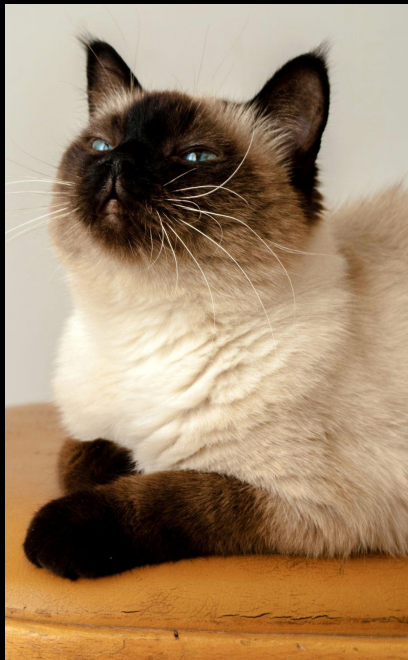


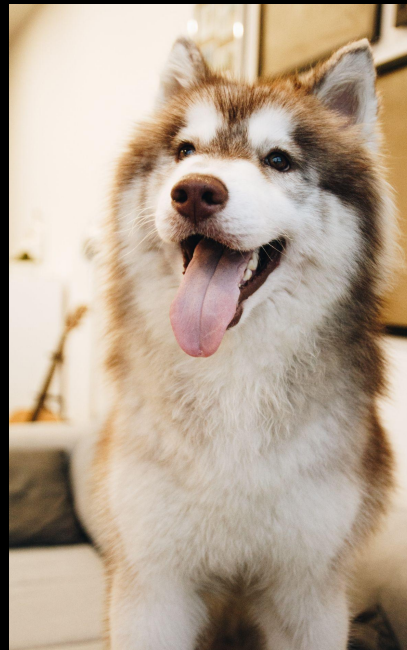
# Image Classification

Nadia Ahmed

# Computer Vision Task: **Image Classification**



“Cat”



“Dog”

# Image Classification:

- **Input:** image
- **Output:** single category for the whole image
- **Why is this interesting?**
  - Useful for image search and other applications.
  - Simple input/output specification, can resize images to the same dimensions, standard loss functions.
  - Requires reasoning about what is happening in the image-analysis and understanding of the image.
  - Useful high level proxy task to learn good representations which transfer to new tasks and domains.
  - Given a big dataset we can train a representation and take that model and fine tune it on another task and we get a way with fewer data on the target task b/c we have pretrained on a much larger dataset

# Feature Learning vs Feature Engineering



- **Edge detector**: early attempts of hand-engineering object detection algorithms failed as object shapes and appearances are too hard to describe
- Need ML and **labeled datasets** to learn relationships

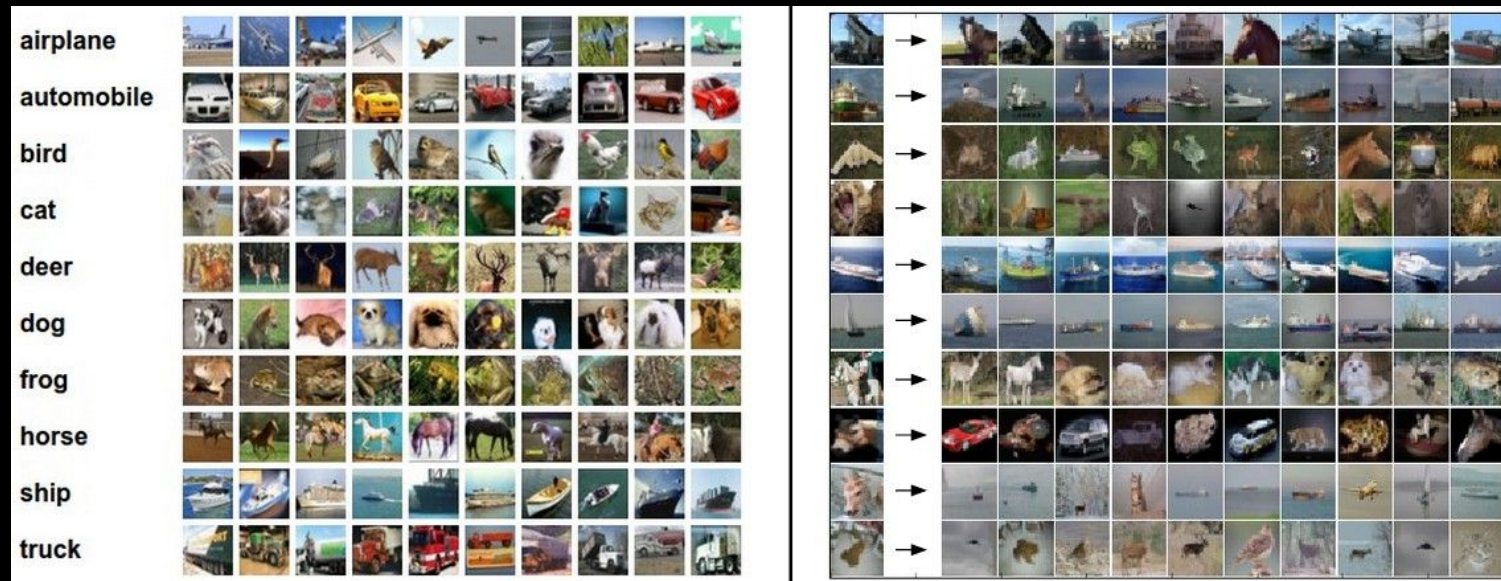


\*source: [introtodeeplearning](http://introtodeeplearning.com)

# The Data Driven Approach

- Popular Classification Single Label Datasets
  - **MNIST**: 60k training w/ labels 10k test 28x28, 10 categories
  - **Caltech101**: 101 categories
  - **ImageNet**: 15,000,000 images w/ 22,000 categories
  - **CIFAR-10**: 60,000 images 10 categories
- An Aside and Preview of Coming Attractions: **Transfer Learning**
  - Models **pre-trained** on large datasets **generalize**
  - Can get CNNs **pre-trained** on **generic tasks** (like ImageNet)
  - Can fine-tune (**retrain**) only the **last layers** on little data of a new task which leads to state of the art performance
  - **Usual** standard approach today-can download model from internet

# Simple Models



- Nearest Neighbor Classifier (find the image distance)

$$d(I_1, I_2) = \sum_p |I_1(p) - I_2(p)|_1$$

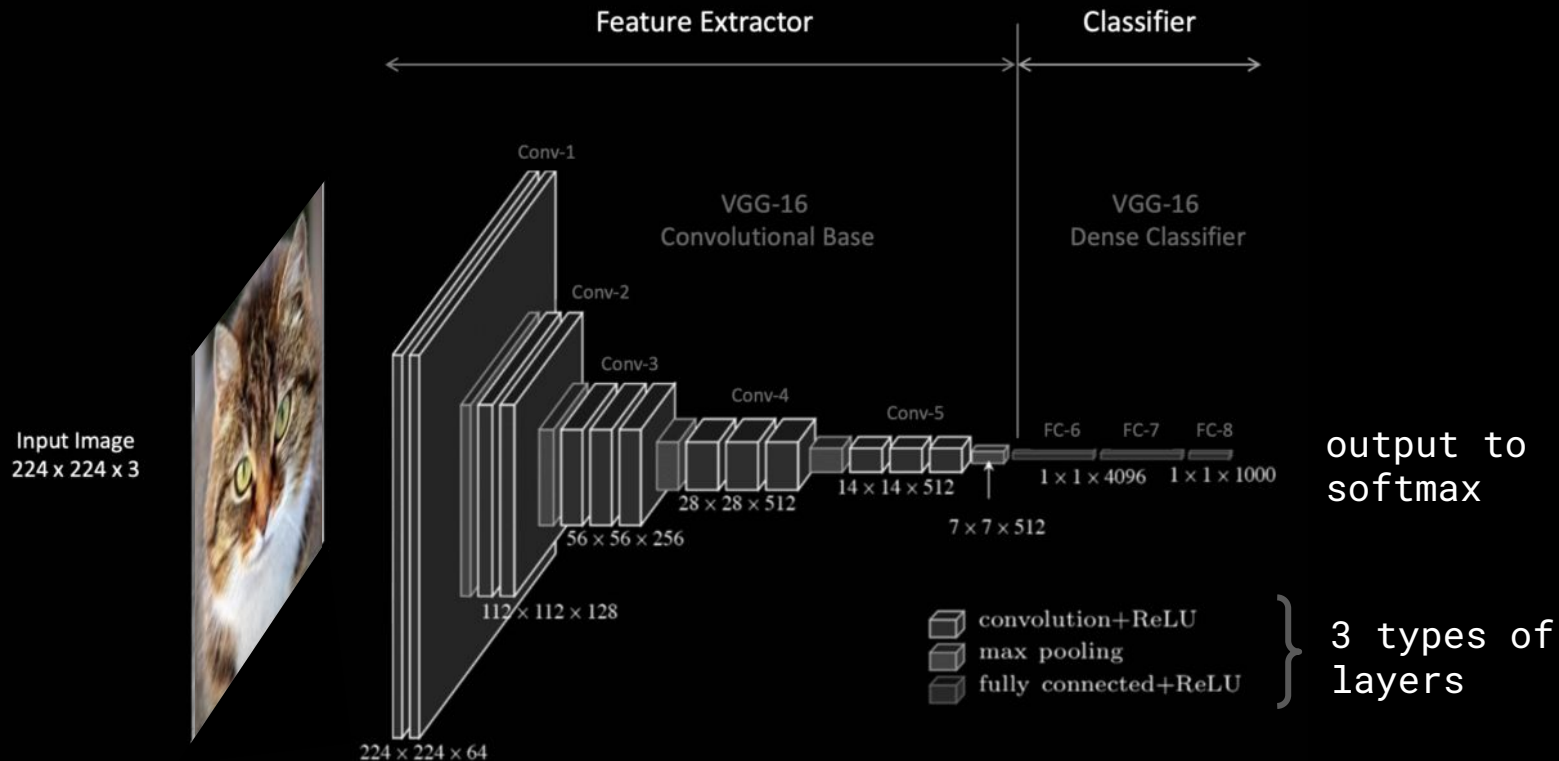
# Simple Models

test image					training image					pixel-wise absolute value differences				
56	32	10	18		10	20	24	17		46	12	14	1	
90	23	128	133		8	10	89	100		82	13	39	33	
24	26	178	200	-	12	16	178	170	=	12	10	0	30	→ 456
2	0	255	220		4	32	233	112		2	32	22	108	

