

Computer Vision

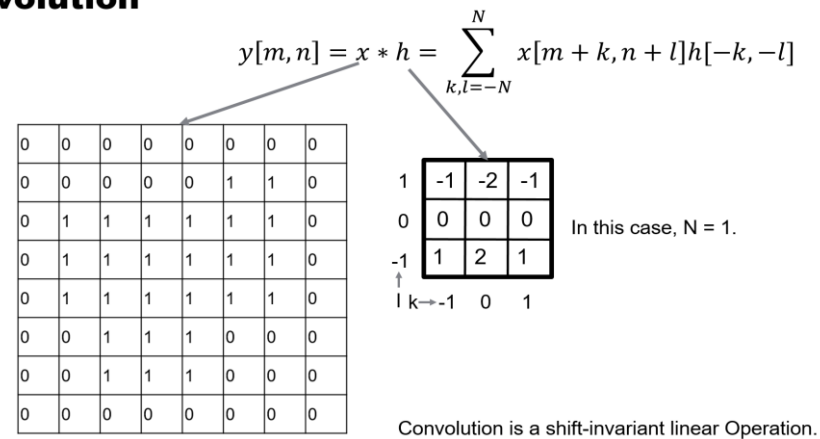


Lecture 4: Machine Learning

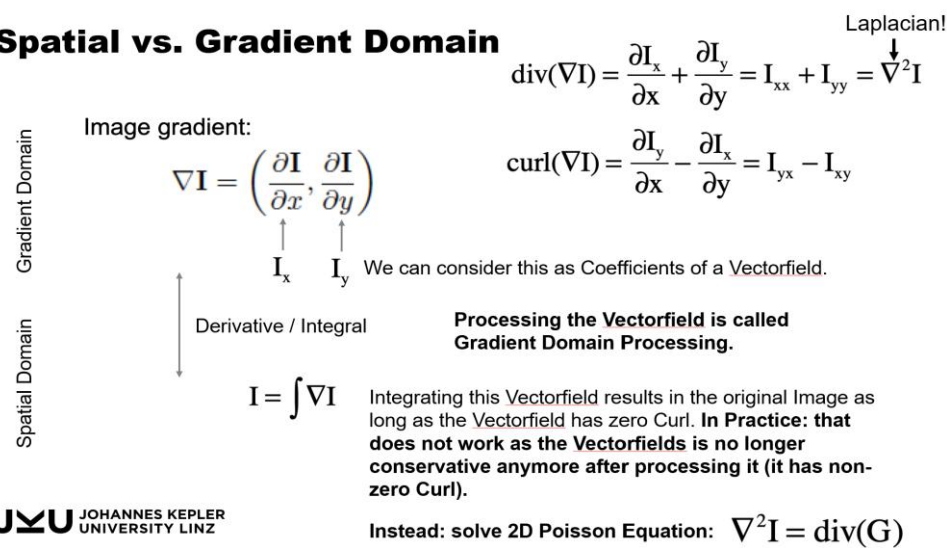
Oliver Bimber

Last Week: Digital Image Processing

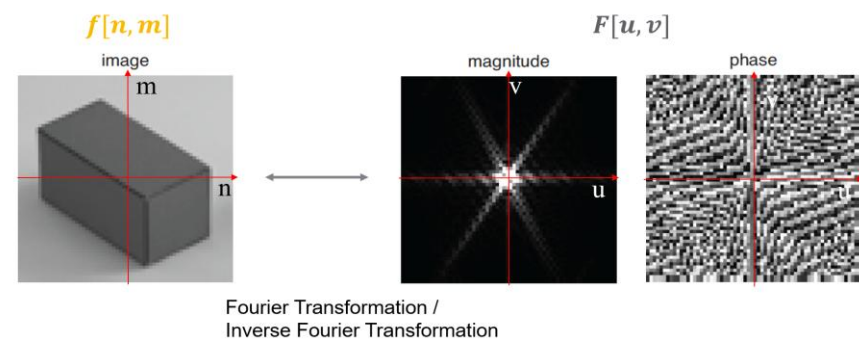
Convolution



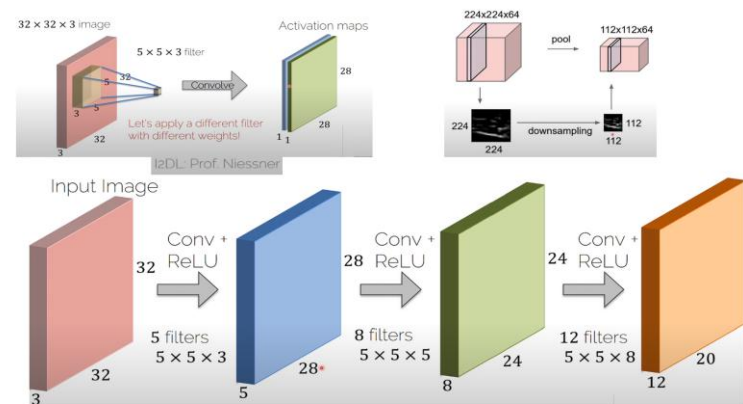
Spatial vs. Gradient Domain



Spatial vs. Frequency Domain



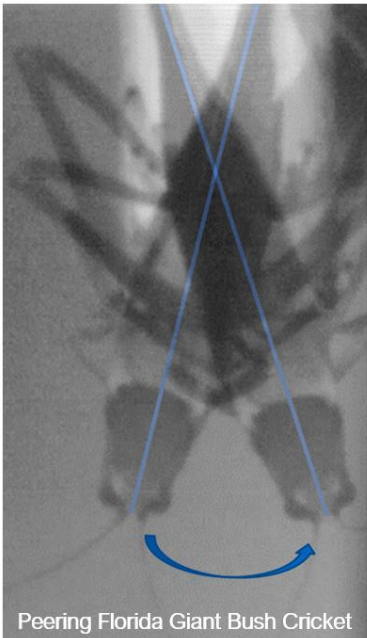
Convolutional Neural Networks (CNNs)



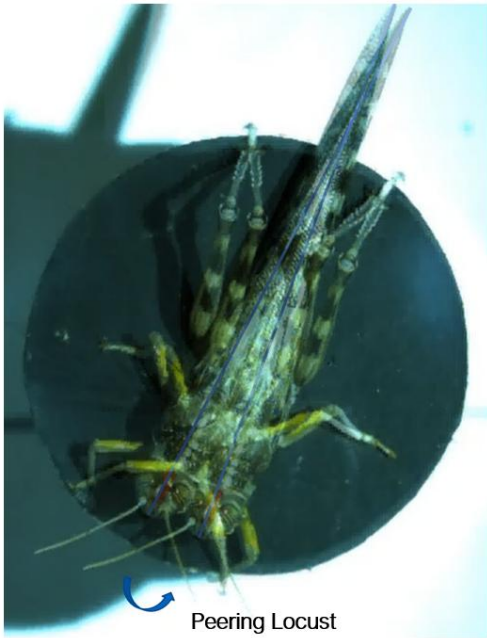
Course Overview

CW	Topic	Date	Place	Lab
41	Introduction and Course Overview	07.10.2025	Zoom	Lab 1
42	Capturing Digital Images	14.10.2025	Zoom	Lab 2
43	Digital Image Processing	21.10.2025	Zoom	Assignment 1
→ 44	Machine Learning	28.10.2025	Zoom	
45	Feature Extraction	04.11.2025	Zoom	Open Lab 1
46	Segmentation	11.11.2025	Zoom	Assignment 2
47	Optical Flow	18.11.2025	Zoom	Open Lab 2
48	Object Detection	25.11.2025	Zoom	Assignment 3
49	Multi-View Geometry	02.12.2025	Zoom	Open Lab 3
50	3D Vision	09.12.2025	Zoom	Assignment 4
3	Trends in Computer Vision	13.01.2026	Zoom	
4	Q&A	20.01.2026	Zoom	Open Lab 4
5	Exam	27.01.2026	HS1 (Linz), S1/S3 (Vienna), S5 (Bregenz)	
9	Retry Exam	24.02.2026	tba	

Research Example: How Robot Dogs See the Unseeable



Peering Florida Giant Bush Cricket



Peering Locust



Peering Quadruped Robot



Metal Fence and Plants

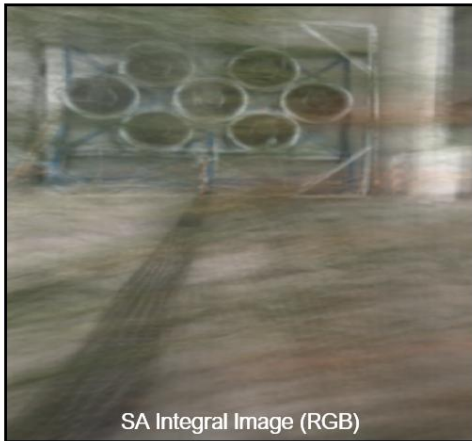


Camera-Equipped Head-Unit



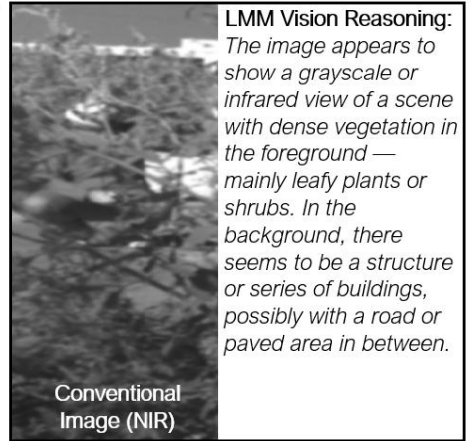
Conventional Image (RGB)

LMM Vision Reasoning:
The image shows some green plants with thin stems and fine leaves growing in front of a metallic fence or railing. Behind the plants, there is an industrial or commercial-looking structure and a concrete floor. The setting appears to be an industrial facility or a technical building exterior.



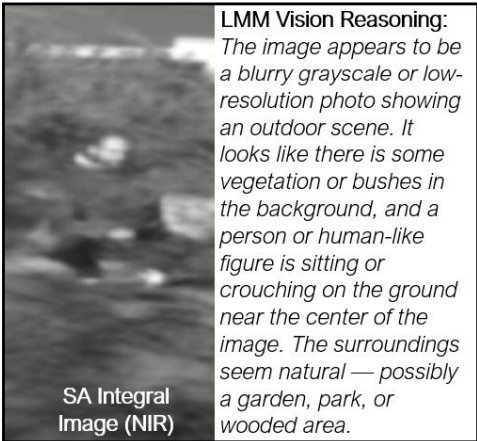
SA Integral Image (RGB)

LMM Vision Reasoning:
The image appears to be quite blurry, but some features can still be made out. It seems to show an industrial or abandoned facility, possibly a warehouse or a factory. In the center, there's a large structure with circular openings or pipes — maybe ventilation ducts, tubes, or barrels stacked in a frame. The surroundings include a metal framework.



Conventional Image (NIR)

LMM Vision Reasoning:
The image appears to show a grayscale or infrared view of a scene with dense vegetation in the foreground — mainly leafy plants or shrubs. In the background, there seems to be a structure or series of buildings, possibly with a road or paved area in between.



SA Integral Image (NIR)

LMM Vision Reasoning:
The image appears to be a blurry grayscale or low-resolution photo showing an outdoor scene. It looks like there is some vegetation or bushes in the background, and a person or human-like figure is sitting or crouching on the ground near the center of the image. The surroundings seem natural — possibly a garden, park, or wooded area.

What does ☆ do?

