward Pass**

a,b

$H$

$W$

$\cdots$  $H$

r,c

a,b

$W$

**Backward Pass** k$(0,0)$

**Backward Pass** $k(2,2)$

$H$

$W$

r,c

$H$

$W$

**Does this look familiar?**

**Cross-correlation between upstream gradient and input!**

(until $k_1 \times k_2$ output)

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

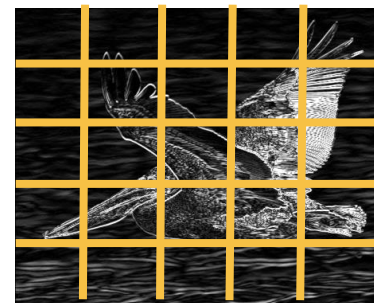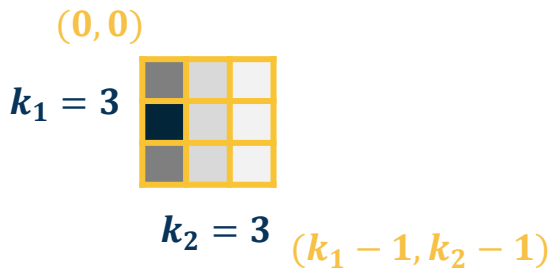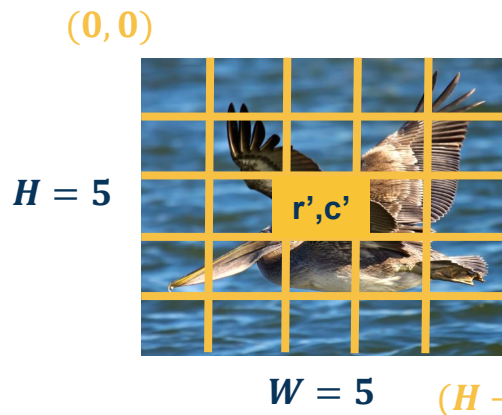**Forward and Backward Duality**

Georgia Tech

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

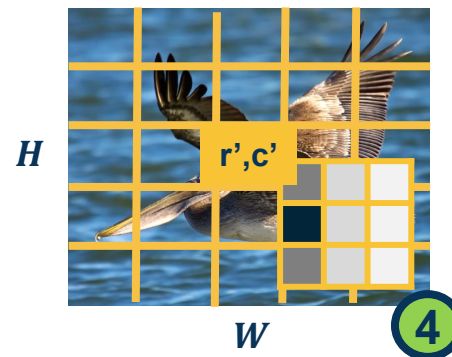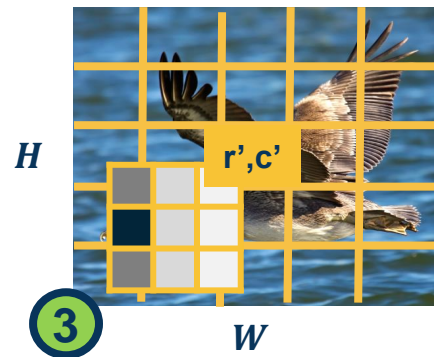Gradient for input (to pass to prior layer)

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

**What does this input pixel affect at the output?**

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

$(0, 0)$

$H = 5$

r',c'

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

$(0, 0)$

$k_1 = 3$

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

Georgia Tech

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

$H = 5$

$W = 5$

r',c'

$k_1 = 3$

$k_2 = 3$

r',c'

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

**Extents at the Output**

Chain rule for affected pixels (sum gradients):

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

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



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

H = 5

W = 5

r',c'

r',c'

**Summing Gradient Contributions**

Chain rule for affected pixels (sum gradients):

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

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

$x(r', c') * k(0,0) \Rightarrow y(r', c')$
$x(r', c') * k(1,1) \Rightarrow ?$

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



$H = 5$

$W = 5$

r',c'

r',c'

**Summing Gradient Contributions**

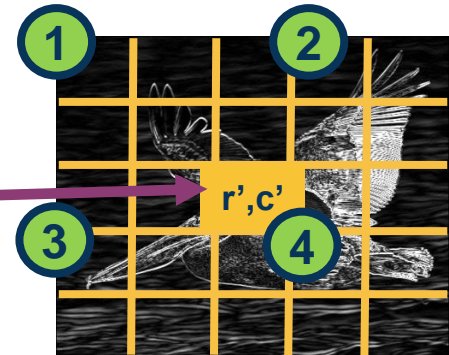# Chain rule for affected pixels (sum gradients):