- Nearest Neighbor Classifier (find the image distance)

$$d(I_1, I_2) = \sum_p |I_1(p) - I_2(p)|_1$$

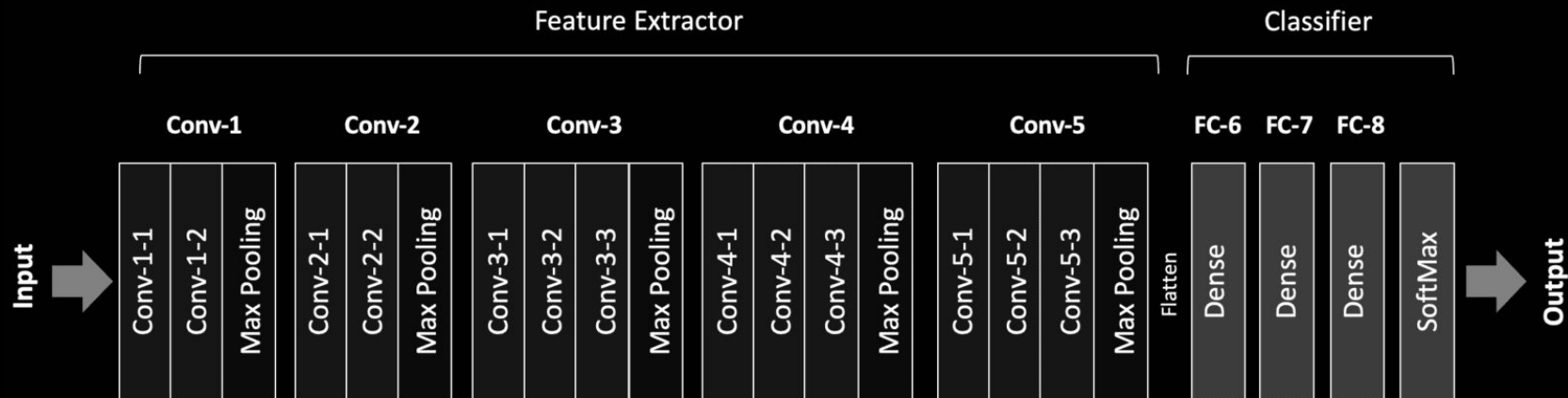
# Convolutional Neural Networks





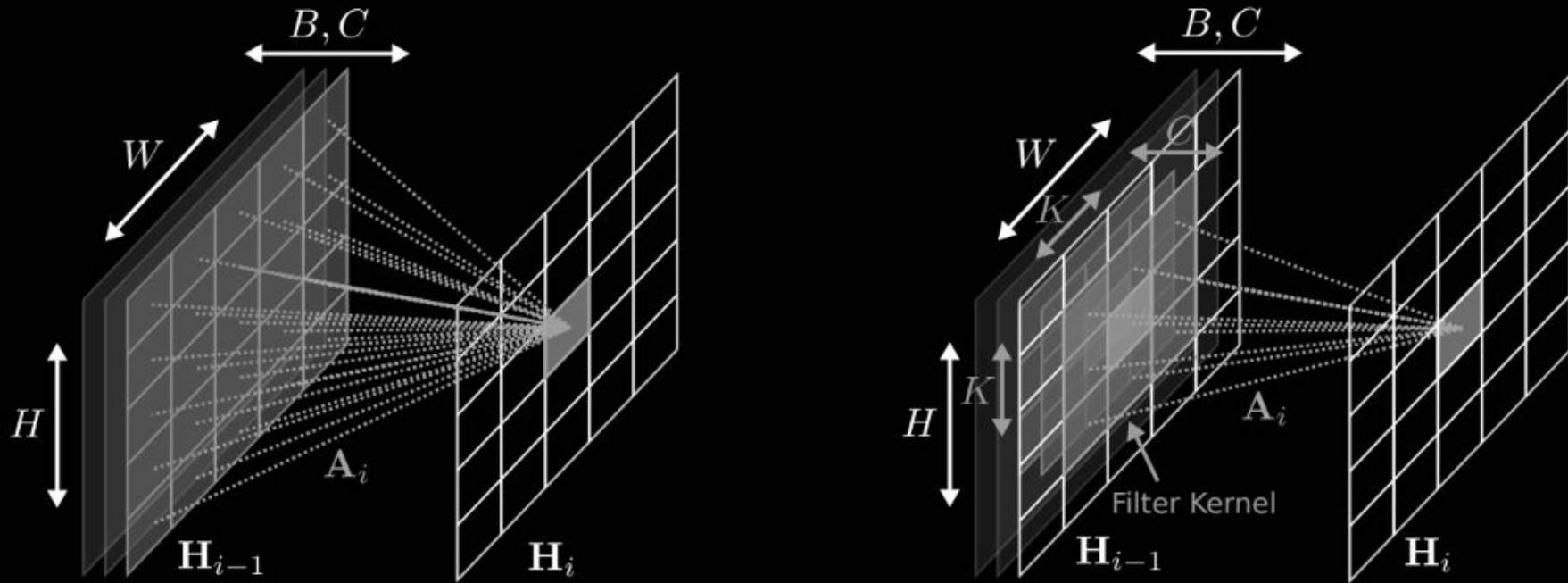
# Convolutional Neural Networks

- Convolutional layers
- Downsampling layers
- Fully Connected layers



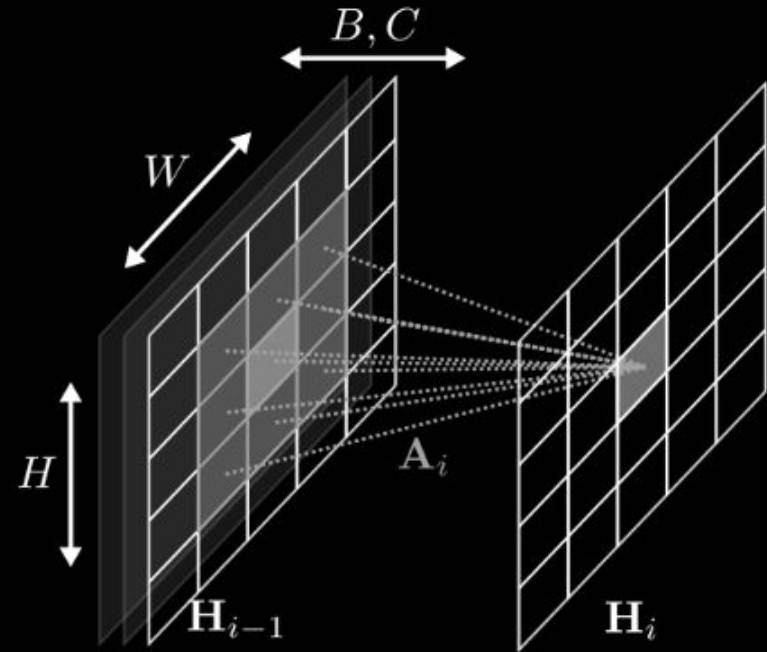
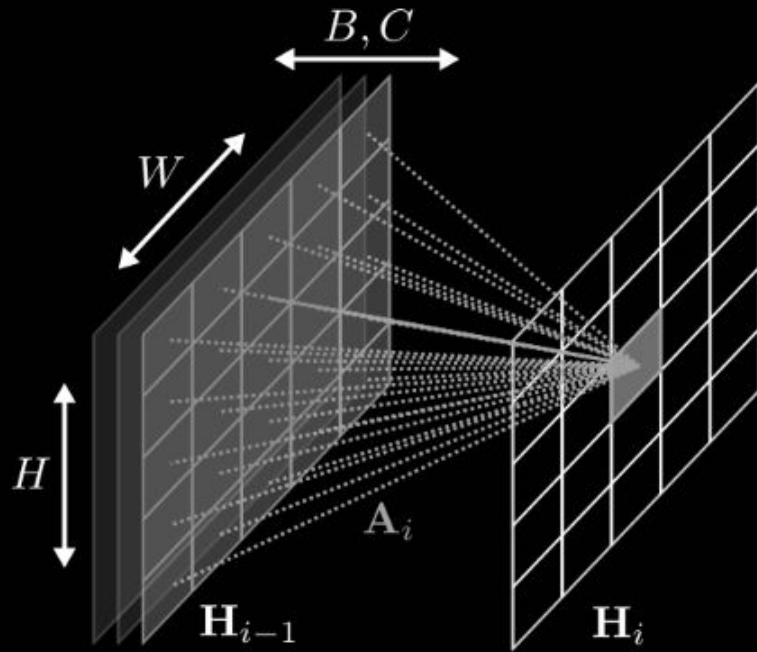
# Convolutional Layer