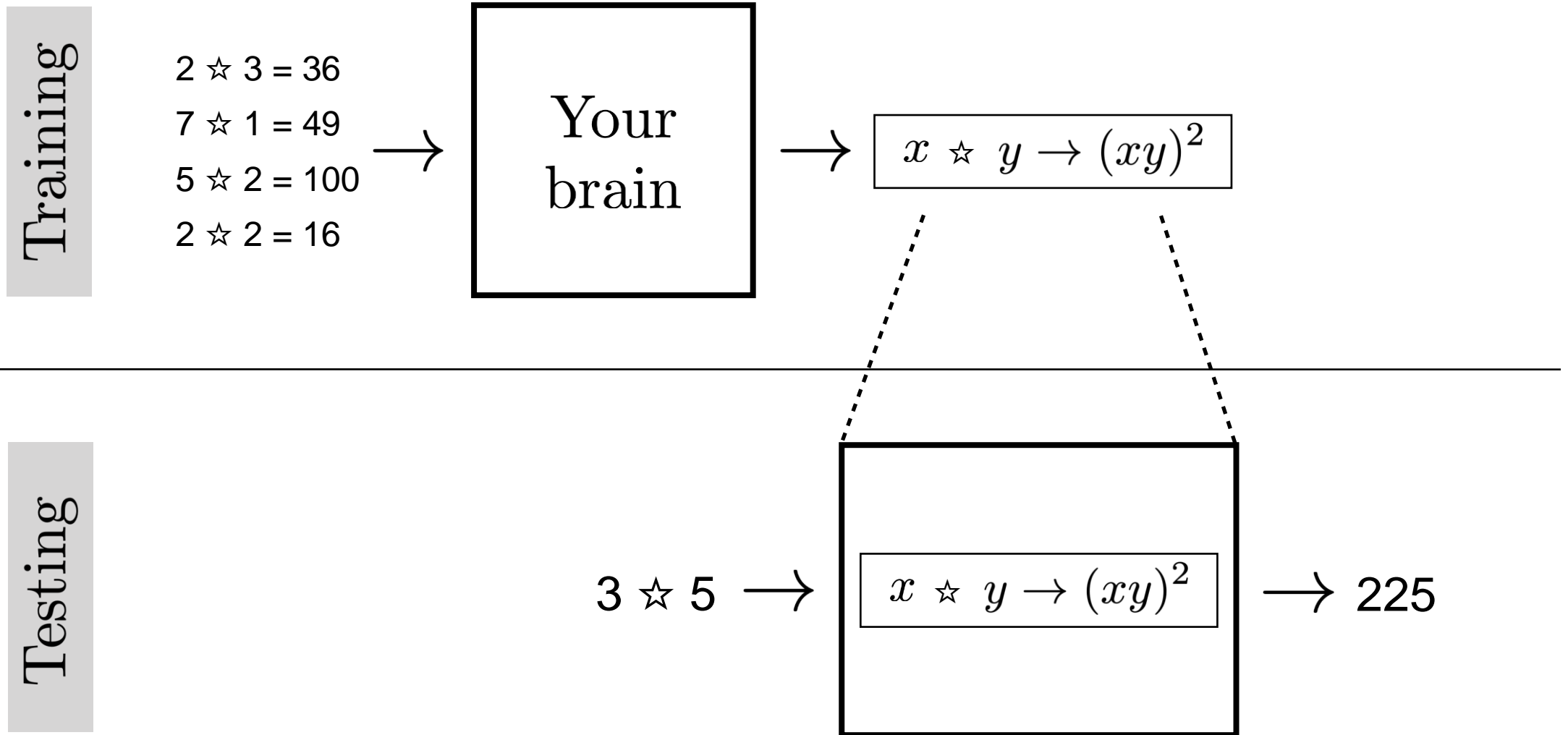
$$2 \star 3 = 36$$

$$7 \star 1 = 49$$

$$5 \star 2 = 100$$

$$2 \star 2 = 16$$

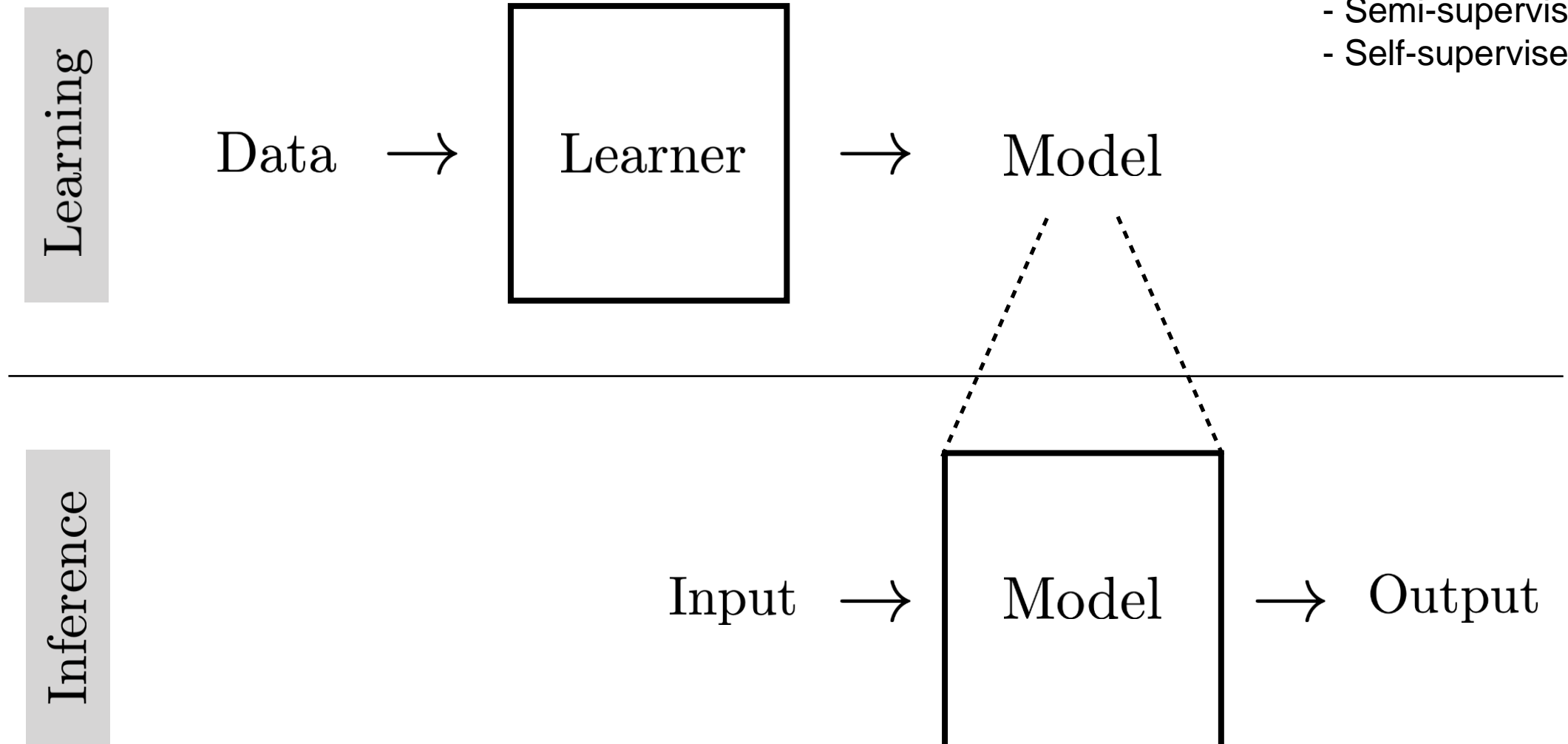
What does \star do?



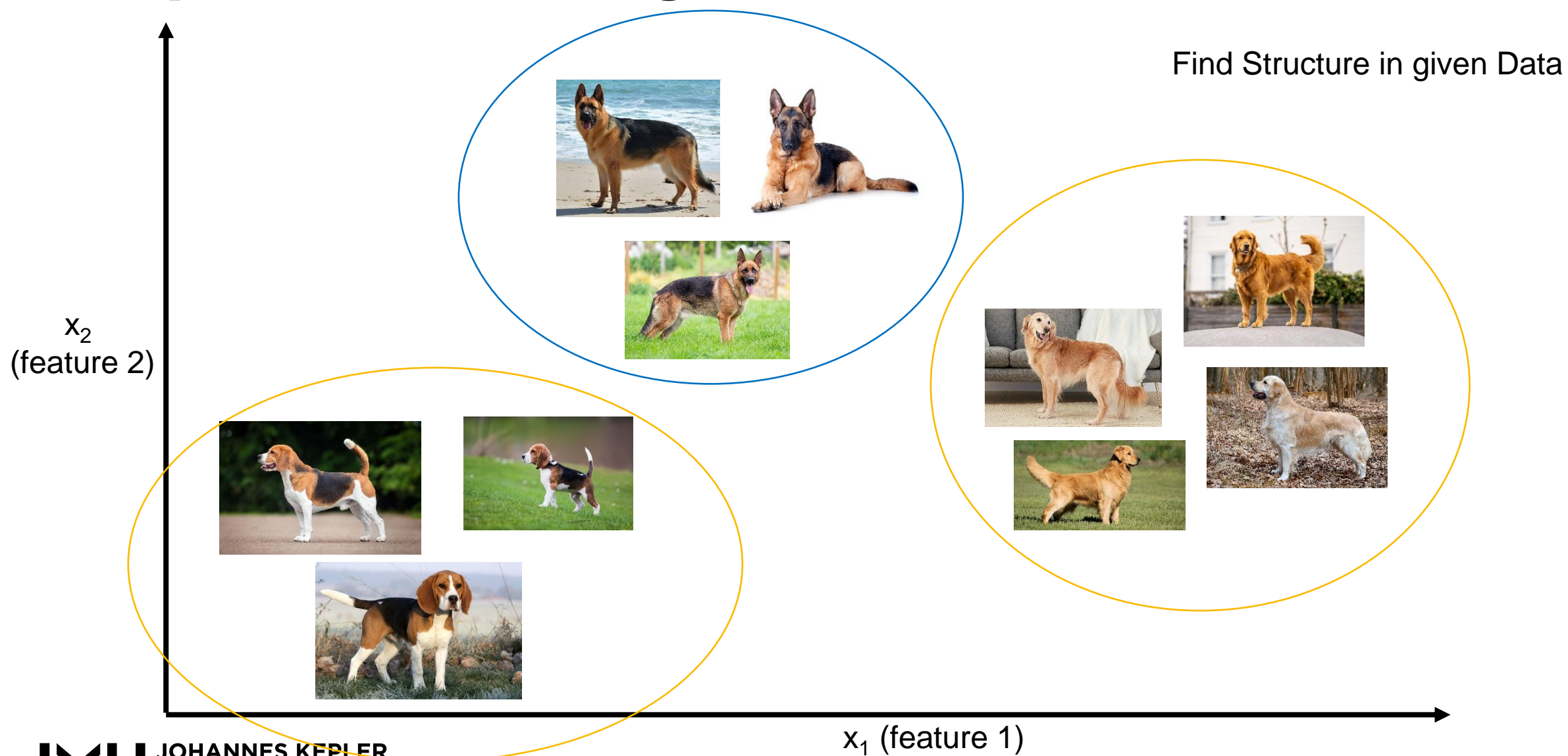
Machine Learning

ML Classes:

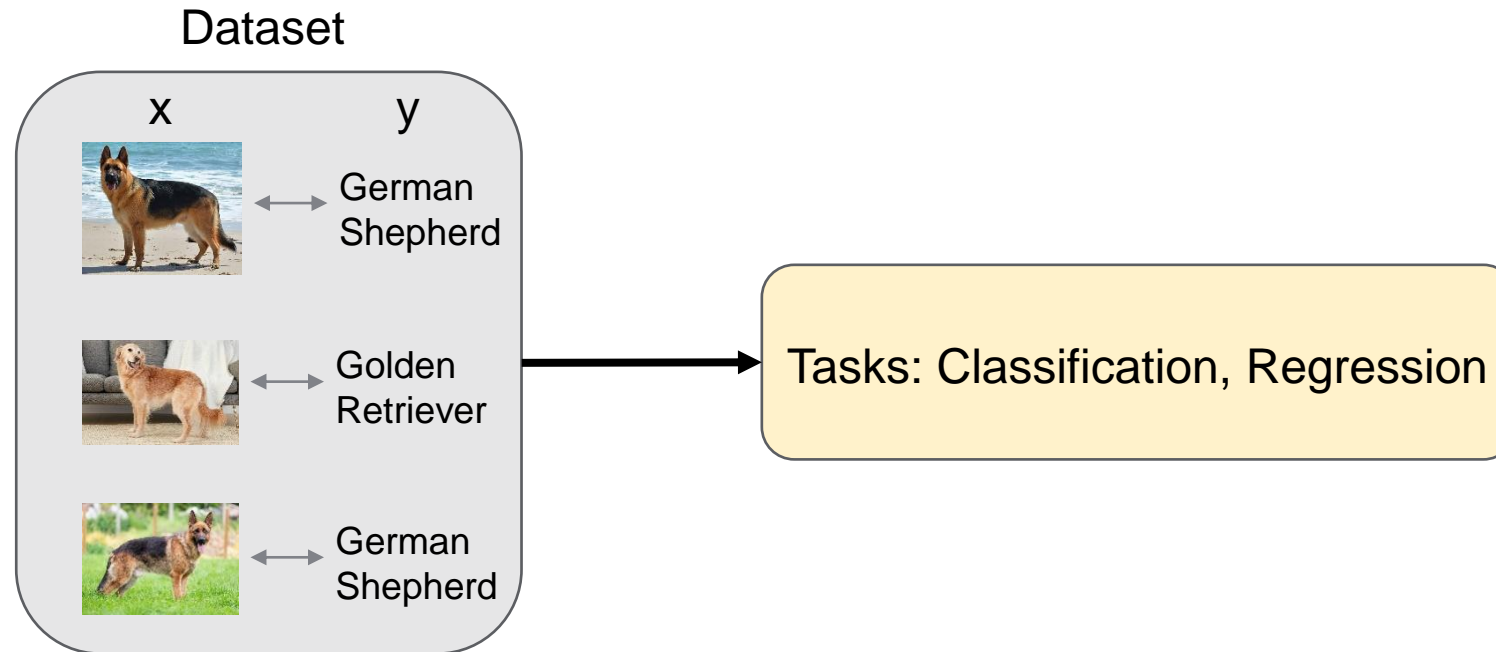
- Unsupervised
- Supervised
- Semi-supervised
- Self-supervised



Unsupervised Learning

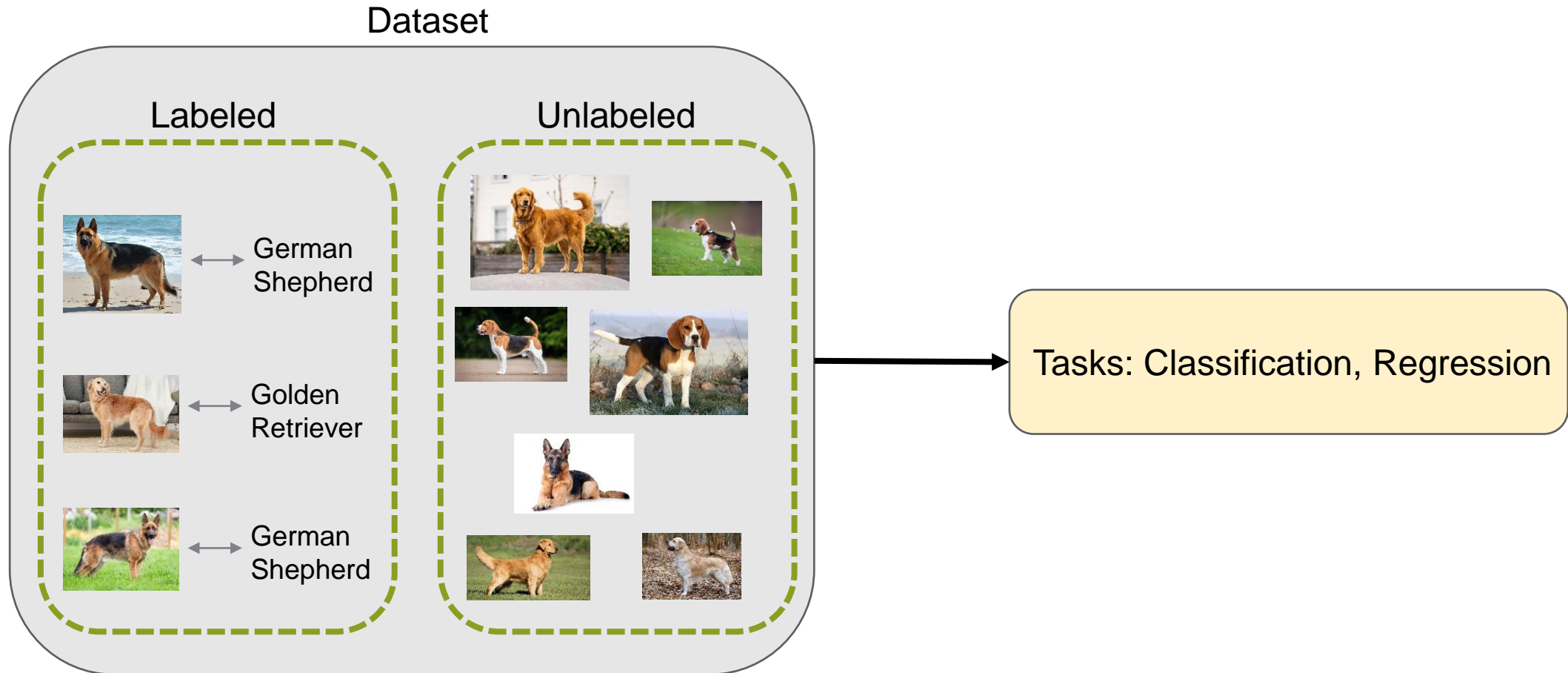


Supervised Learning



Use labeled Data for a Prediction Task.

