

IEEE ICPRS 2025

[HOME](#) [About](#) [Registration](#) [Travel](#) [Committee](#) [FAQs](#) [Keynotes](#) [Workshops](#) [Sponsors](#) [Previous](#) [Programme](#) [Hotels](#) [Call](#) [Contact](#)



[About The Event](#)

Dr. Rodrigo Salas Fuentes
rodrigo.salas@uv.cl

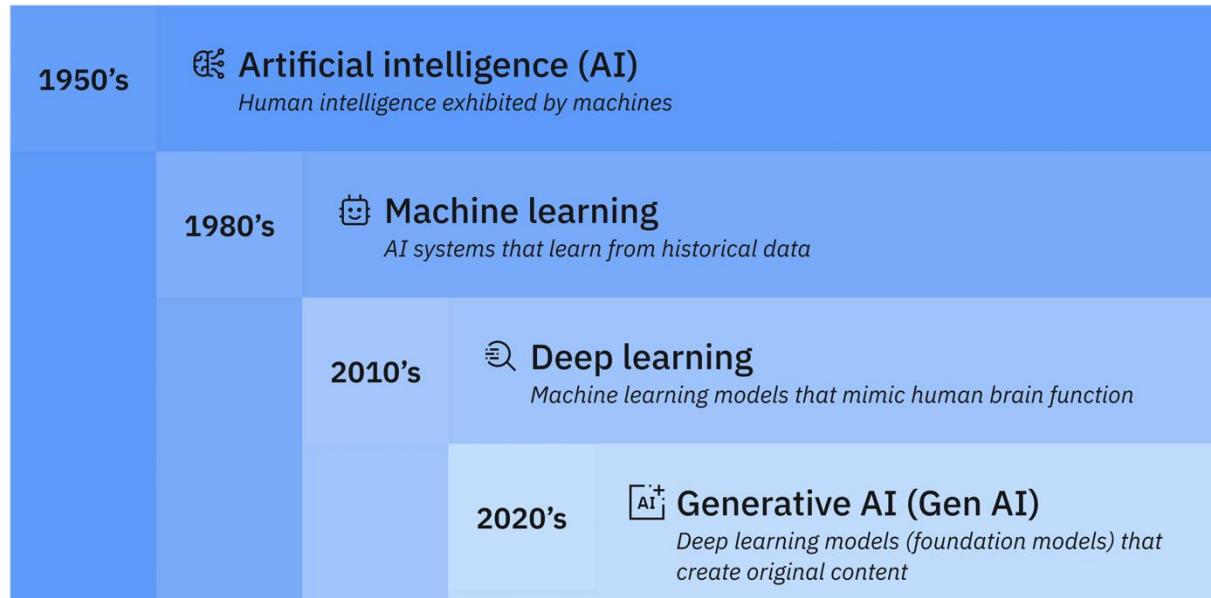


Explainable AI (XAI) for Image Analysis with Deep Learning

Artificial Intelligence

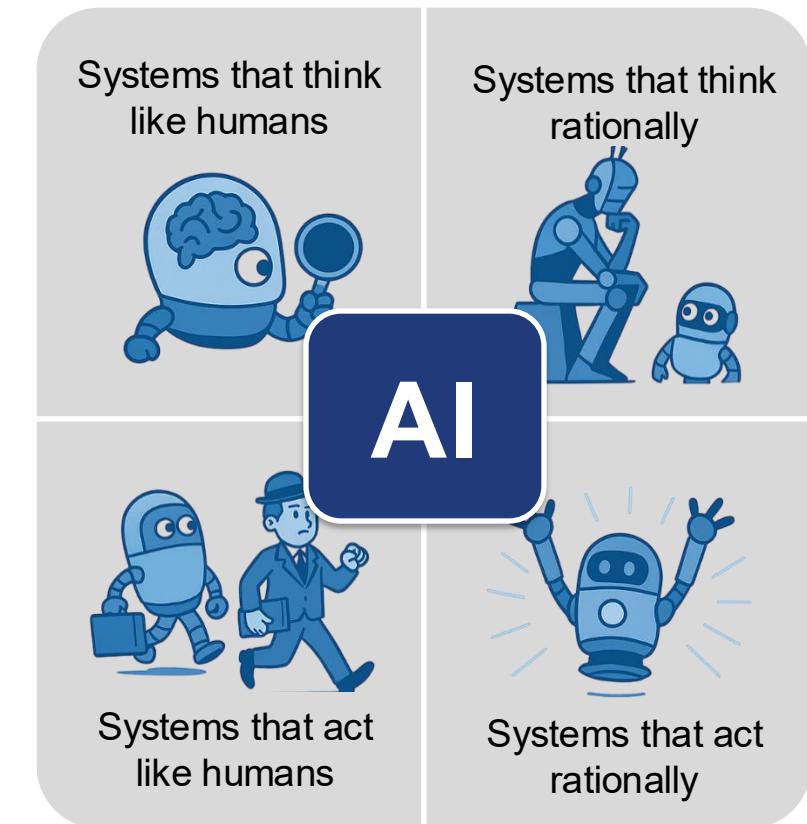
AI

What is Artificial Intelligence



<https://www.ibm.com/mx-es/think/topics/sentient-ai>

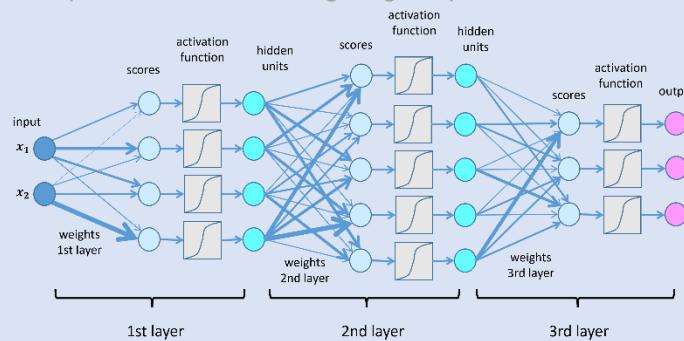
Artificial Intelligence (AI) is technology that enables computers and machines to simulate human learning, comprehension, problem solving, decision making, creativity and autonomy.



Stuart Russell y Peter Norvig, 2020

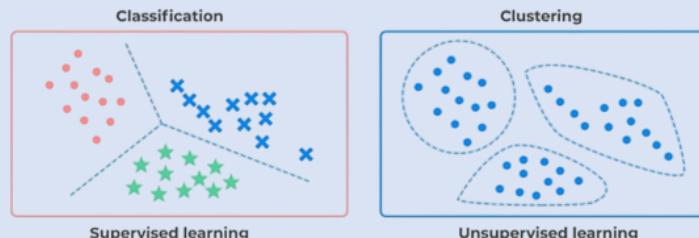
AI Methods

<https://lamarr-institute.org/blog/deep-neural-networks/>