# Fully Connected vs. Convolutional Layers

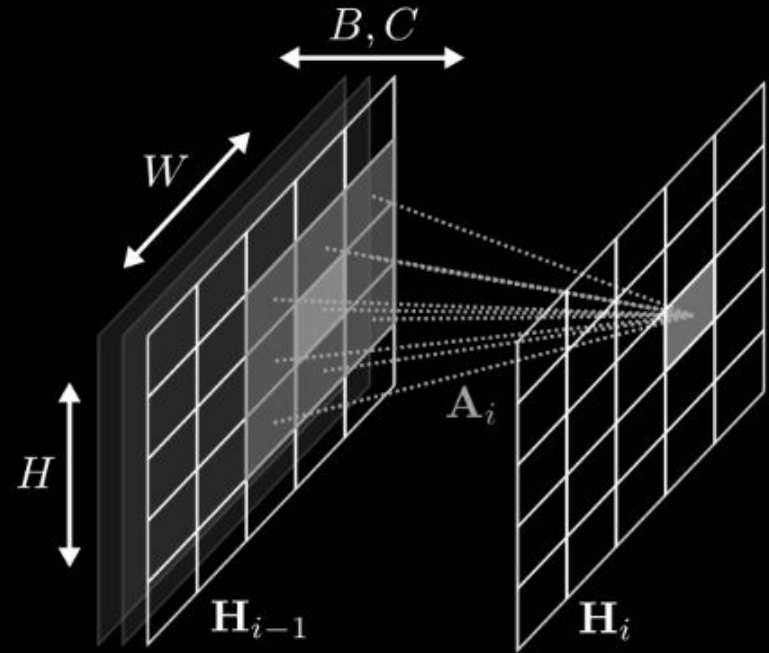
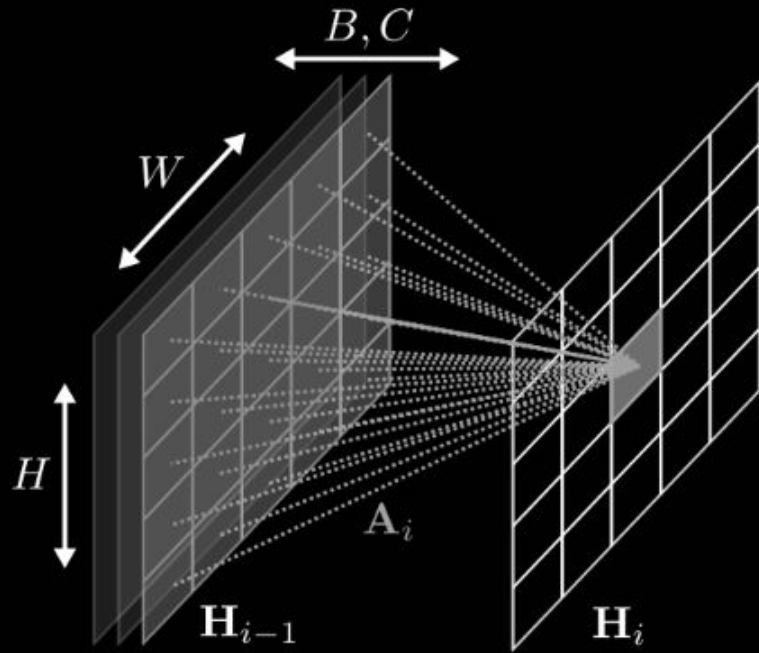


- Compared to fully conn. layers (left), convolutional layers (right) share weights which reduces parameters significantly.
- Convolutional layers are **translationally equivariant**.

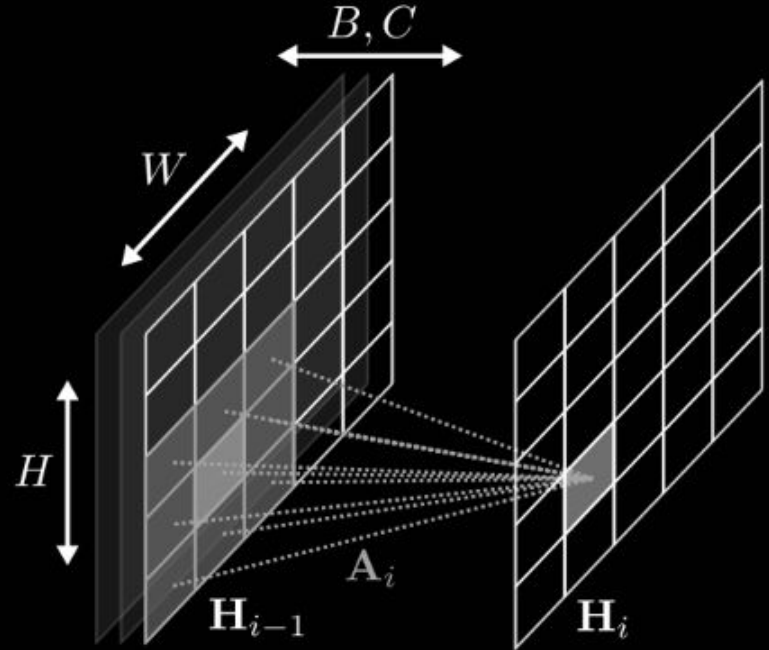
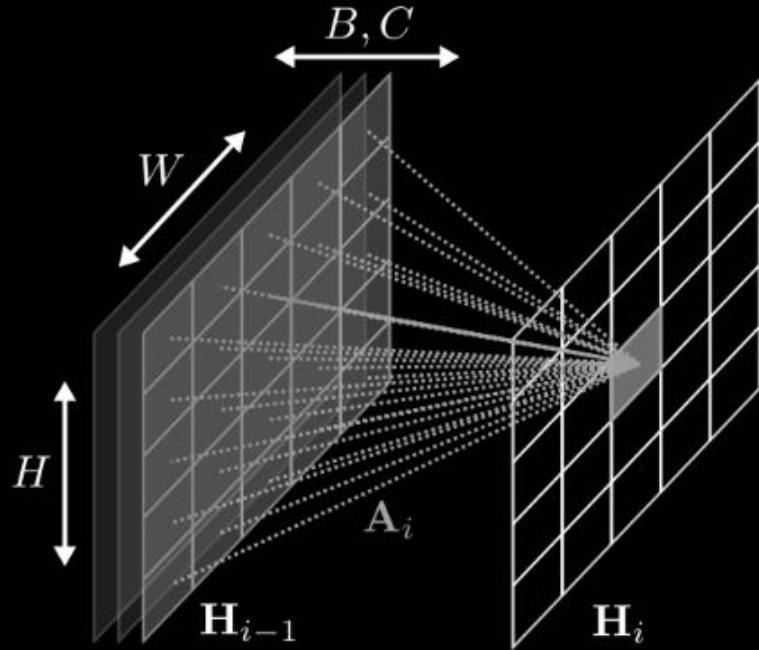
# Fully Connected vs. Convolutional Layers



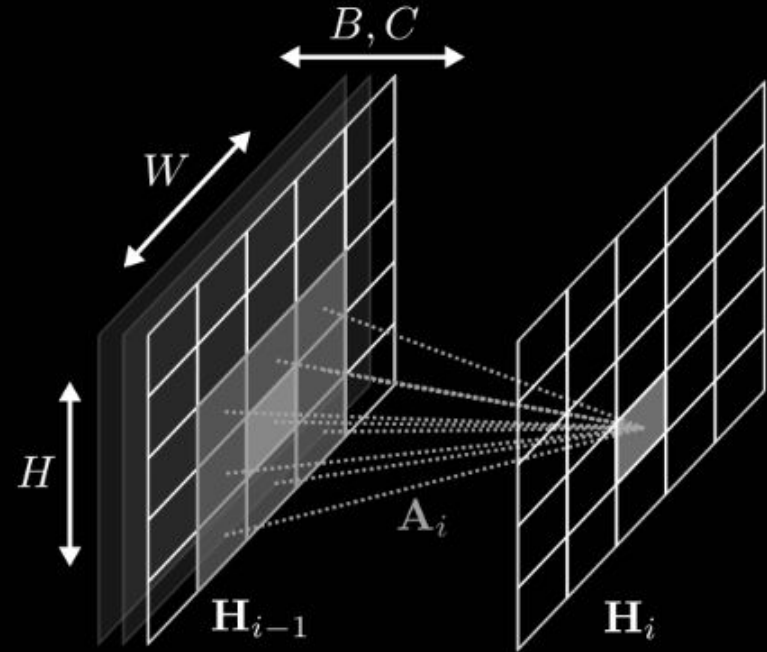
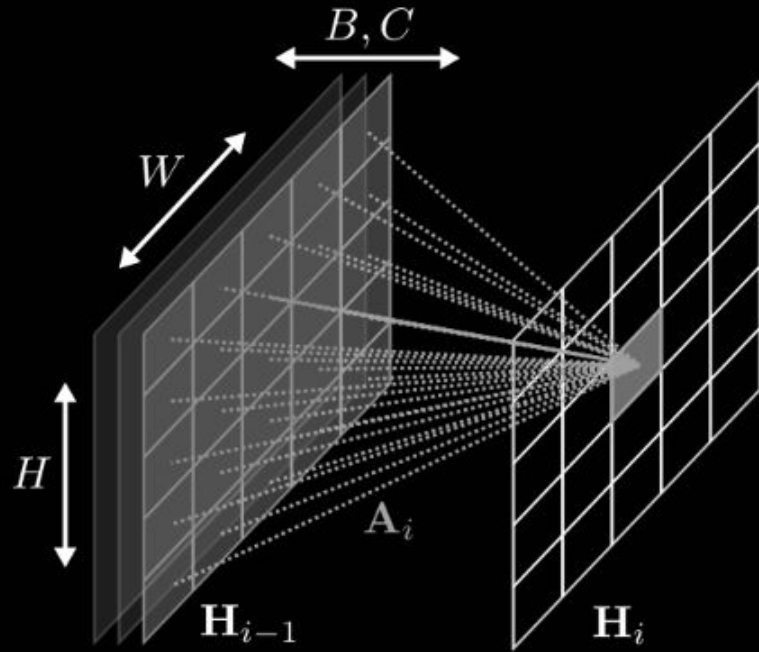
# Fully Connected vs. Convolutional Layers