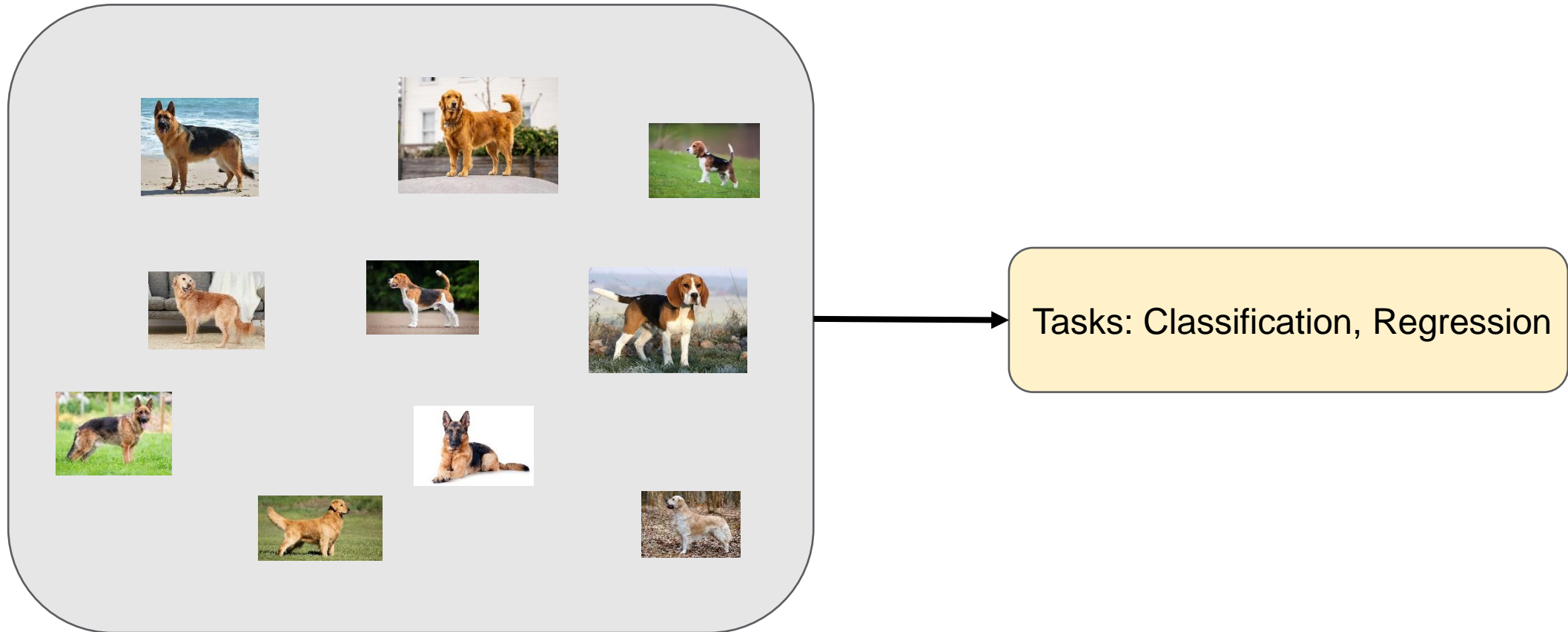
Semi-Supervised Learning



Use labeled Data + unlabeled Data for a Prediction Task.

Self-Supervised Learning

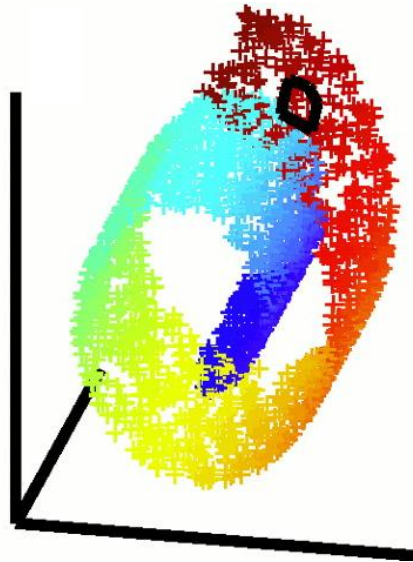
Dataset



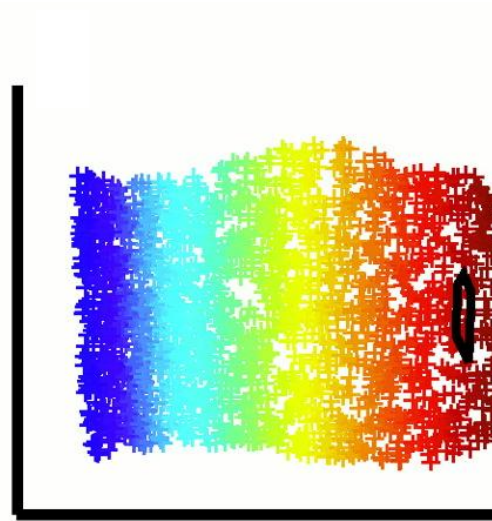
No external Labels. Instead, uses Properties of the Data itself as supervisory Signals. Example (learn how to solve a puzzle): take unlabeled images, rotate them, and train a model to guess the rotation. It needs to understand how the objects look upright → helps to learn how to classify.

Dimensionality Reduction

Given Data Points in d Dimensions, convert them to Data Points in $k < d$ Dimensions with minimal loss of Information.

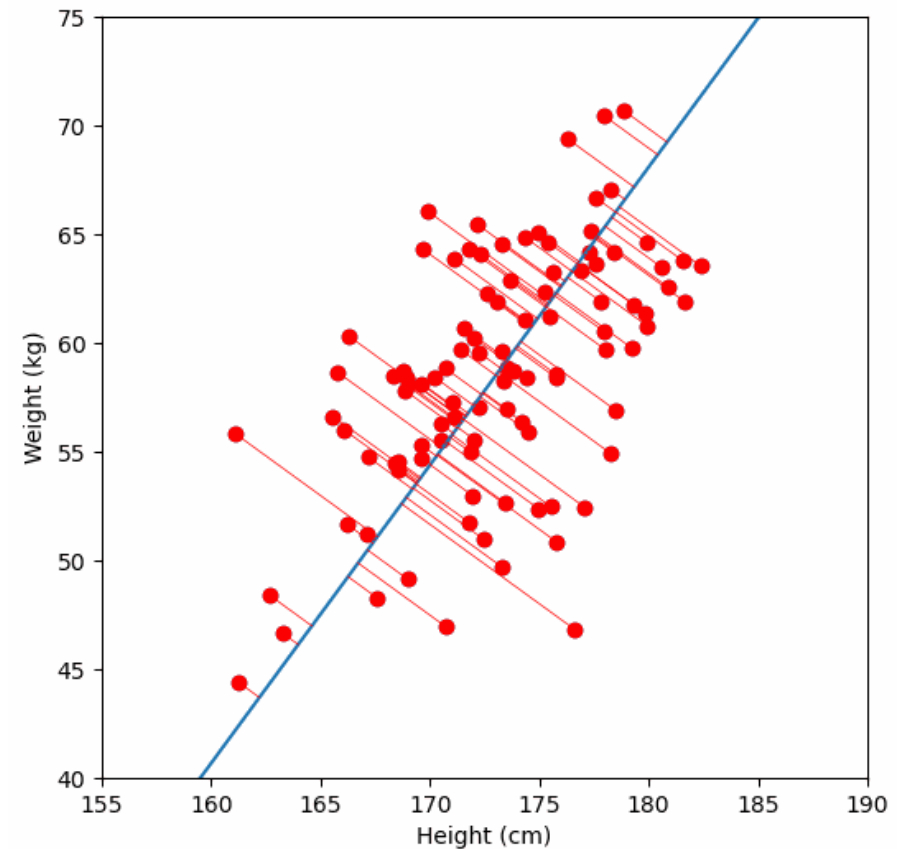
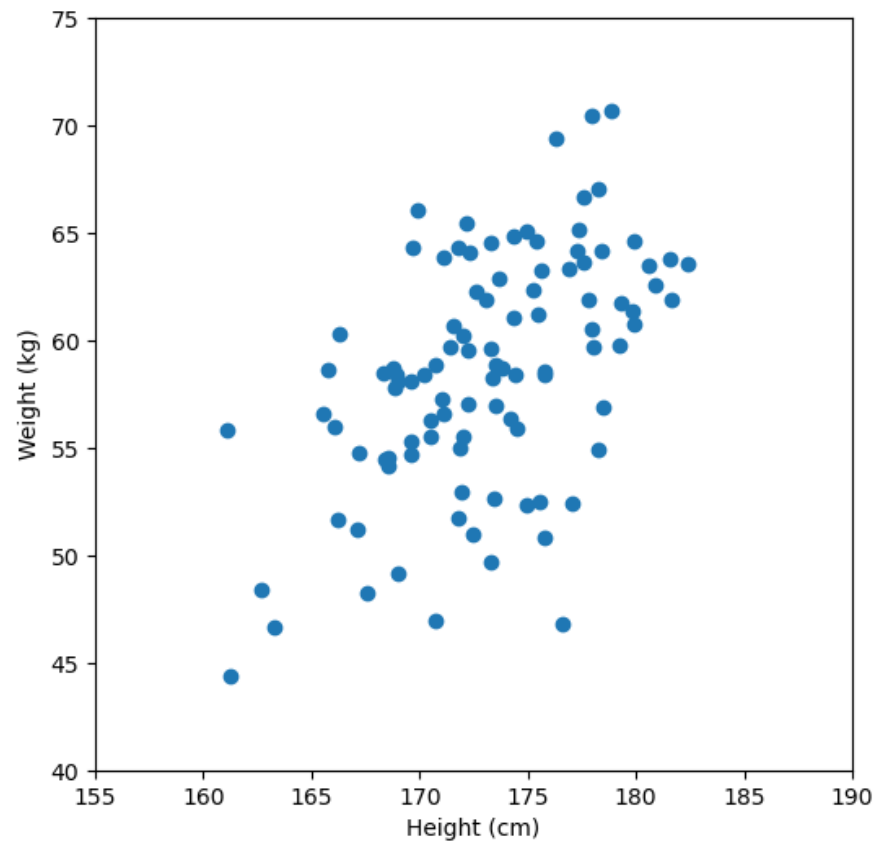


Original
Data in 3D



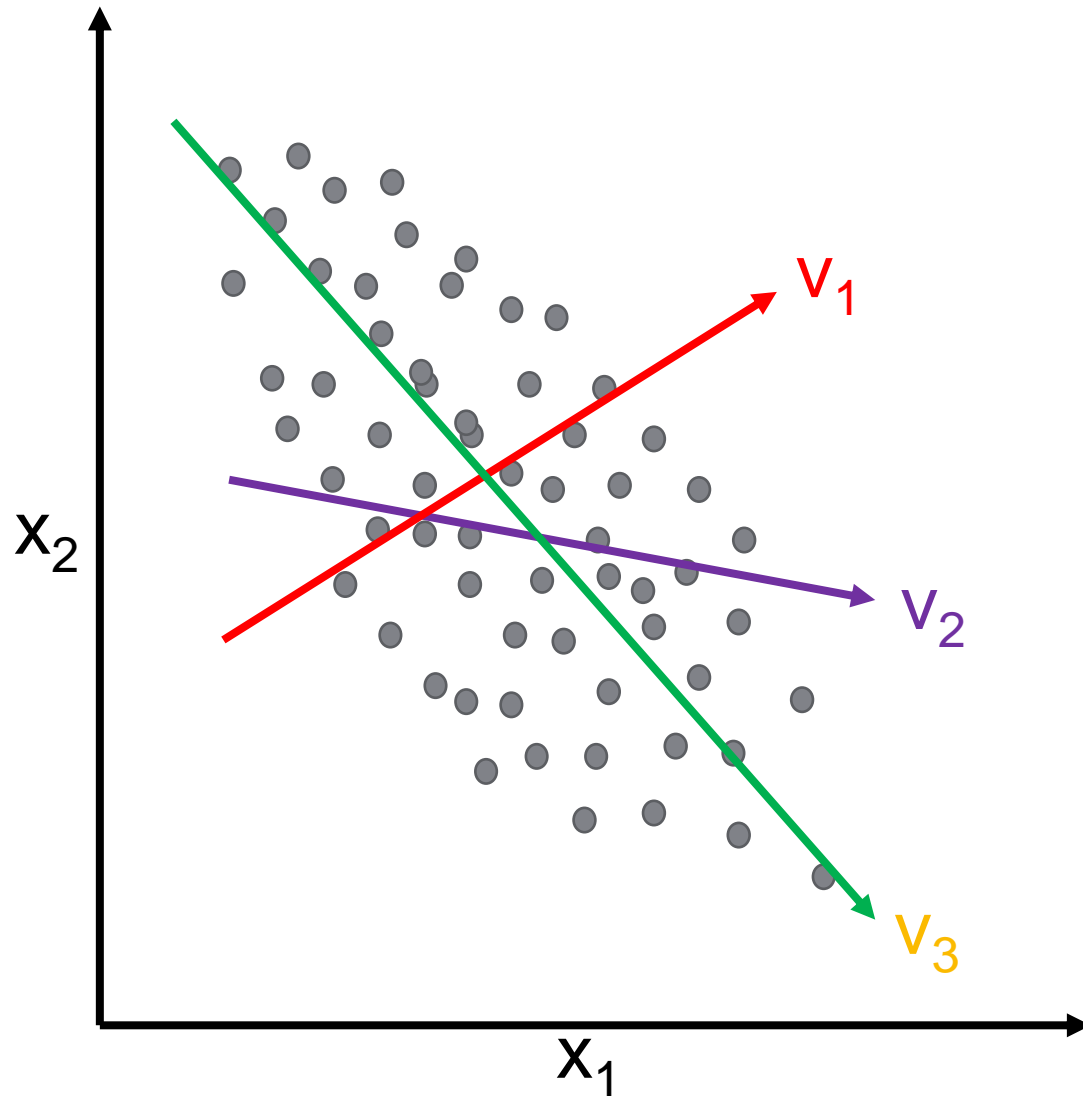
Same data
in 2D

Dimensionality Reduction: Example



Correlation Body Weight vs. Body Size

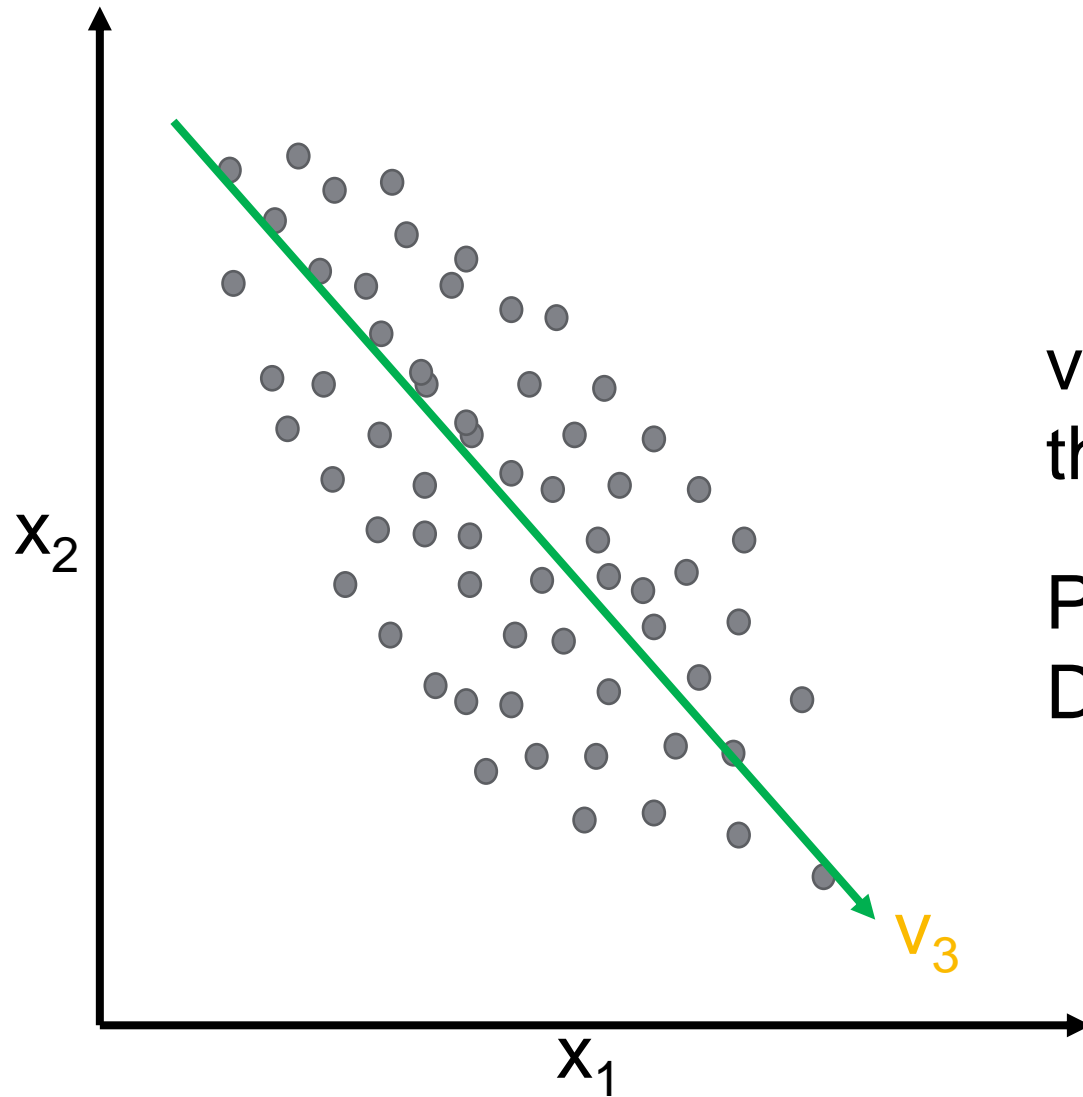
What is the best Direction?