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

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

$$x(r', c') * k(0, 0) \Rightarrow y(r', c')$$
$$x(r', c') * k(1, 1) \Rightarrow y(r'-1, c'-1)$$
$$\dots$$
$$x(r', c') * k(a, b) \Rightarrow y(r'-a, c'-b)$$

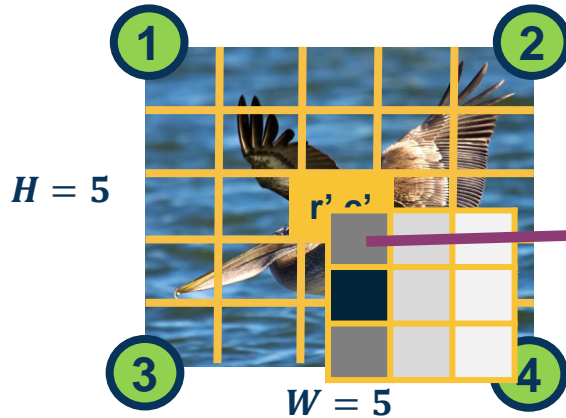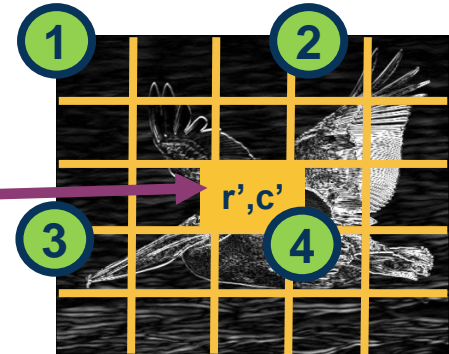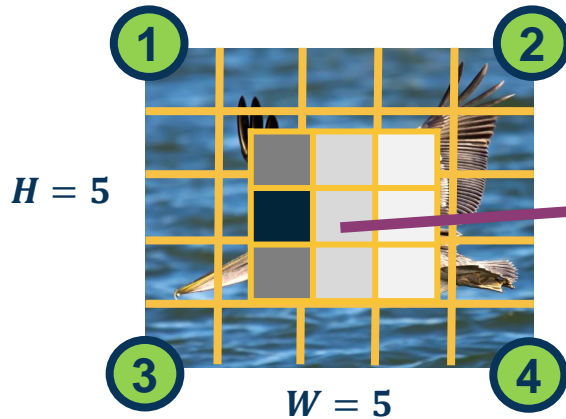$$(r' - k_1 + 1, c' - k_2 + 1)$$



$H = 5$

$W = 5$

r',c'

Georgia Tech

Chain rule for affected pixels (sum gradients):
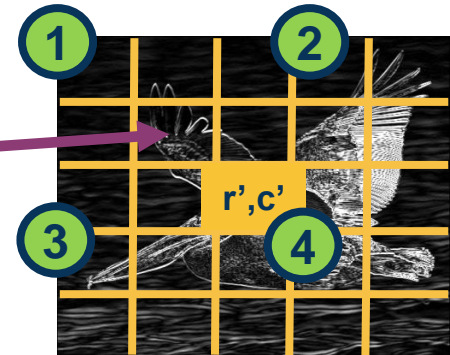
$$\frac{\partial L}{\partial x(r', c')} = \sum_{Pixels\ p} \frac{\partial L}{\partial y(p)} \frac{\partial y(p)}{\partial x(r', c')}$$

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

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

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



$H = 5$

$W = 5$

r',c'

Georgia Tech

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

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

Plug in what we actually wanted :

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

What is $\dfrac{\partial y(r' - a, c' - b)}{\partial x(r', c')} = k(a, b)$     (we want term with $x(r', c')$ in it; this happens when $a' = a$ and $b' = b$)

Georgia Tech

Plugging in to earlier equation:

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

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

**Does this look familiar?**

**Convolution between upstream gradient and kernel!**

**(can implement by flipping kernel and cross- correlation)**

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

**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)

Georgia
Tech