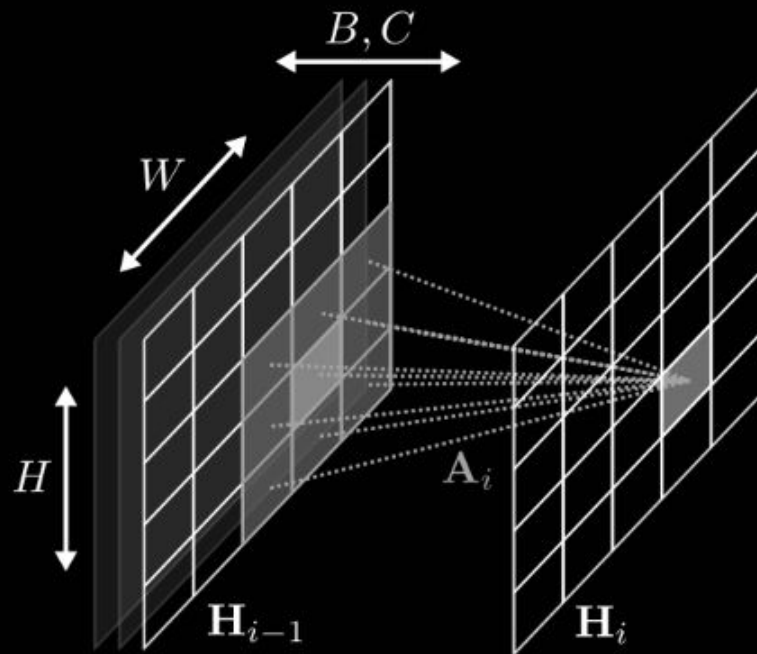
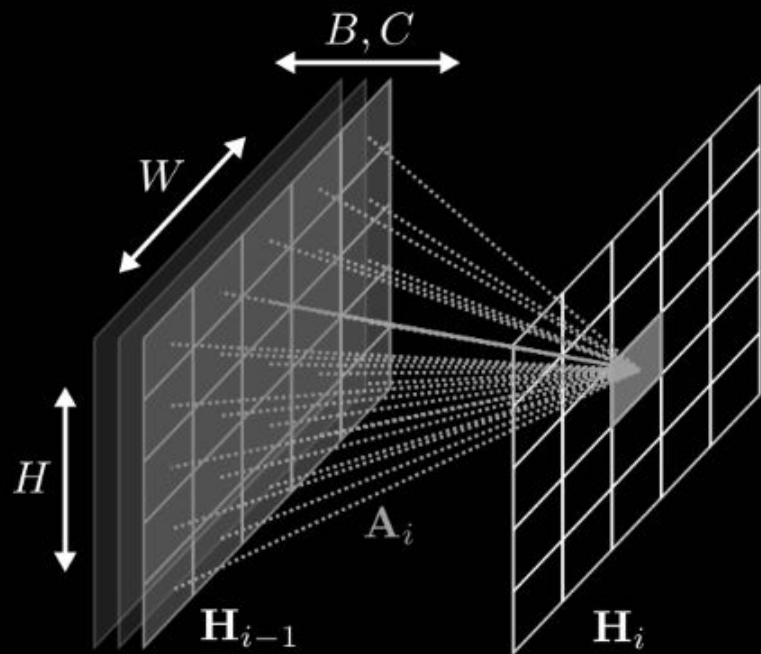
# Fully Connected vs. Convolutional Layers



# Fully Connected vs. Convolutional Layers

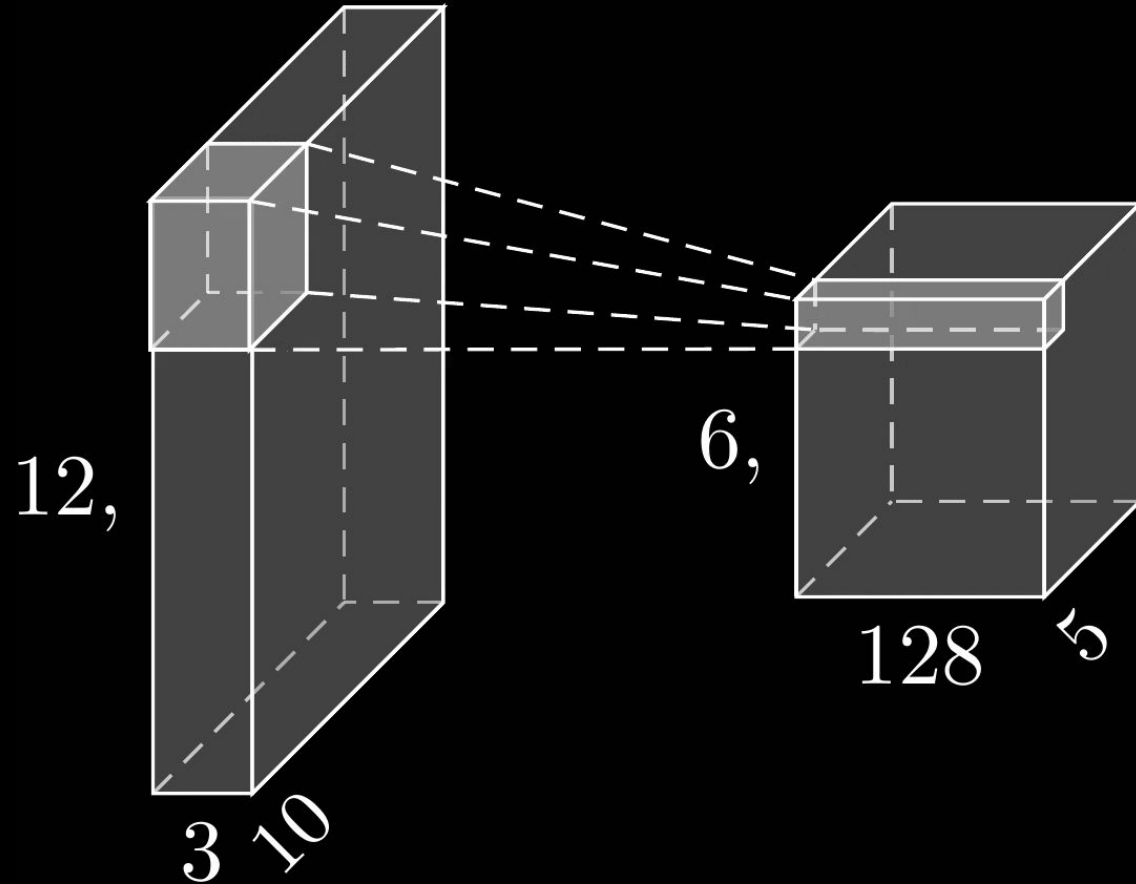


# Fully Connected vs. Convolutional Layers





# Convolution is Happening in 3D

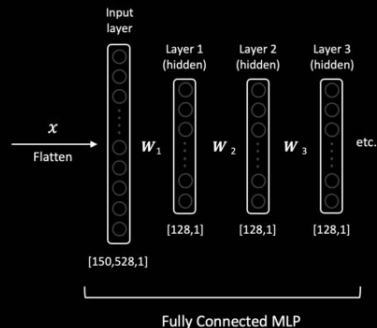
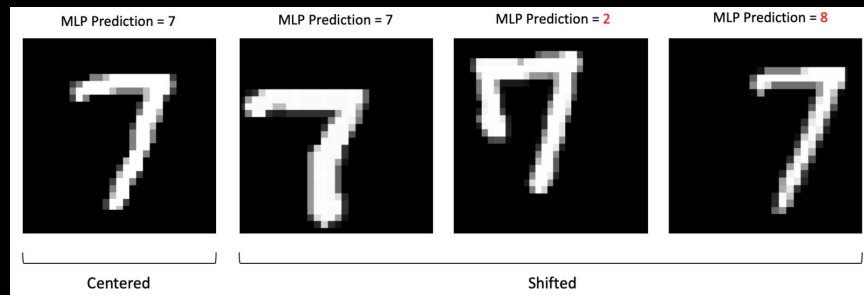


\*source: [C. White](#)

# Fully Connected MultiLayer Perceptron (MLP) vs Convolutional Layers

## Multilayer Perceptron:

- Not translational invariant
- Prone to overfitting
  - Since every pixel connects to every neuron the number of weights becomes HUGE which means training takes longer which means the number of parameters to tune becomes larger!