What is the “best” direction to represent this data: V_1 , V_2 , or V_3 ?

Why?

What is the best Direction?



v_3 explains most of the **Variance** in the Data.

Points are most “spread out” in this Direction.

Principle Component Analysis (PCA)

The Variance along \mathbf{v} of all the Point Projections is:

$$\begin{aligned} \text{var}(\mathbf{v}) &= \sum_x \|(x - \bar{x})^T \cdot \mathbf{v}\|^2 \\ &= \sum_x \mathbf{v}^T (x - \bar{x})(x - \bar{x})^T \mathbf{v} \\ &= \mathbf{v}^T \left[\sum_x (x - \bar{x})(x - \bar{x})^T \right] \mathbf{v} \\ &= \mathbf{v}^T \mathbf{C} \mathbf{v} \quad , \text{ where } \mathbf{C} = X^T X \text{ is the} \\ &\quad \text{covariance matrix of } X \end{aligned}$$

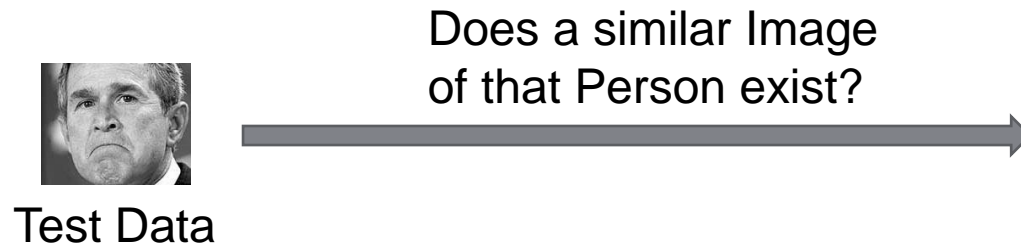
\mathbf{v} = Eigenvectors of \mathbf{C}
 $\text{var}(\mathbf{v})$ = Eigenvalues of \mathbf{v}

PCA is used to compute $\mathbf{v}, \text{var}(\mathbf{v})$

We are interested in \mathbf{v} with the largest $\text{var}(\mathbf{v})$!

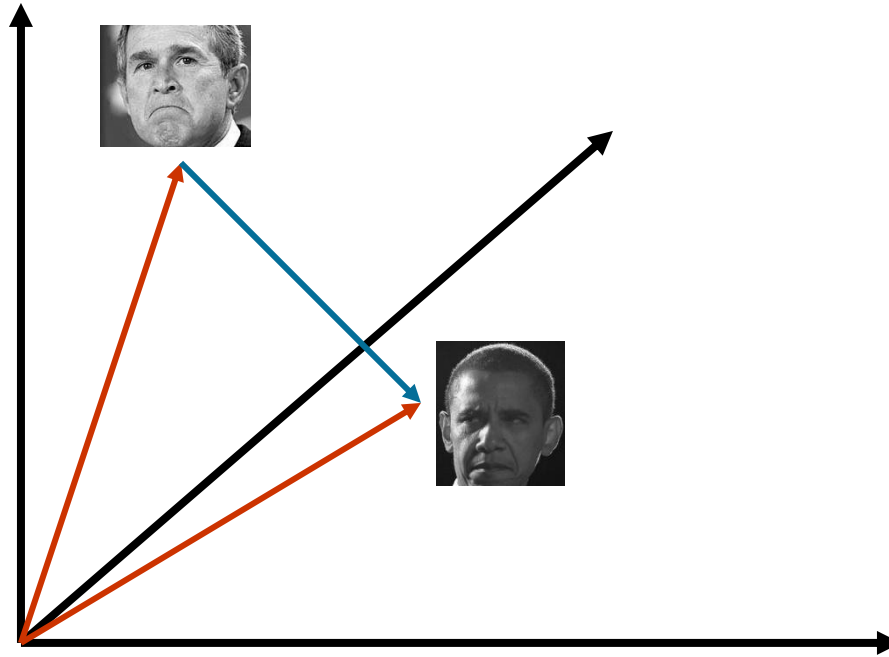
This works also in higher dimensions.

Example: Eigenfaces



Training Data

Example: Eigenfaces



Here, an Image represents a Point in a High Dimensional space

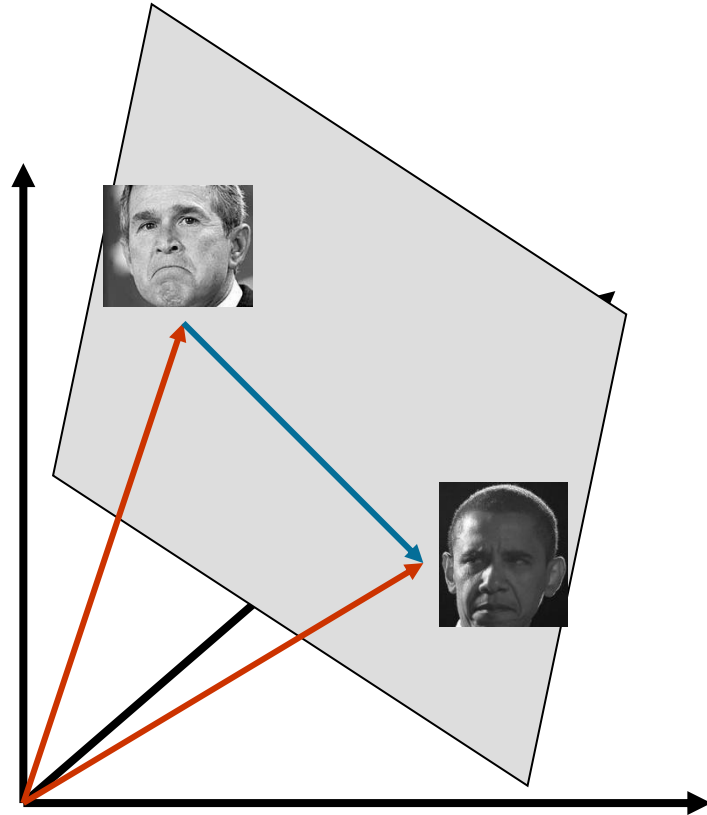
- An $N \times M$ (binary) Image is a Point in \mathbb{R}^{NM}
- We can define Vectors in this Space

Face Space

- When viewed as Vectors of Pixel Values, Face Images are extremely high-dimensional
 - 100 x 100 Image = 10,000 Dimensional
- But very few 10,000-dimensional Vectors are valid Faces
- Let's find a lower-dimensional Representation using PCA!



Dimensionality Reduction



The set of Faces is a “Subspace” of the Set of Images

- Suppose it is K-dimensional
- We can find the best Subspace using PCA
- This is like fitting a “Hyperplane” to the Set of Faces spanned by vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_K$

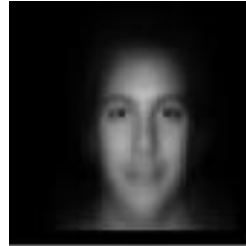
$$\text{Any Face: } \mathbf{x} \approx \bar{\mathbf{x}} + a_1 \mathbf{v}_1 + a_2 \mathbf{v}_2 + \dots + a_k \mathbf{v}_k$$

Eigenfaces

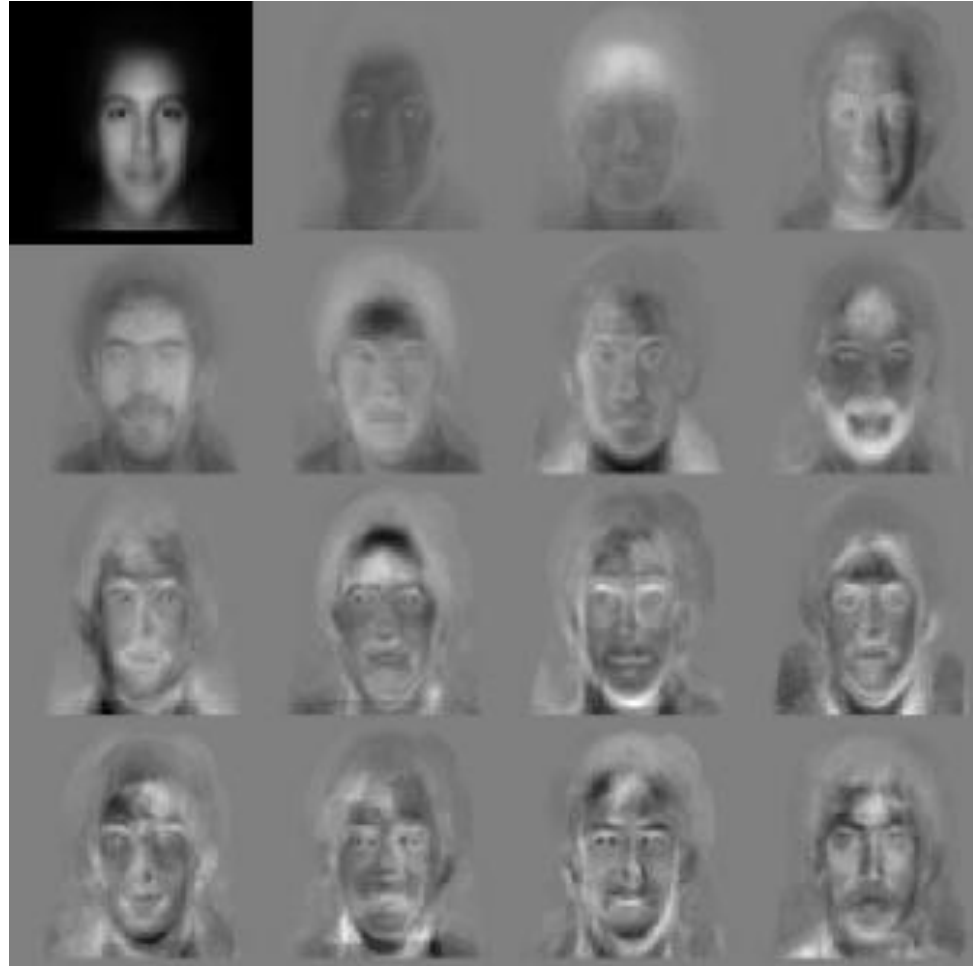
PCA extracts a set of vectors $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \dots$

Each one of these Vectors is a Direction in Face space. What do these look like?

Mean face



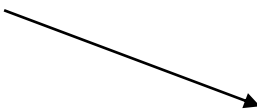
Eigenvectors



Projecting onto Eigenfaces

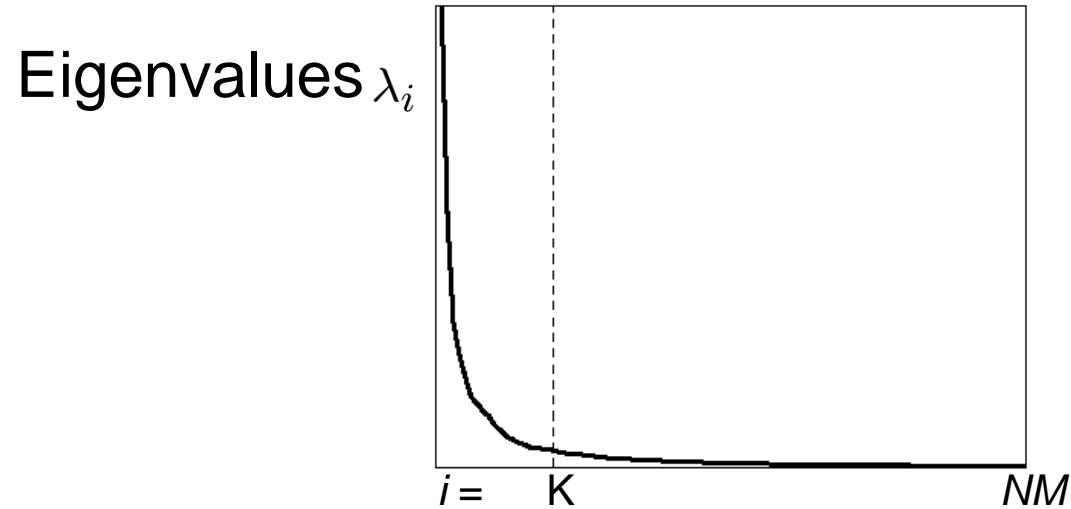
The Eigenfaces $\mathbf{v}_1, \dots, \mathbf{v}_K$ span the Space of Faces. A Face is converted to Eigenface coordinates by

$$\mathbf{x} \rightarrow \left(\underbrace{(\mathbf{x} - \bar{\mathbf{x}}) \cdot \mathbf{v}_1}_{a_1}, \underbrace{(\mathbf{x} - \bar{\mathbf{x}}) \cdot \mathbf{v}_2}_{a_2}, \dots, \underbrace{(\mathbf{x} - \bar{\mathbf{x}}) \cdot \mathbf{v}_K}_{a_K} \right)$$

$$\mathbf{x} \approx \bar{\mathbf{x}} + a_1 \mathbf{v}_1 + a_2 \mathbf{v}_2 + \dots + a_K \mathbf{v}_K$$




Choosing the Dimension K



How many Eigenfaces to use? Look at the Decay of the Eigenvalues

- The Eigenvalue tells you the amount of Variance “in the Direction” of that Eigenface
- Ignore Eigenfaces with low Variance