### Number of Trainable Weights

$$W_1 = 150,528 \times 128 = 19,267,584 \text{ Weights}$$

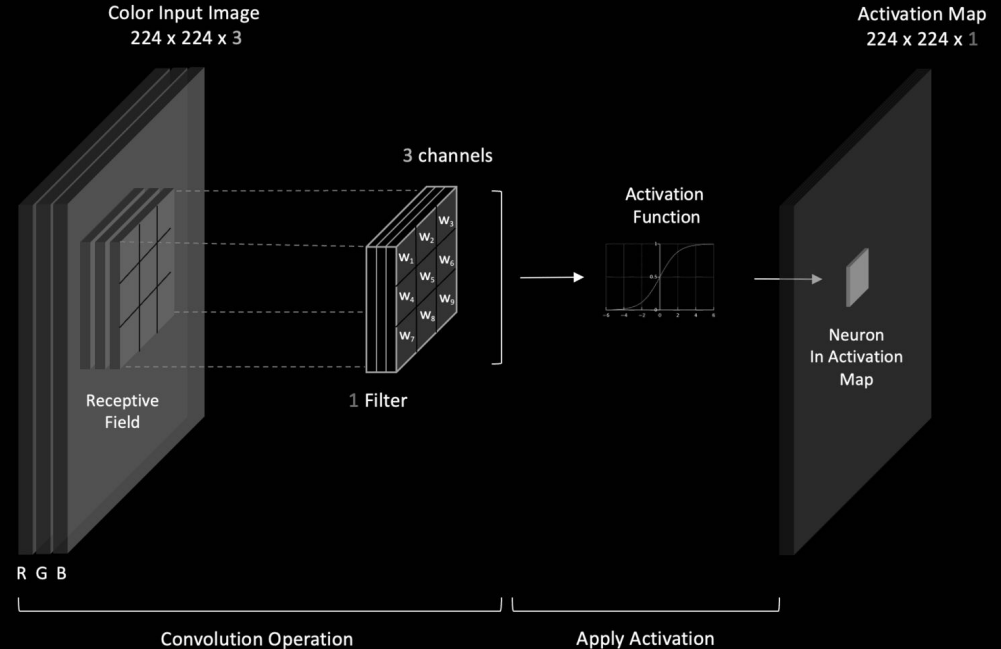
$$W_2 = 19,267,584 \times 128 \sim 2.4 \text{ Billion}$$

$$W_3 = 2.48 \times 128 > 300 \text{ Billion}$$

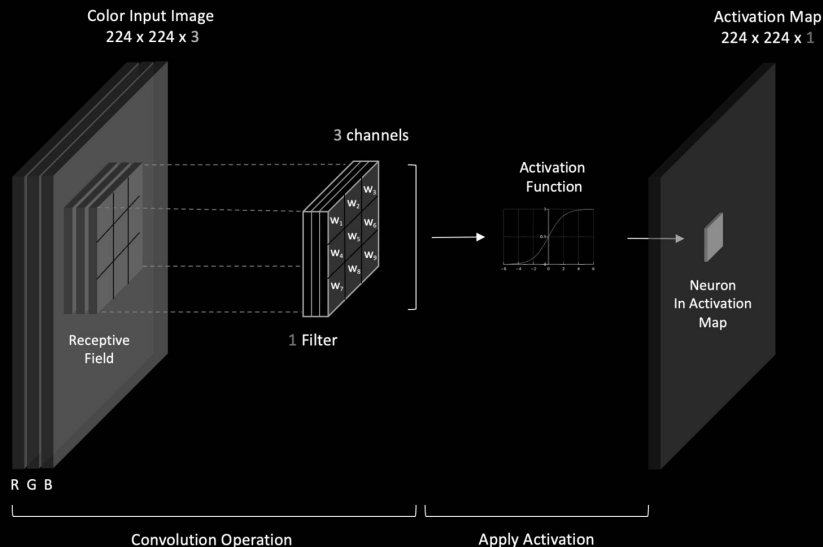
# Fully Connected MultiLayer Perceptron (MLP) vs Convolutional Layers

## ConvNet:

- Translational equivariance (invariance for classification)
  - parameter sharing: same weights are used to process different parts of the input image
- Eyes of Network

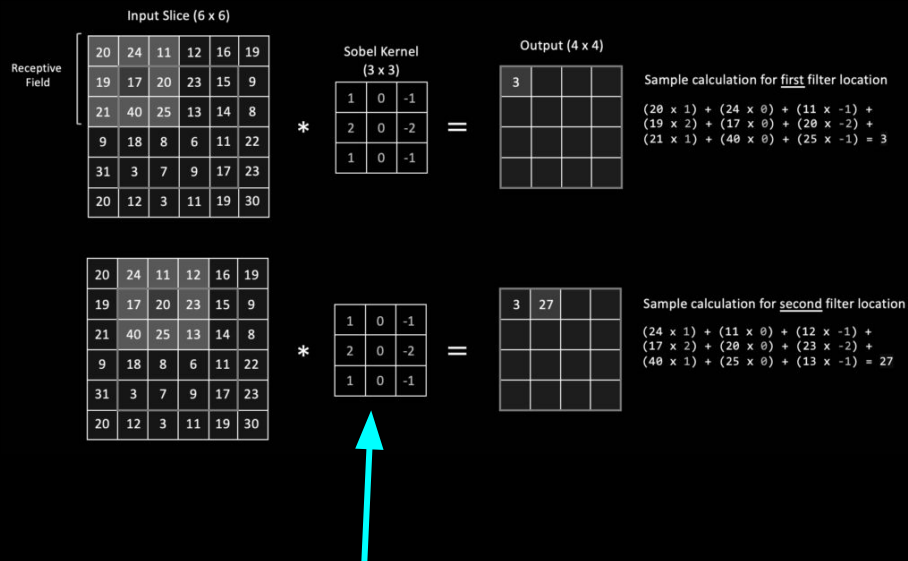


# Convolutional Layers



- **filter/kernel** moves across the input, and at each filter location, a convolution operation is performed, which produces a single number
- This value is then passed through an **activation function**, and the output from the activation function populates the corresponding entry in the output, also known as an **activation map**
- the **activation map** is a summary of features from the input via the convolution process
- elements in the **kernel** are weights that are learned by the network during training
- Sliding the filter one pixel at a time corresponds to a **stride** of one
- The region of the filter performs the operation on is called the **receptive field**

# Convolutional Layers



\*\* Sobel Filter for math demonstration purposes only.  
Actual elements in the kernel are weights that are learned by the network during training

- **filter/kernel** moves across the input, and at each filter location, a convolution operation is performed, which produces a single number
- This value is then passed through an **activation function**, and the output from the activation function populates the corresponding entry in the output, also known as an **activation map**
- the **activation map** is a summary of features from the input via the convolution process
- elements in the **kernel** are weights that are learned by the network during training
- Sliding the filter one pixel at a time corresponds to a **stride** of one
- The region of the filter performs the operation on is called the **receptive field**

# Convolutional Layer Properties

- The **depth of a filter** (channels) must **match** the depth of the input data (i.e., the number of channels in the input).
- The **spatial size** of the filter is a **design choice** but 3x3 is very common (or sometimes 5x5).
- The number of **filters** is also a **design choice** and dictates the number of activation maps produced in the output.
- Multiple **activations maps** are sometimes collectively referred to as “an activation map containing multiple channels.” But we often refer to each channel as an activation map.
- Each channel in a filter is referred to as a **kernel**, so you can think of a filter as a container for kernels.
- A single-channel filter has just one kernel. Therefore, in this case, the filter and the kernel are one and the same.
- The **weights in a filter/kernel** are initialized to **small random values** and are **learned** by the network during training.
- The **number of trainable parameters** in a convolutional layer depends on the **number of filters**, the **filter size**, and the **number of channels** in the input.

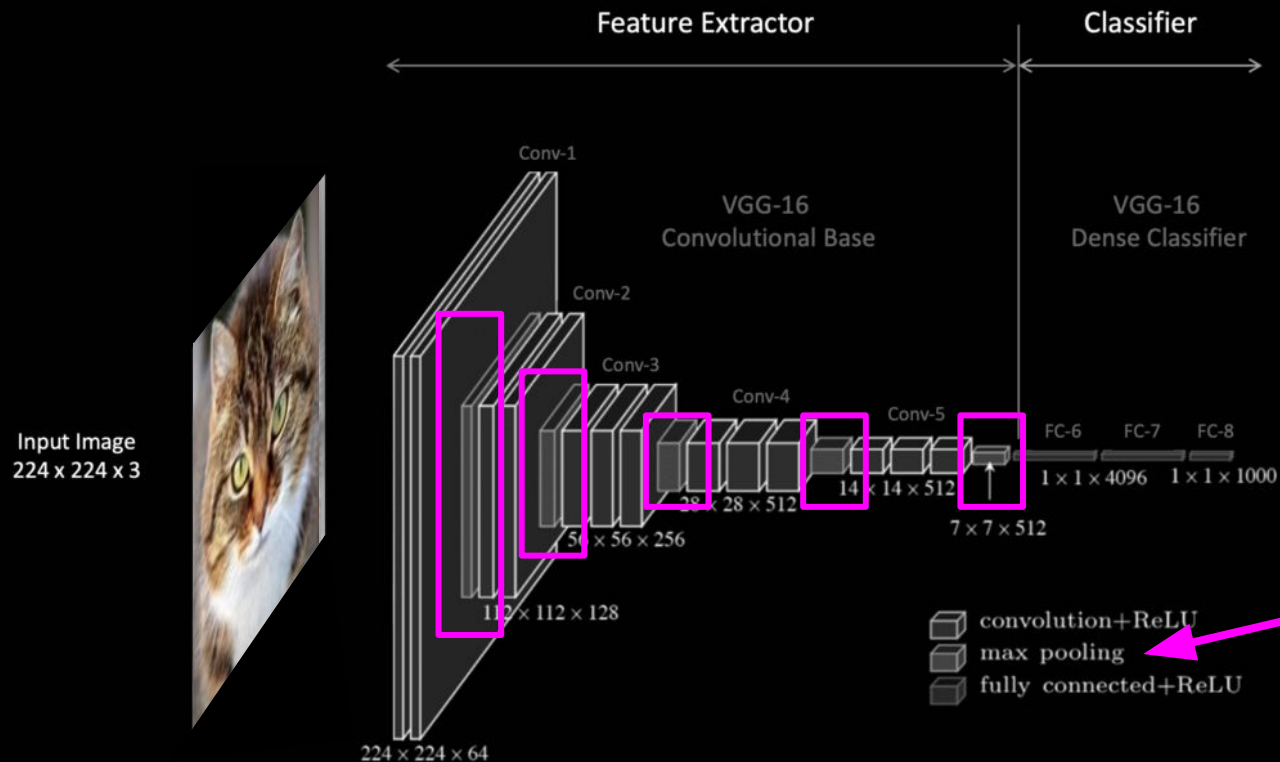
# TLDR: Convolutional Layer

- **Input:** volume of size  $W_1 \times H_1 \times D_1$
- Requires 4 hyperparameters
  - Number of filters  $K$
  - The kernel size/spatial extent  $F$
  - The stride  $S$
  - The amount of zero padding  $P$
- **Output:** volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F + 2P)/S + 1$
  - $H_2 = (H_1 - F + 2P)/S + 1$
- With parameter sharing, it introduces  $F \times F \times D_1$  weights per filter, for a total of  $(F \times F \times D_1) - K$  weights and  $K$  biases.
- In the output volume, the  $d$ th depth slice (of size  $W_2 \times H_2$ ) is the result of performing a valid convolution of the  $d$ th filter over the input volume with a stride of  $S$ , and then offset by  $d$ th bias

**Downsampling Layer**



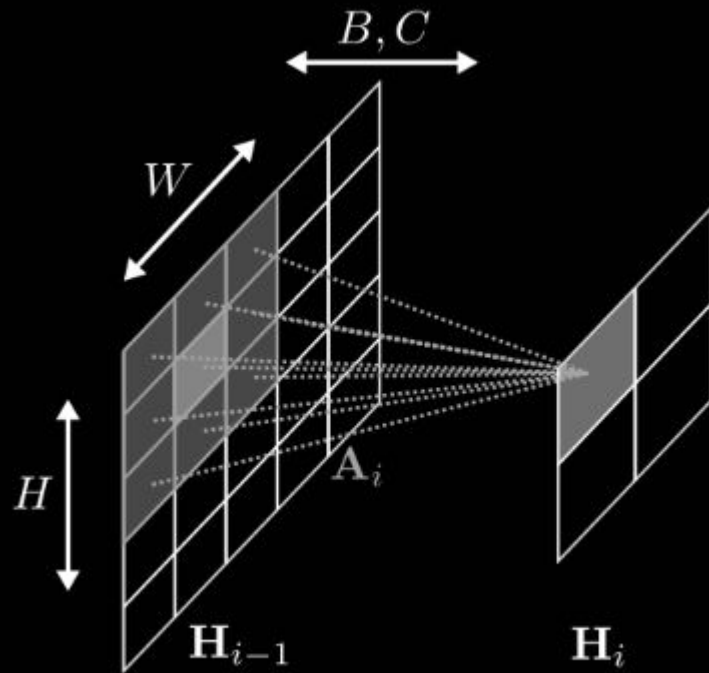
# Downsampling Layer



- reduces the **spatial resolution**
- increases the **receptive field**

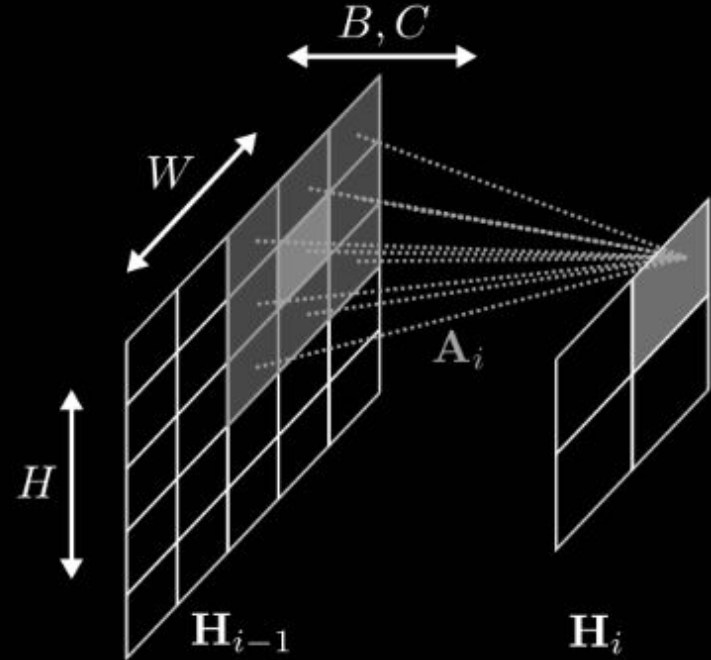
# Downsampling: Pooling Layer

- Stride  $s=2$  and kernel size  $2 \times 2$  reduces spatial dimension by 2
- Pooling has no parameters
  - max pooling
  - min pooling
  - mean pooling
- Pooling is applied to each channel separately preserving the number of channels



# Downsampling: Pooling Layer

- Stride  $s=2$  and kernel size  $2 \times 2$  reduces spatial dimension by 2
- Pooling has no parameters
  - max pooling
  - min pooling
  - mean pooling
- Pooling is applied to each channel separately preserving the number of channels



# Pooling Layer Example: Max Pooling

Single Depth Slice

7	2	5	2
4	5	4	7
3	3	4	2
6	4	8	6

Max Pooling

2x2 Filter & Stride of 2

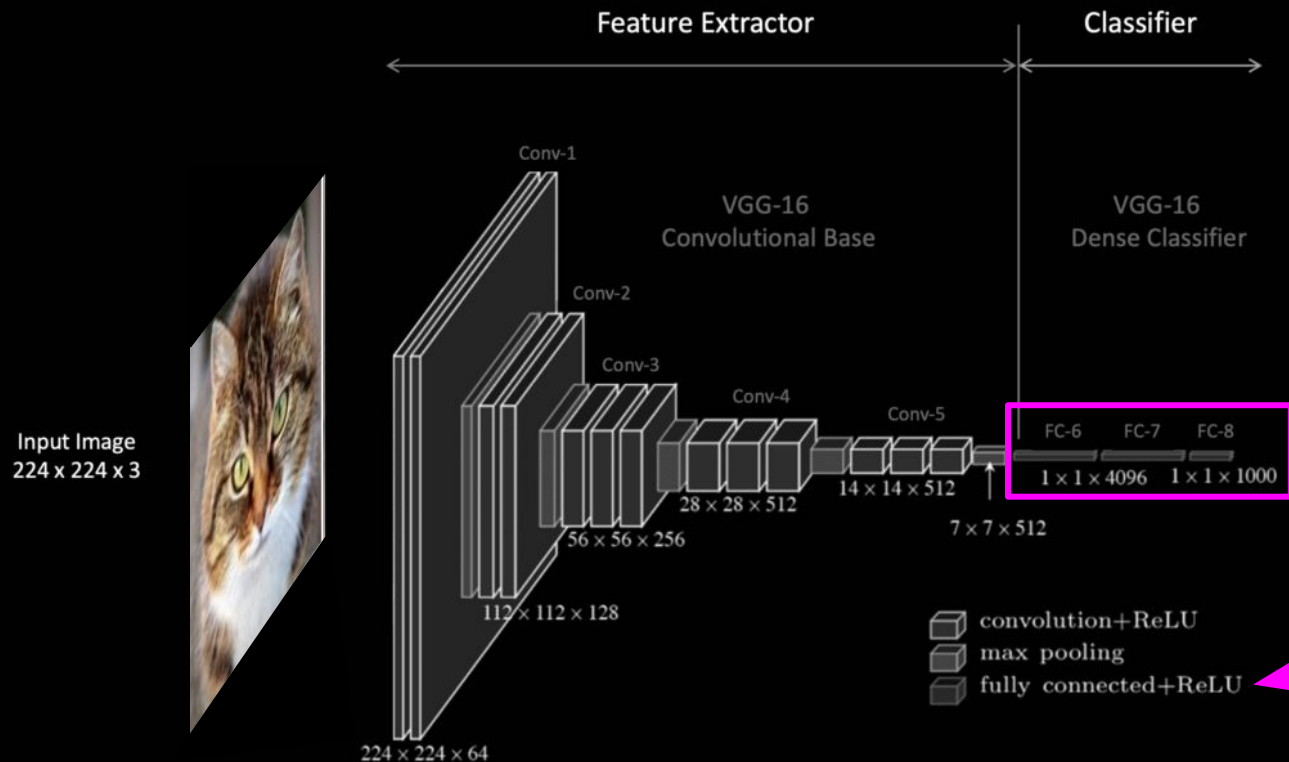
7	7
6	8

# TLDR: Pooling Layer

- **Input:** volume of size  $W_1 \times H_1 \times D_1$
- Requires two hyperparameters:
  - The spatial extent/kernel size  $F$
  - The stride  $S$
- Produces a volume of size  $W_2 \times H_2 \times D_2$ 
  - $W_2 = (W_1 - F) / S + 1$
  - $H_2 = (H_1 - F) / S + 1$
  - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- For Pooling layers, it is not common to pad the input using zero-padding