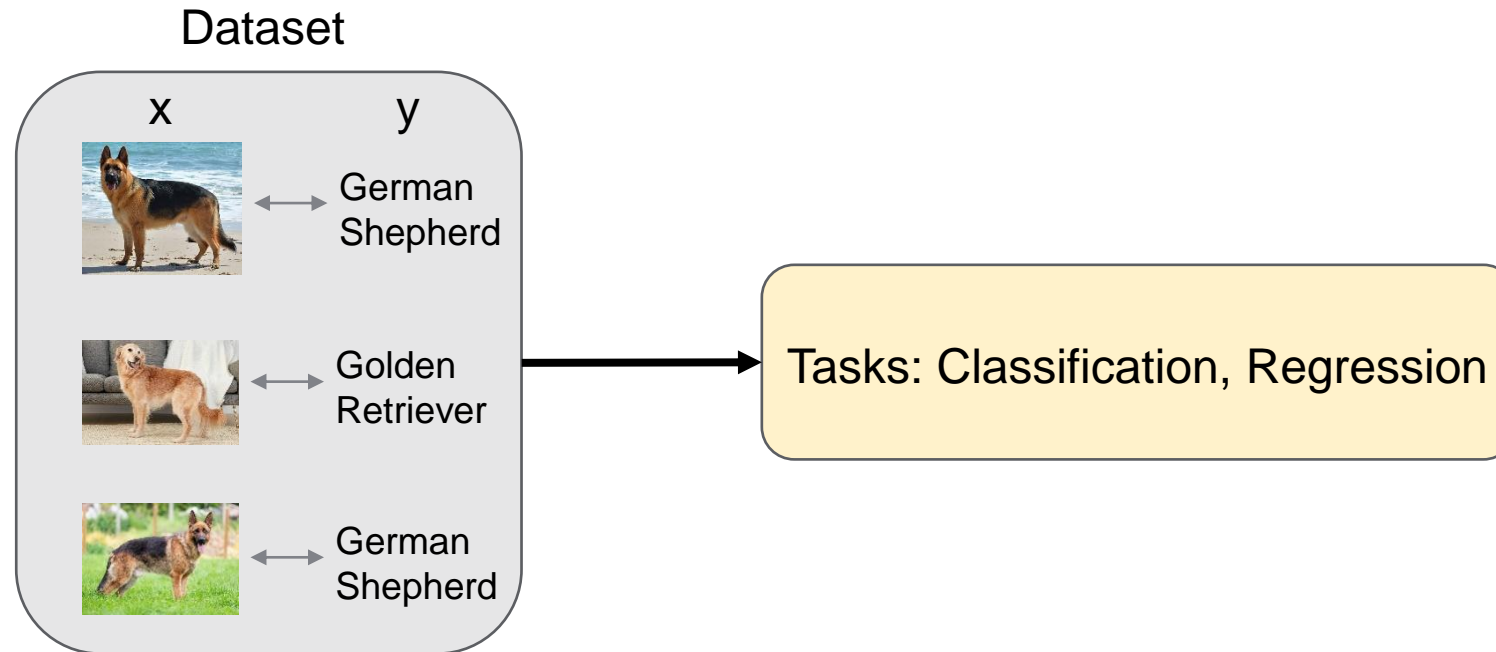
Face Recognition with Eigenfaces

- Project Input Face to Eigenfaces (i.e., represent it by a K-dimensional vector: $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_K$)
- Find the Nearest-Neighbor in Training Image Set (also projected to Eigenfaces)
- Pro: Speed, Size, Fairly Invariant to Resolution
- Cons: Fails quickly if Variation between Training Images and Test Image becomes too large

Dimensionality Reduction Methods

- PCA (Principal Component Analysis):
 - Find Projection that maximize the Variance
- ICA (Independent Component Analysis):
 - Very similar to PCA except that it assumes non-Gaussian features
- Multidimensional Scaling:
 - Find Projection that best preserves Inter-Point Distances
- LDA (Linear Discriminant Analysis):
 - Maximizing the Component Axes for Class-Separation
- ...
- ...

Supervised Learning – Example: Image Classification



Use labeled Data for a Prediction Task.

Learning from Labeled Training Data

Training data

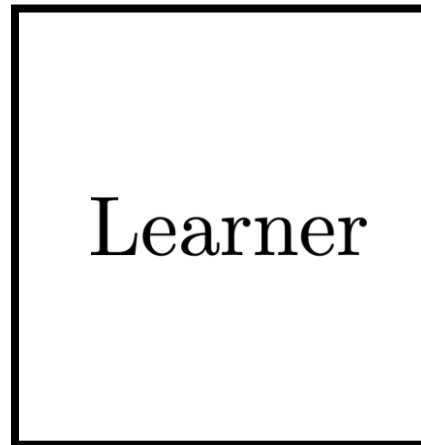
$$\{x^{(1)}, y^{(1)}\}$$

$$\{x^{(2)}, y^{(2)}\}$$

$$\{x^{(3)}, y^{(3)}\}$$

...

→

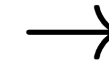
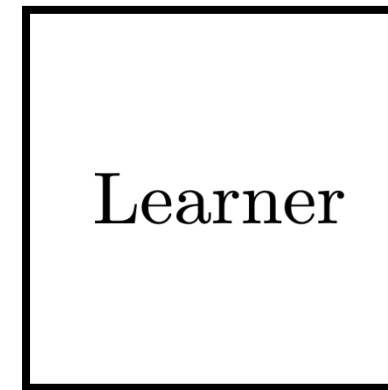
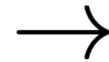
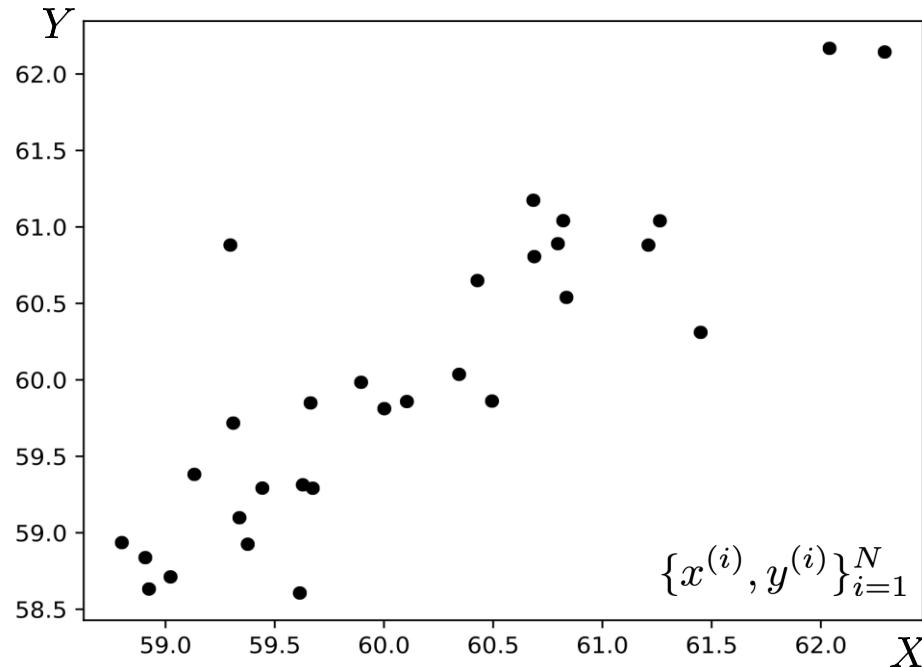


→

$$f : X \rightarrow Y$$

Linear Regression

Training data



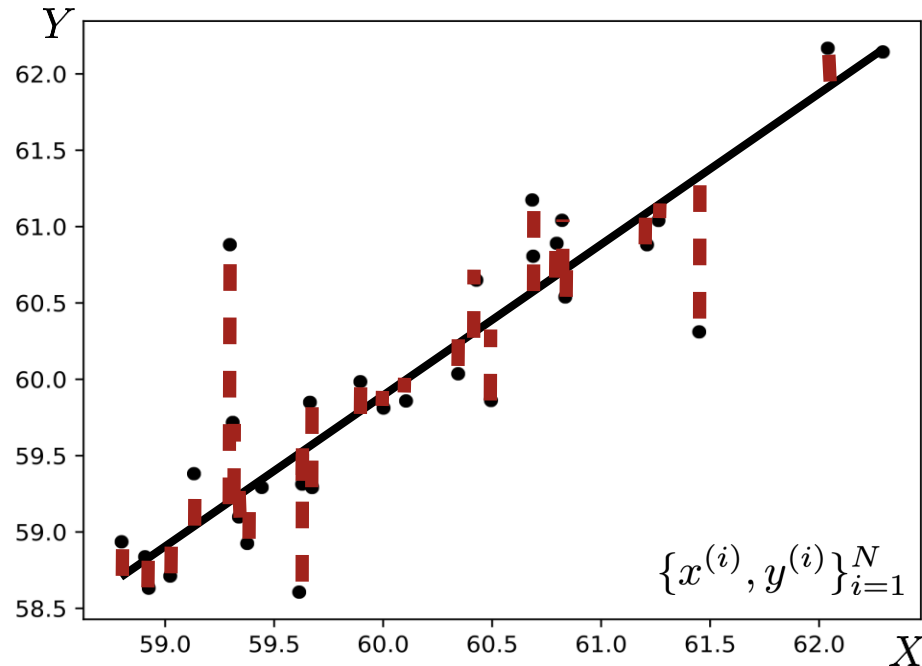
$$f_{\theta}(x) = \theta_1 x + \theta_0$$

Hypothesis space

The relationship between X and Y is roughly linear: $y \approx \theta_1 x + \theta_0$

Linear Regression

Training data



Search for the Parameters, $\theta = \{\theta_0, \theta_1\}$,
that best fit the Data.

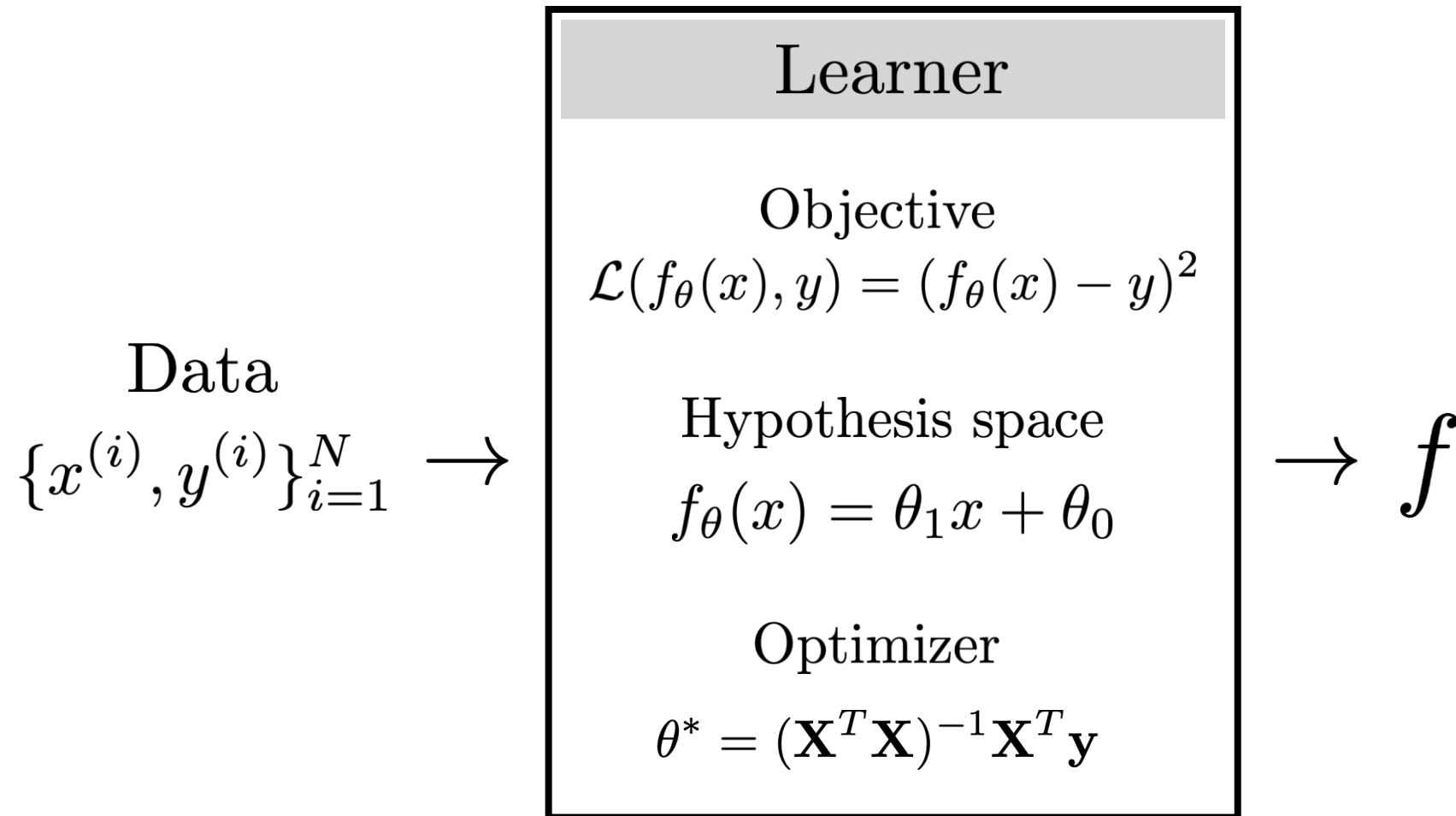
$$f_{\theta}(x) = \theta_1 x + \theta_0$$

Best Fit in what sense?