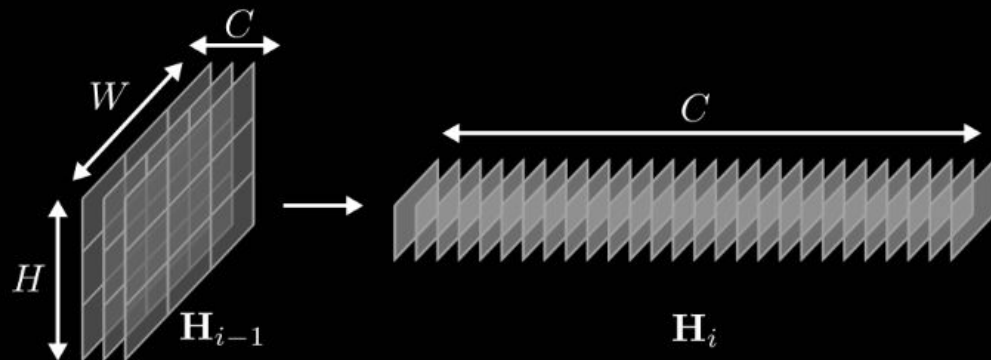
**Fully Connected Layer**

# Fully Connected Layer



- The most memory intensive because all pixels contribute

# Fully Connected Layer

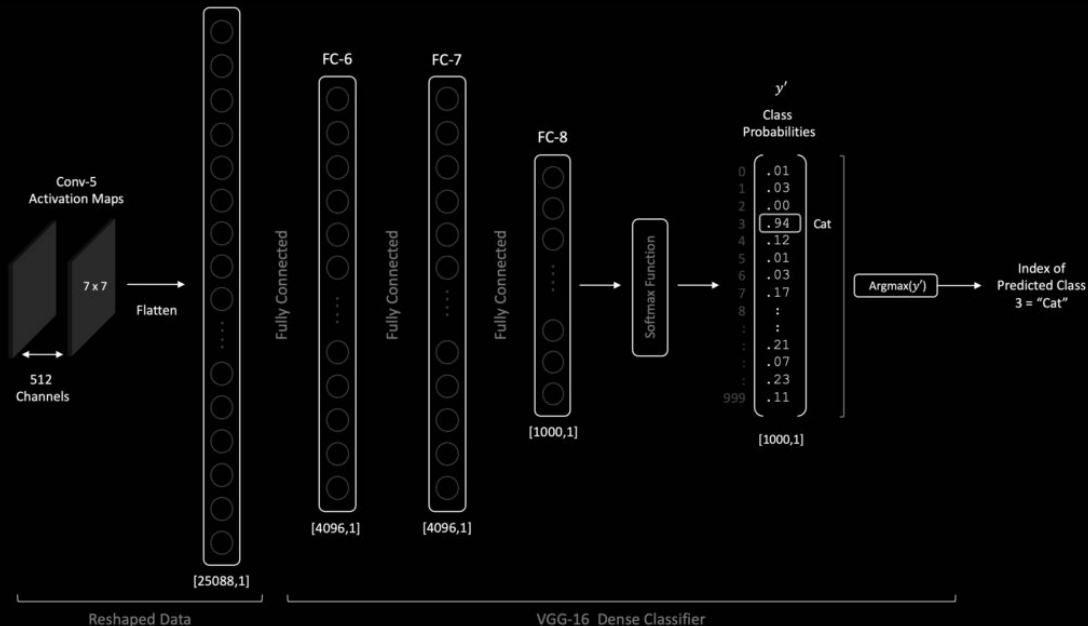


- Reshape  $H_{i-1}[B, X, Y, C]$  into  $H_i[B, C]$
- Now  $X$  and  $Y$  are reshaped into the feature channel dimension  $C$



# Fully Connected Layer

- The number of neurons in the output layer of the network is equal to the number of classes
- Each output neuron is associated with a specific class
- Each neuron's value (after the softmax layer) represents the probability that the associated class for that neuron corresponds to the class associated with the input image.



# The Output

Notice that the output of our fully connected layer becomes the input of our **softmax** function which outputs a one dimensional array that maps probabilities over labels or categories.

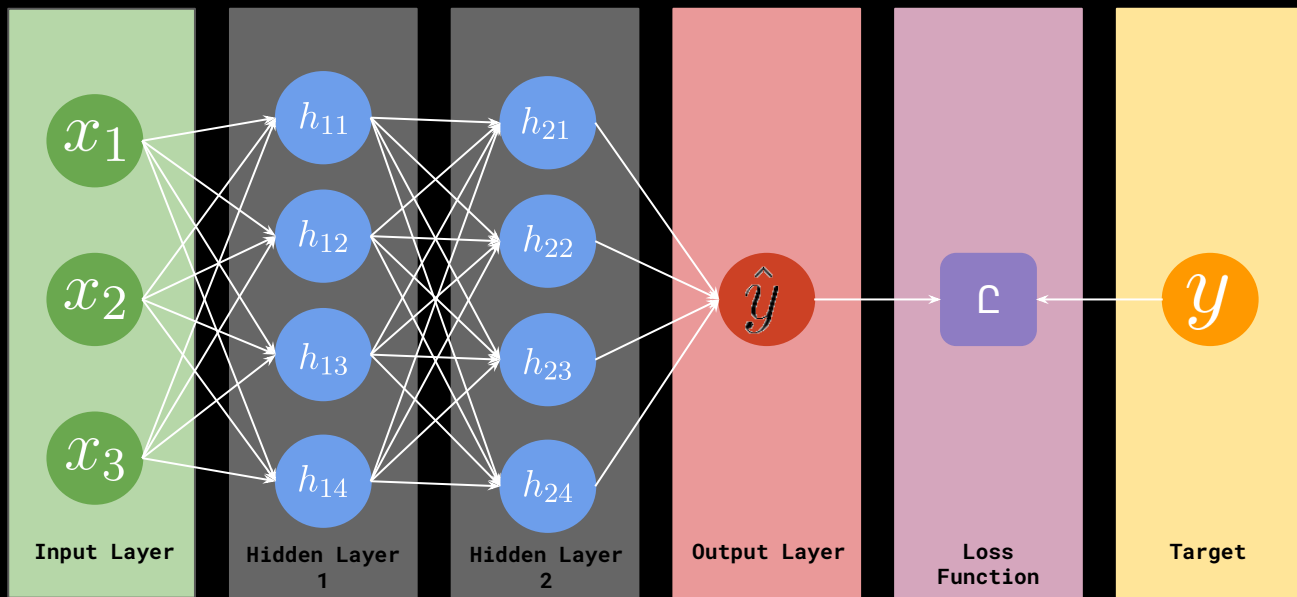
Since image classification results in a label or category, we need to convert that label which is a string to a numeric value that the model can understand.

Categorical data can be

- nominal: **one-hot encoded** (preferred)
- ordinal: label encoded (arbitrary numeric)

# Loss Function

# The Output and Loss



- The **output layer** is the last layer which calculates the output
- The **loss function** compares the result of the output layer to the target
- In **image classification** we use a **softmax** output layer and a **cross entropy loss** to get probability values over labels.

# Model Output: Categorical Distribution

- Obtain a probability distribution over the possible classes that the model can predict
- A vector of scores for each possible class
- The representation is important because it allows us to quantify the model's uncertainty about the predictions and compare it with the ground truth labels.