The Least-Squares **Objective** (aka **Loss**)
says the best Fit is the Function that
minimizes the Squared Error between
Predictions and Target Values:

$$\mathcal{L}(\hat{y}, y) = (\hat{y} - y)^2 \quad \hat{y} \equiv f_{\theta}(x)$$

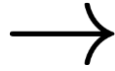
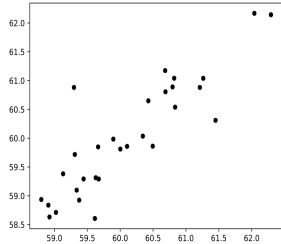
Linear Regression



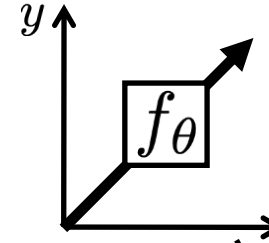
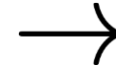
Linear Regression

Training

Data



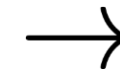
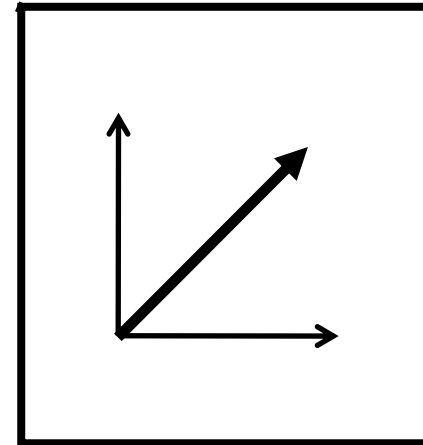
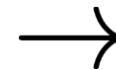
Learner



We could assume more complex Functions (like Polynomial Regression) – but still very limited in terms of **generalization...**

Testing

Input

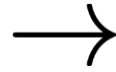
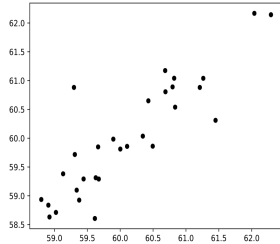


Output

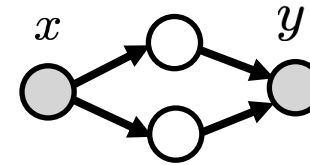
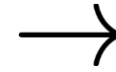
Deep Learning

Training

Data

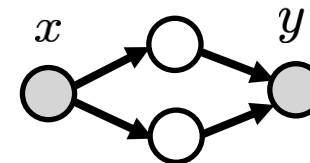
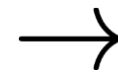


Learner



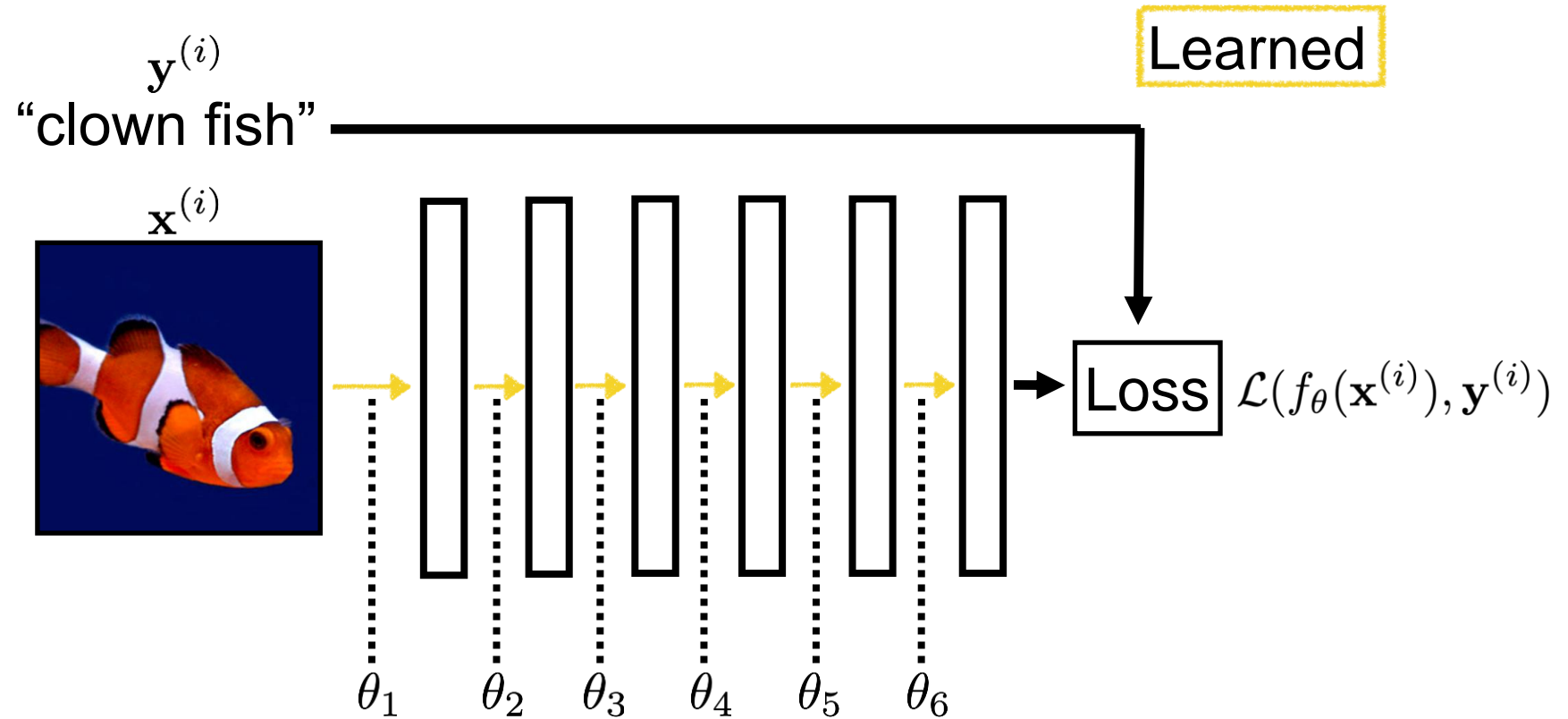
Testing

Input



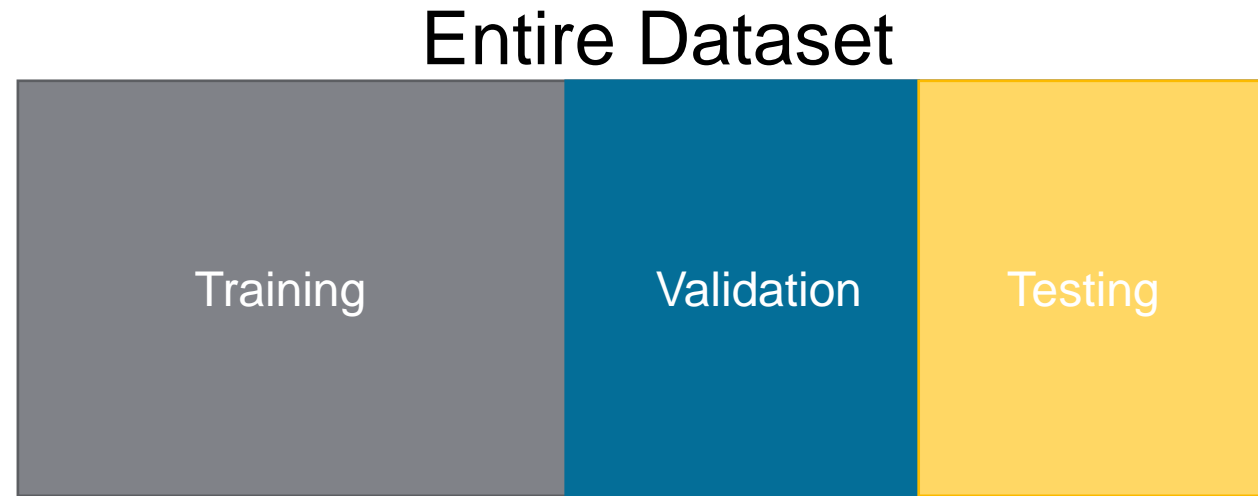
Output

Deep Learning



$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$$

Deep Learning

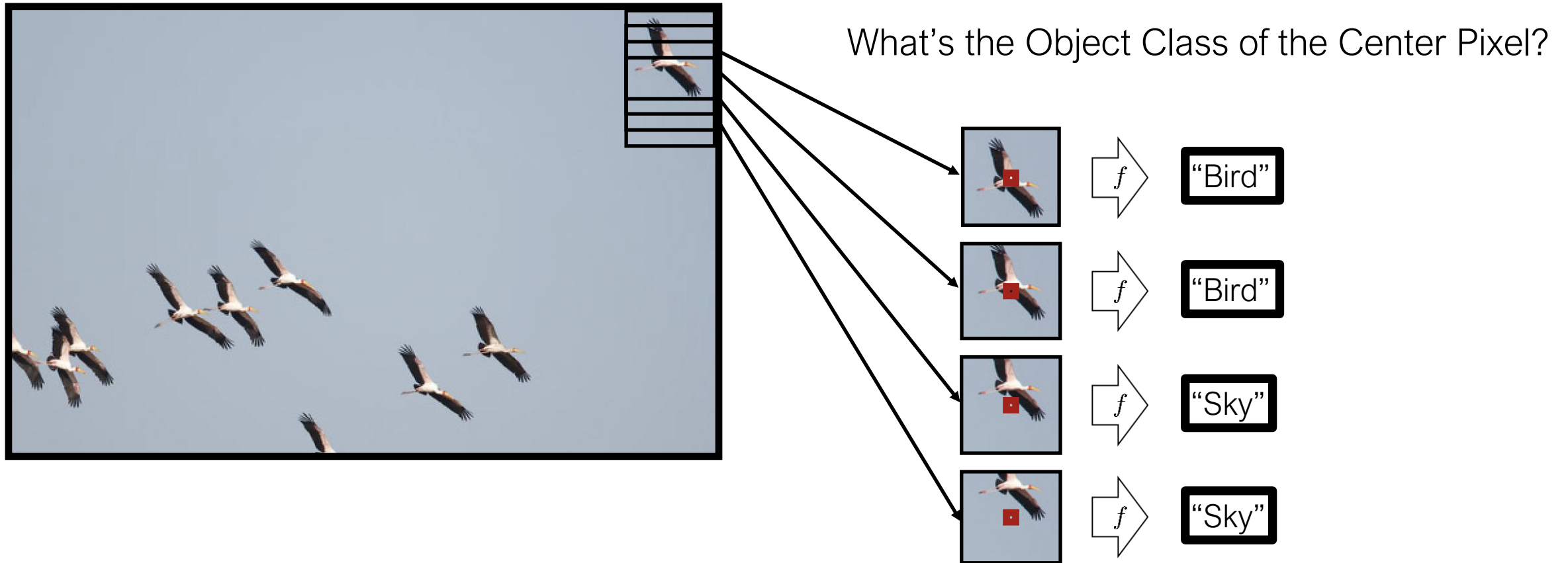


Training: Used to train Model Parameters

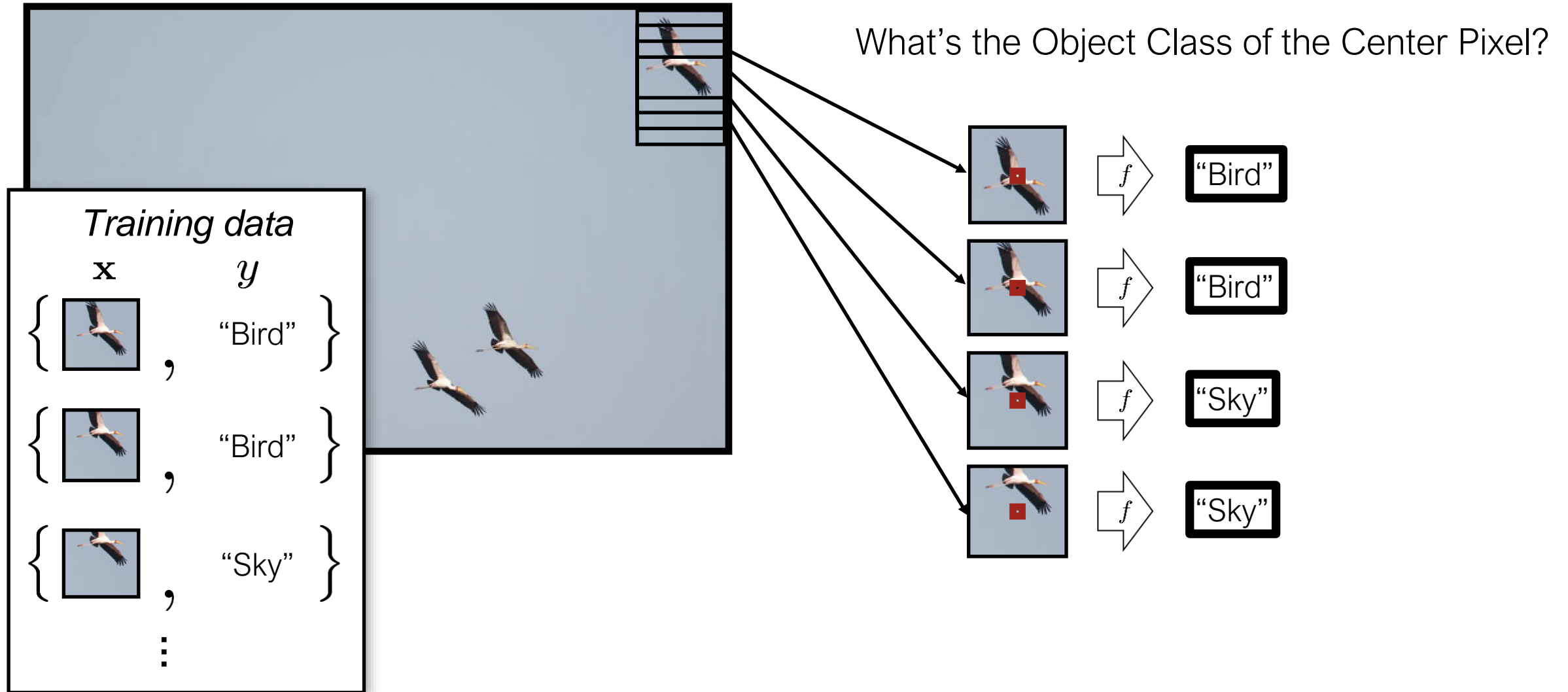
Validation: Used to tune Hyperparameters (e.g., λ in regularized regression)

Testing: Used to report final Accuracy (shouldn't be touched before!)

Classification with Convolution



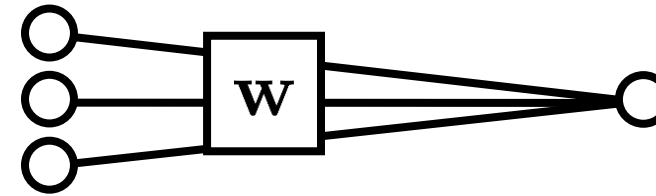
Classification with Convolution



Classification with Convolution



W computes a weighted Sum of all Pixels in the Patch

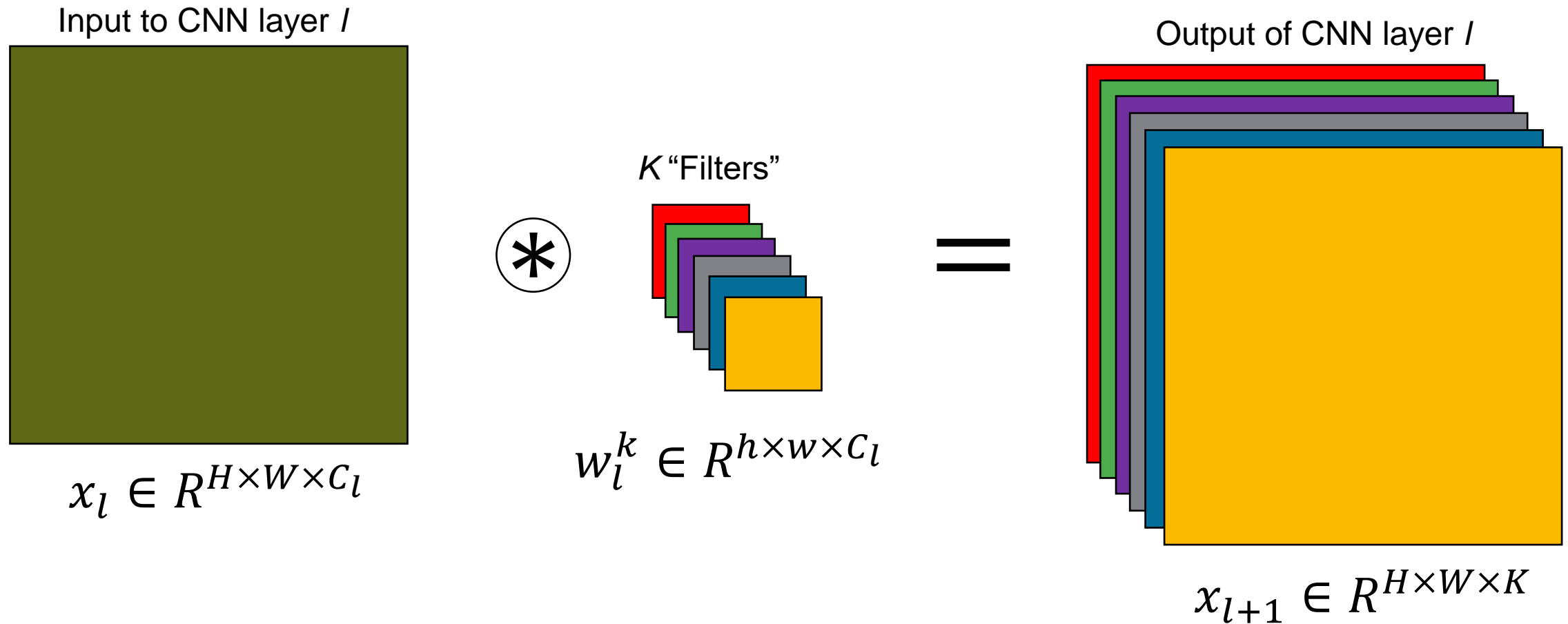


W is a Convolutional Kernel applied to the Full Image!

It gives highest Response if the Kernel matches the Image Content

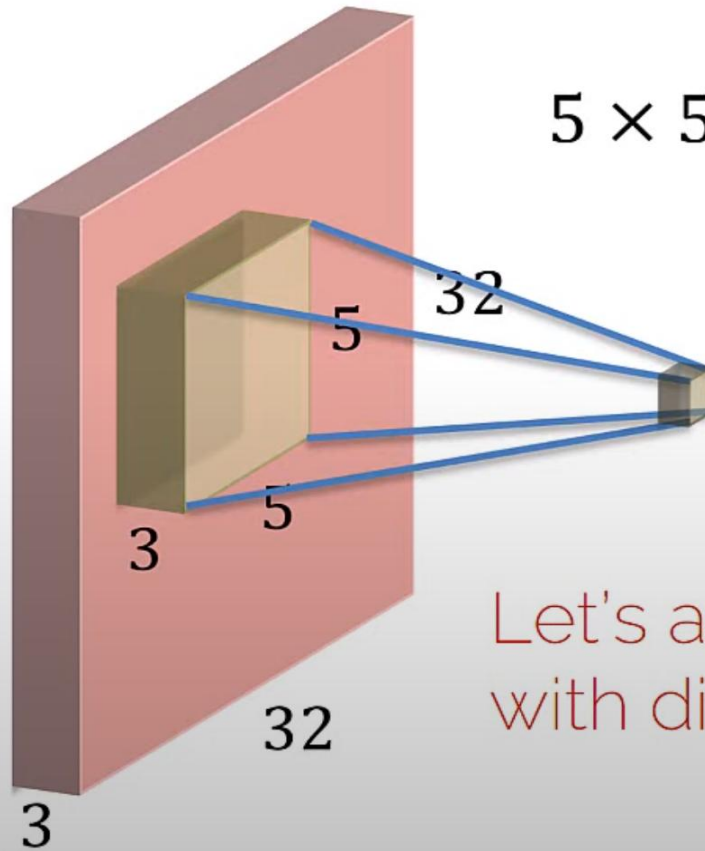
It is Shift Invariant. But how about Scale, Rotation, other Poses, etc.?

Convolutional Neural Networks (CNNs)



Convolutional Layer

$32 \times 32 \times 3$ image

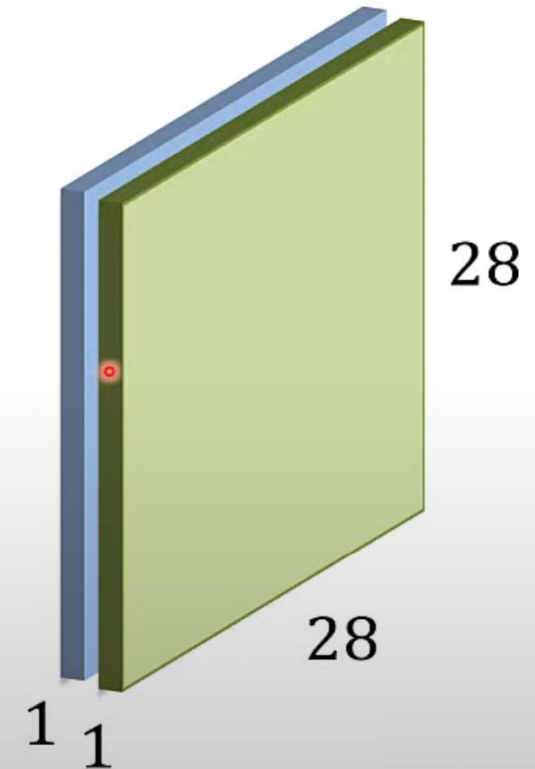


$5 \times 5 \times 3$ filter



Let's apply a different filter
with different weights!

Activation maps



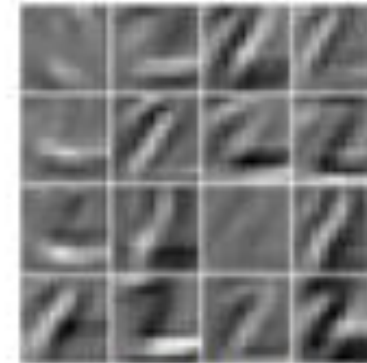
Activation (Feature) Maps



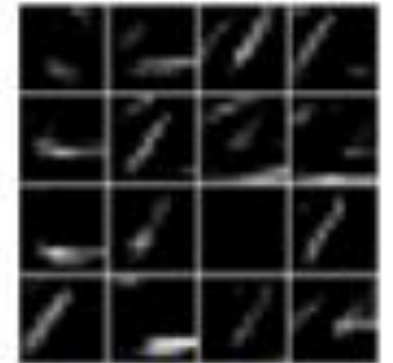
conv1



relu1



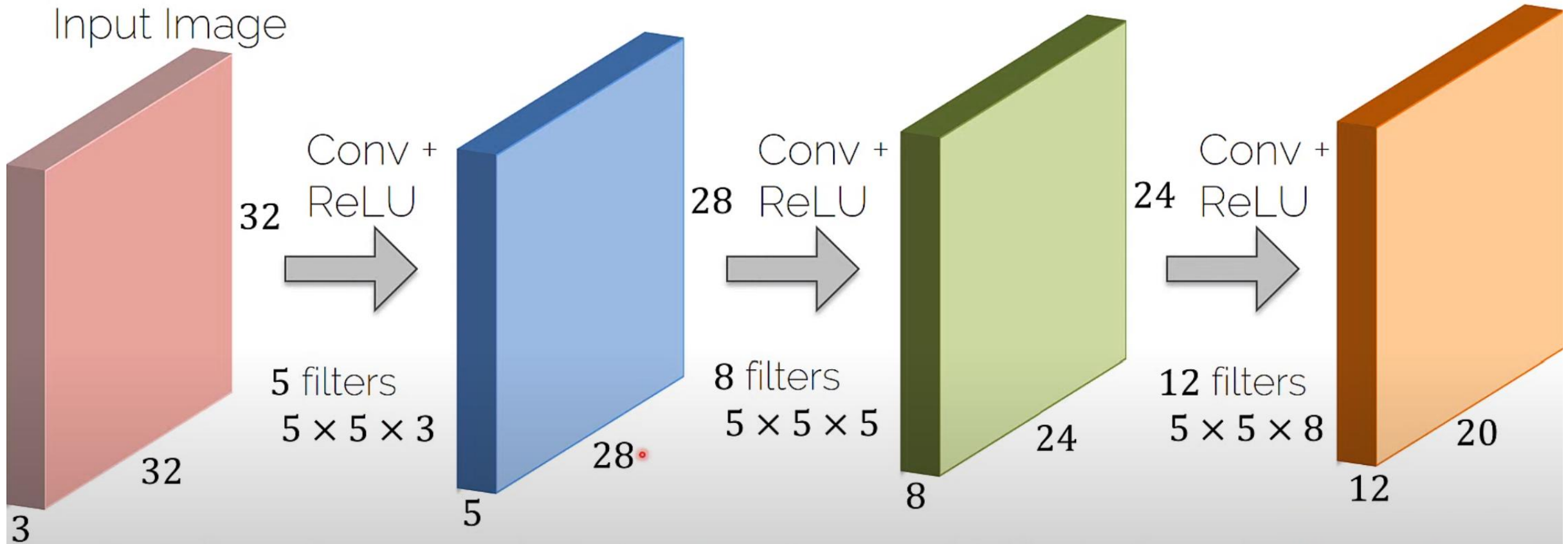
conv2



relu2

- Each Layer can be thought of as a set of C **Feature Maps** aka **Channels**
- Each Feature Map is an NxM Image
- **ReLU** (Rectified Linear Unit, $f(x)=\max(0,x)$) Activation Function is used to introduce Nonlinearity in a Neural Network, helping mitigate the Vanishing Gradient Problem

Multiple Convolutional Layers

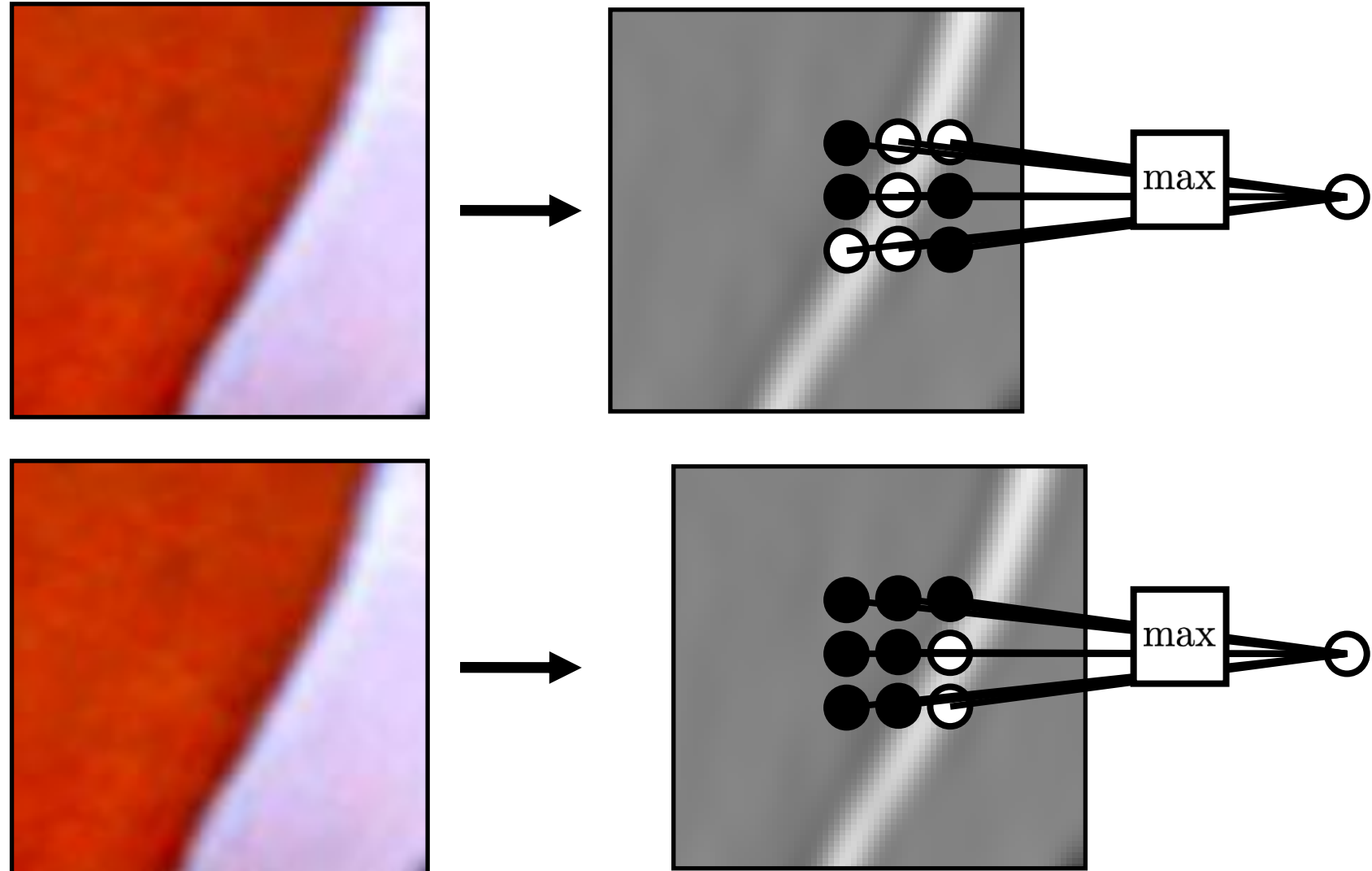


Pooling

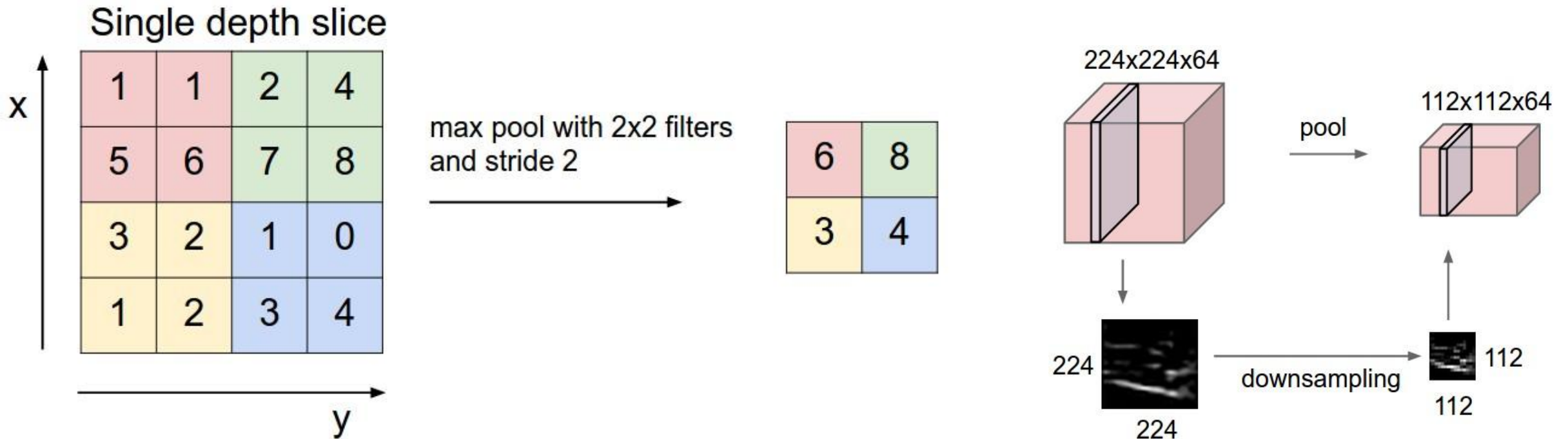
Max pooling
(usually 2x2 or 3x3)