**Categorical Distribution:** Index label

$$p(y = c) = \mu_c$$

- $\mu_c$ : probability for class c

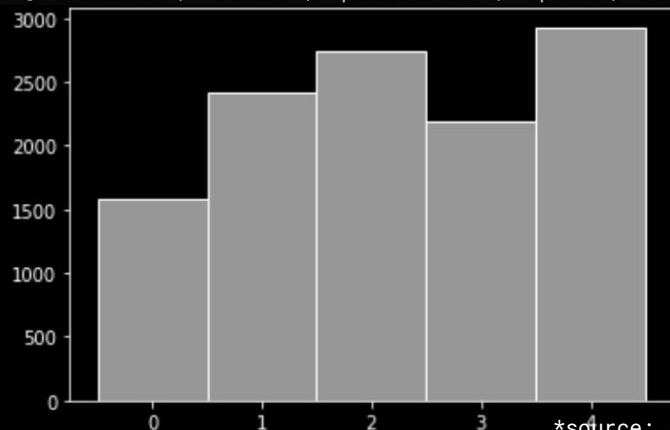
One Hot Encoding Notation:

$$p(y) = \prod_{c=1}^C \mu_c^{y_c} \quad y_c \in 0, 1$$

- $y$ : "one-hot" vector with
- $y = (0, \dots, 0, 1, \dots, 0)^T$  with all zeros except for one (the true class)




`classes = ['filaments', 'flares', 'prominences', 'quiet', 'sunspots']`



\*source: [A. Geiger](#)  
S. Chatterjee

# One-Hot Vector Representation

Class	Label	$y(\text{indexed})$	$y(\text{one-hot})$
	"Cat"	1	$[1, 0]^T$
	"Dog"	2	$[0, 1]^T$

- One hot vector  $y$  w/ binary elements
- Index  $c$  where  $y_c = 1$  determines the correct class

# Categorical Distribution & Cross-Entropy Loss

Let  $p_{model}(y|x, w) = \prod_{c=1}^C f_w^{(c)}(x)^{y_c}$  be a **Categorical distribution**

Maximizing the log-likelihood leads to the **cross-entropy** loss:

$$\begin{aligned}\hat{\mathbf{w}}_{ML} &= \operatorname{argmax}_{\mathbf{w}} \sum_{i=1}^N \log p_{model}(\mathbf{y}_i | \mathbf{x}_i, \mathbf{w}) \\ &= \operatorname{argmax}_{\mathbf{w}} \sum_{i=1}^N \log \prod_{c=1}^C f_{\mathbf{w}}^{(c)}(\mathbf{x}_i)^{y_{i,c}} \\ &= \operatorname{argmin}_{\mathbf{w}} \underbrace{\sum_{i=1}^N \sum_{c=1}^C -y_{i,c} \log f_{\mathbf{w}}^{(c)}(\mathbf{x}_i)}_{\text{CE Loss}}\end{aligned}$$

[Math explanation: A. Webb](#)

The target  $\mathbf{y} = (0, \dots, 0, 1, 0, \dots, 0)^\top$  is a “one-hot” vector with  $y_c$  as it's  $C$ th element.

\*source: [A. Geiger](#)

# Softmax

How can we ensure that  $f_w^{(c)}(x)^{y_c}$  predicts a valid **categorical distribution**?

- We must guarantee 2 things:

$$f_w^{(c)}(x) \in [0, 1] \qquad \sum_{c=1}^C f_w^{(c)}(x) = 1$$

- Use **softmax** to do this:

$$\text{softmax}(\mathbf{x}) = \left( \frac{\exp(x_1)}{\sum_{k=1}^C \exp(x_k)}, \dots, \frac{\exp(x_C)}{\sum_{k=1}^C \exp(x_k)} \right)$$

- Let the score vector **s** denote the network output after the last affine layer, then:

$$f_{\mathbf{w}}^{(c)}(\mathbf{x}) = \frac{\exp(s_c)}{\sum_{k=1}^C \exp(s_k)} \quad \Rightarrow \quad \log f_{\mathbf{w}}^{(c)}(\mathbf{x}) = s_c - \log \sum_{k=1}^C \exp(s_k)$$

\*source: [A. Geiger paperswcode](#)



# Log Softmax

Let the score vector **s** denote the network output after the last affine layer.

$$f_{\mathbf{w}}^{(c)}(\mathbf{x}) = \frac{\exp(s_c)}{\sum_{k=1}^C \exp(s_k)} \Rightarrow \log f_{\mathbf{w}}^{(c)}(\mathbf{x}) = s_c - \log \sum_{k=1}^C \exp(s_k)$$

- If  $c$  is the correct class label we would like to maximize the log softmax
  - we would like  $s_c$  for the correct class to increase
  - we would like  $s_k$  to decrease (negative)
  - The second term can be approximated to  $\log \sum_{k=1}^C \exp(s_k) \approx \max_k s_k$  since  $\exp(s_k)$  insignificant for all  $s_k < \max_k s_k$
  - The loss is always strongly penalized for incorrect predictions
  - If the correct class has the largest score, both terms cancel and it will contribute very little to the overall training cost.

# Softmax

- Softmax responds to differences between inputs
- It is invariant to adding the same scalar to all its inputs:

$$\text{softmax}(x) = \text{softmax}(x + c)$$

- Another variant of this is (taking into account floating point approximations required by computation):

$$\text{softmax}(x) = \text{softmax}(x - \max_{k=1,\dots,L} x_k)$$





- It allows for accurate computation even when  $x$  is large

# Cross Entropy Loss

The cross entropy loss for a single training sample  $(x, y) \in \chi$ :

$$\sum_{c=1}^C -y_c \log f_w^{(c)}(x)$$

Example: Suppose 4 classes and 4 training samples:

Input $x$	Label $y$	Predicted $s$	softmax( $s$ )	CE Loss
 "Cat"	$[1, 0, 0, 0]^T$	$[+3, +1, -1, -1]^T$	$[0.85, 0.12, 0.02, 0.01]^T$	0.16
 "Dog"	$[0, 1, 0, 0]^T$	$[+3, +3, +1, +0]^T$	$[0.46, 0.46, 0.06, 0.02]^T$	0.78
 "Bunny"	$[0, 0, 1, 0]^T$	$[+1, +1, +1, +1]^T$	$[0.25, 0.25, 0.25, 0.25]^T$	1.38
 "Horse"	$[0, 0, 0, 1]^T$	$[+3, +2, +3, -1]^T$	$[0.42, 0.15, 0.42, 0.01]^T$	4.87 ← contributes most to loss!

\*source: [A. Geiger](#)



dog

cat

bird

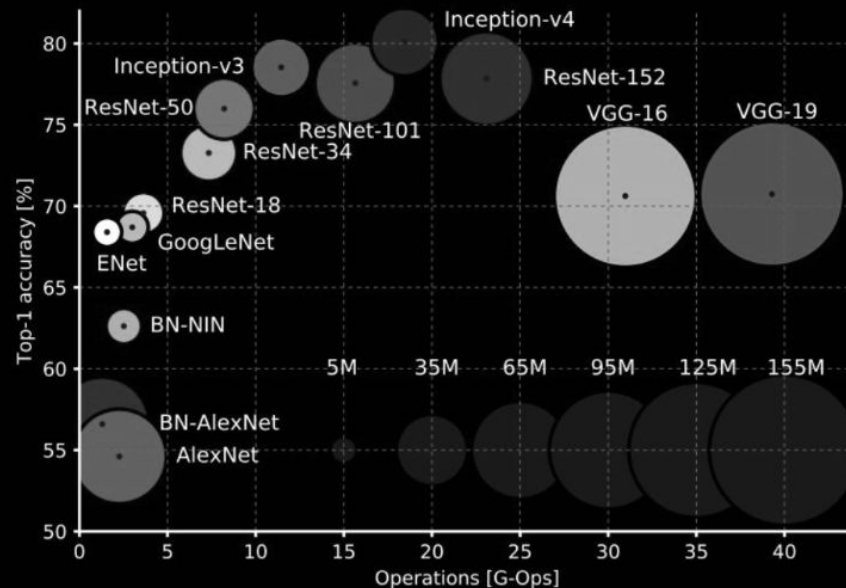
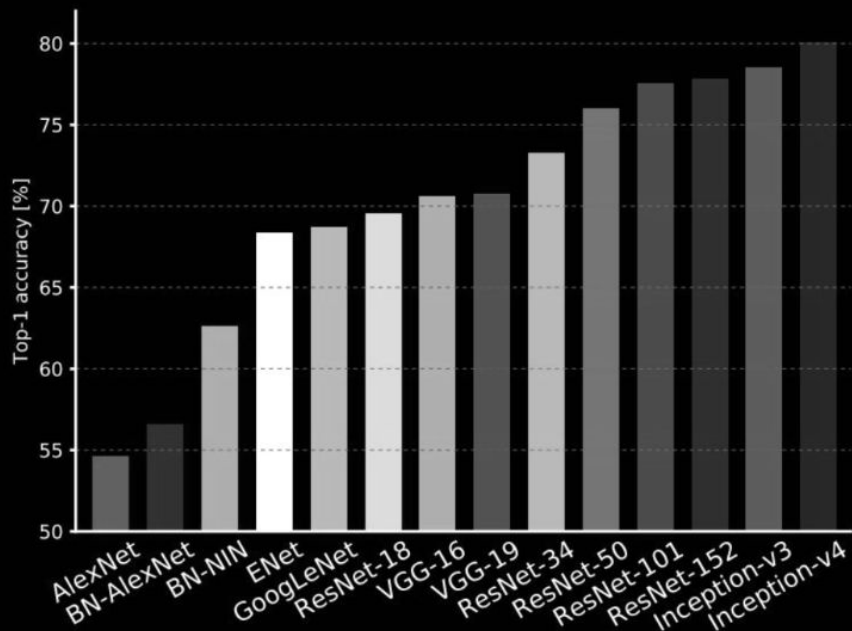
frog

ship

# Popular Image Classification Architectures

- [LeNet-5](#): 2Conv, 2Pool, 2Fully Connected (MNIST)
- [AlexNet](#): 8 layers, ReLus, dropout, augmentation
- [VGG](#) (slowest b/c of number of parameters): 16, 19 layer versions (uses 3x3 receptive field instead of 7x7)
- [Inception/GoogLeNet](#): 22 layers, inception modules
- [ResNet](#): 152 layers use residuals w/ skip connections

# Accuracy vs Complexity



\*source: [A. Geiger](#)