$$y_j = \max_{j \in \mathcal{N}(j)} h_j$$

provides large Response
regardless of exact Position of
Edge

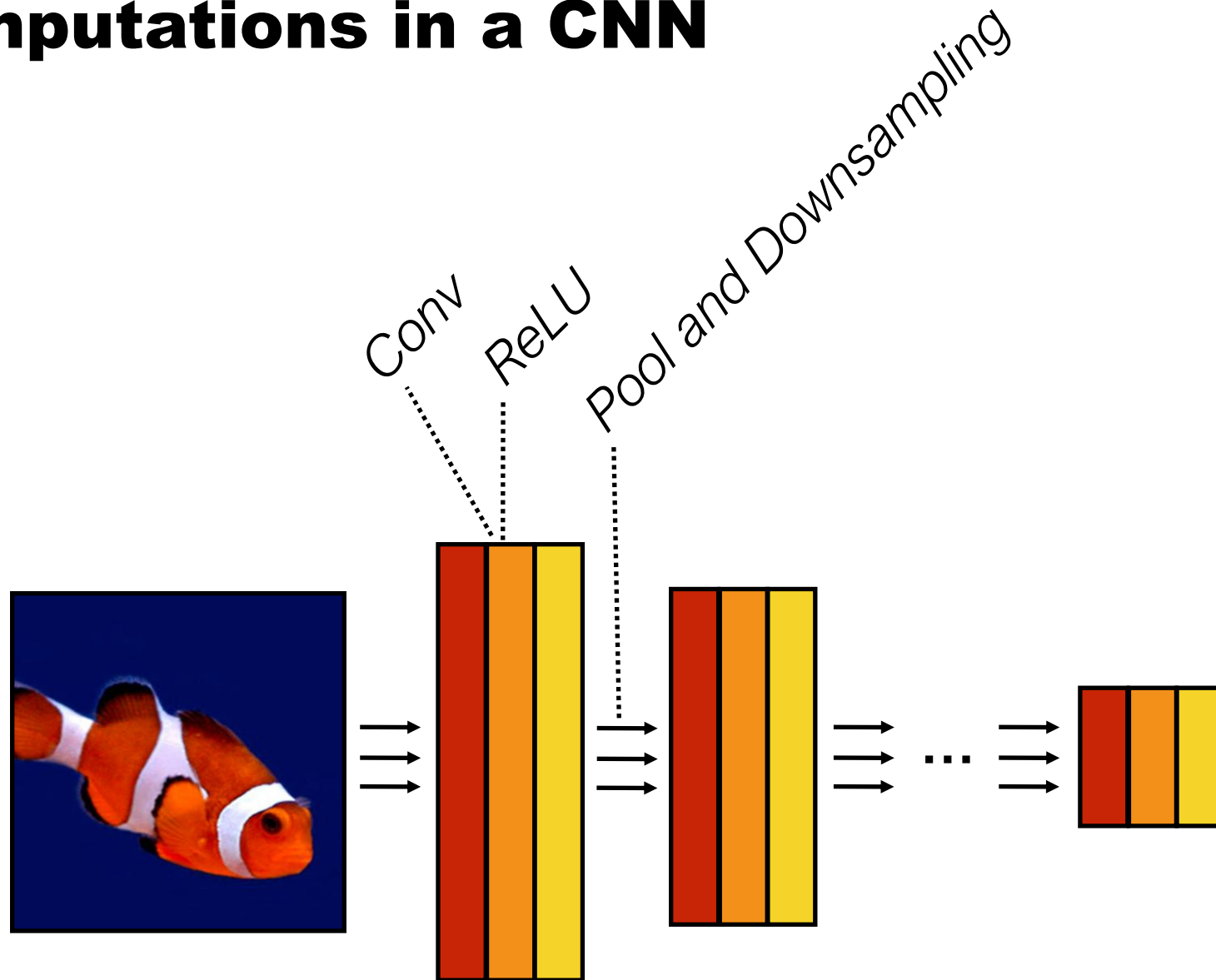


Pooling and Down-Sampling

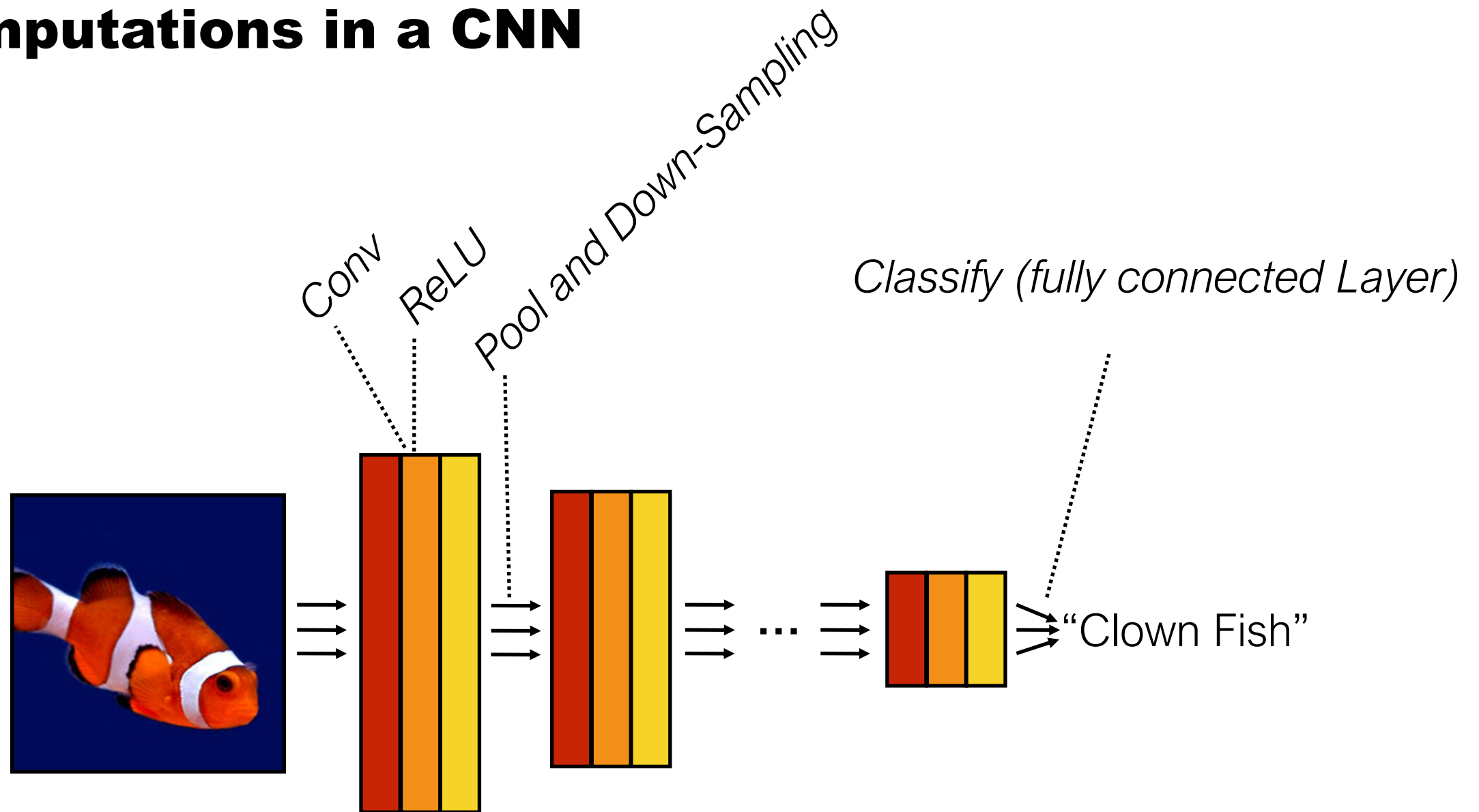


Stride: shift of kernel (for normal Convolution it is 1)

Computations in a CNN

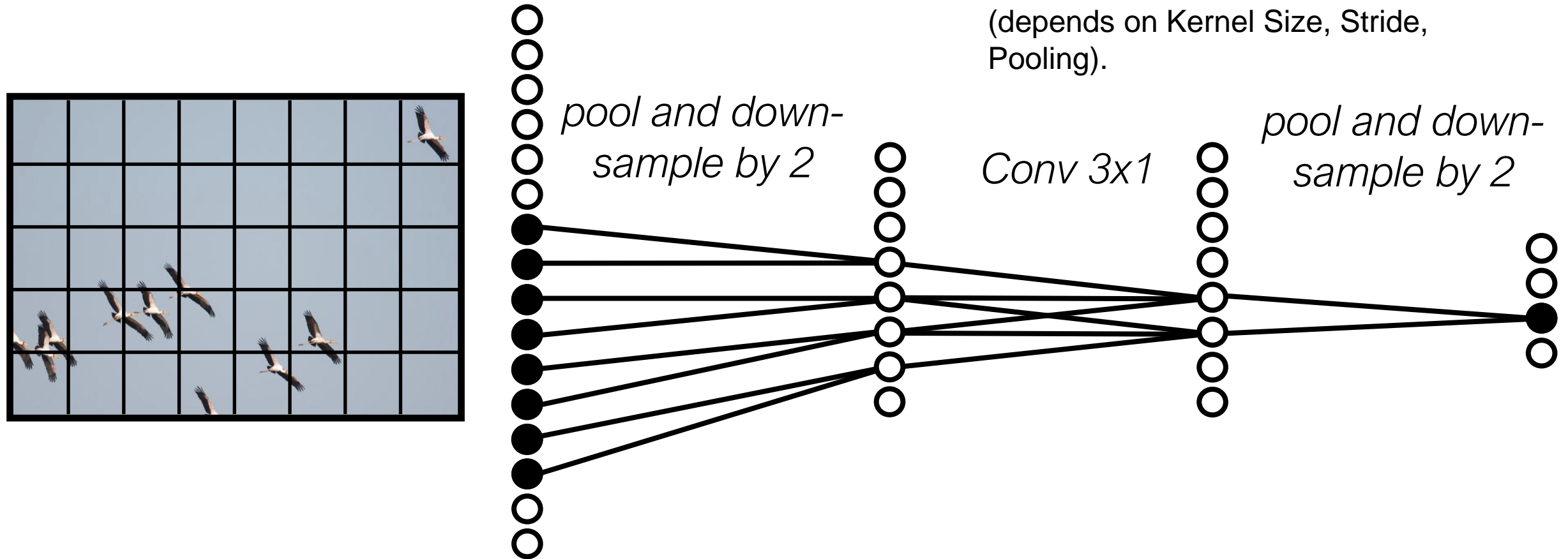


Computations in a CNN



Receptive Field

Receptive Field is the Region of the Image that a particular Neuron in a Convolutional Layer is “looking at” or taking into account when making its Predictions or Feature Extractions (depends on Kernel Size, Stride, Pooling).



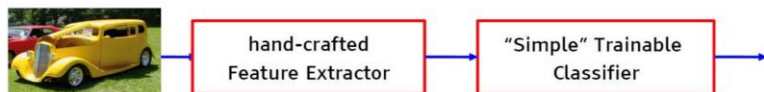
Course Overview

CW	Topic	Date	Place	Lab
41	Introduction and Course Overview	07.10.2025	Zoom	Lab 1
42	Capturing Digital Images	14.10.2025	Zoom	Lab 2
43	Digital Image Processing	21.10.2025	Zoom	Assignment 1
44	Machine Learning	28.10.2025	Zoom	
→ 45	Feature Extraction	04.11.2025	Zoom	Open Lab 1
46	Segmentation	11.11.2025	Zoom	Assignment 2
47	Optical Flow	18.11.2025	Zoom	Open Lab 2
48	Object Detection	25.11.2025	Zoom	Assignment 3
49	Multi-View Geometry	02.12.2025	Zoom	Open Lab 3
50	3D Vision	09.12.2025	Zoom	Assignment 4
3	Trends in Computer Vision	13.01.2026	Zoom	
4	Q&A	20.01.2026	Zoom	Open Lab 4
5	Exam	27.01.2026	HS1 (Linz), S1/S3 (Vienna), S5 (Bregenz)	
9	Retry Exam	24.02.2026	tba	

Next Week: Feature Extraction

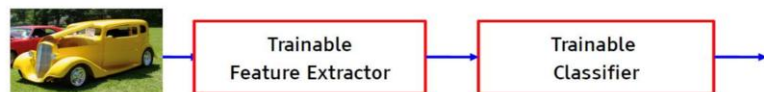
Model-Based vs. Learning-Based Feature Extraction

- Fixed engineered features (or kernels) + trainable classifier

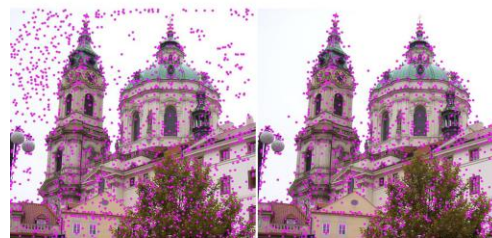


VS.

- End-to-end learning / feature learning / deep learning



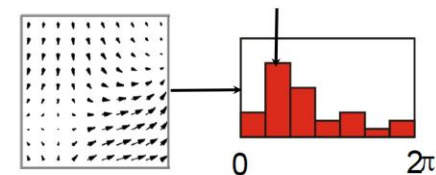
Example: Scale-Invariant Feature Transform (SIFT)



unfiltered

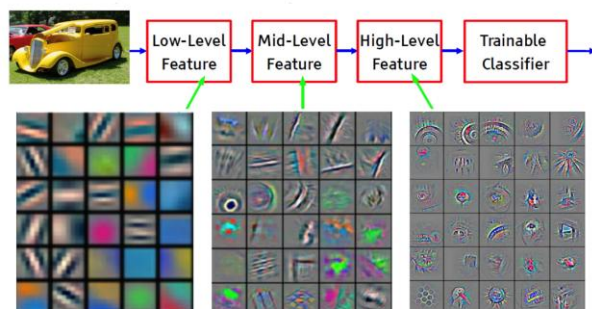
filtered

- (3) Filter out Features in low-contrast Regions (Noise)



- (4) Determine Feature (i.e., Gradient) Orientations and sort them into Histogram (largest Bin = main Orientation)

From Low-Level to High-Level Features

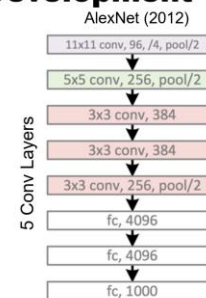


Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

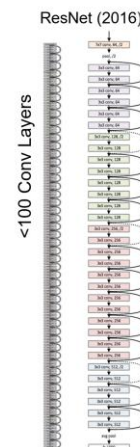
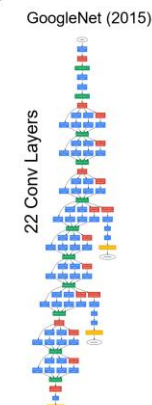
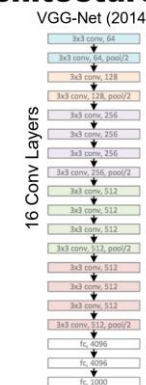


Series 2014

Development of Architectures



...going really deep...



Thank You