Deep Learning

<https://www.superannotate.com/blog/image-classification-basics>



Machine Learning

<https://devopedia.org/natural-language-processing>

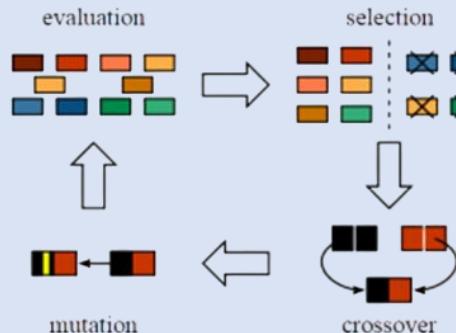
Understanding Language

 "Literally ur facebook message app is useless, you only want it to increase profit. Please fix yourself. Its sad @facebook"

- Emotion: Frustrated
- Tone: Negative, Subjective
- Organization: Facebook
- Product: Messenger App
- Adjectives: "useless", "sad"
- Language: English, Informal

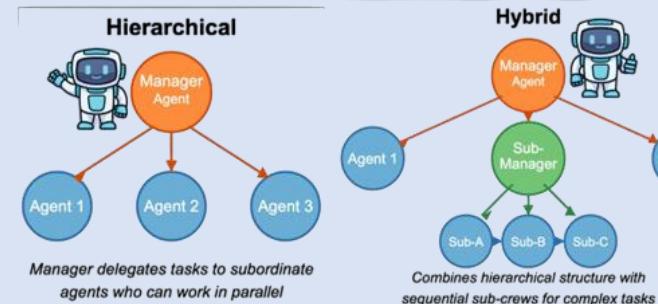
Natural Language Processing

<https://doi.org/10.1145/3446132.3446142>



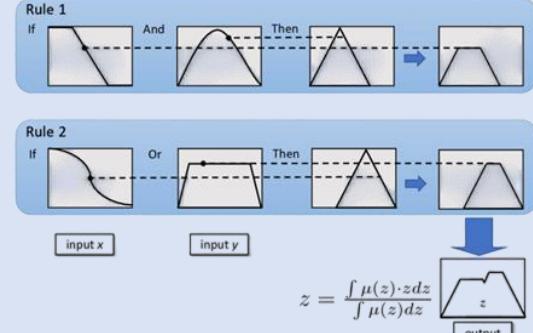
Metaheuristics

<https://l1nq.com/at8FT>



Multi-Agent Systems

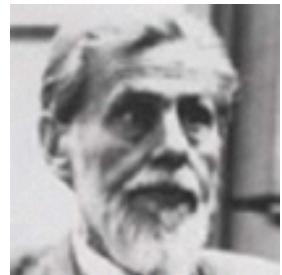
<http://dx.doi.org/10.1016/j.ins.2018.09.005>



Fuzzy Inference Systems

Artificial Neural Networks ANN

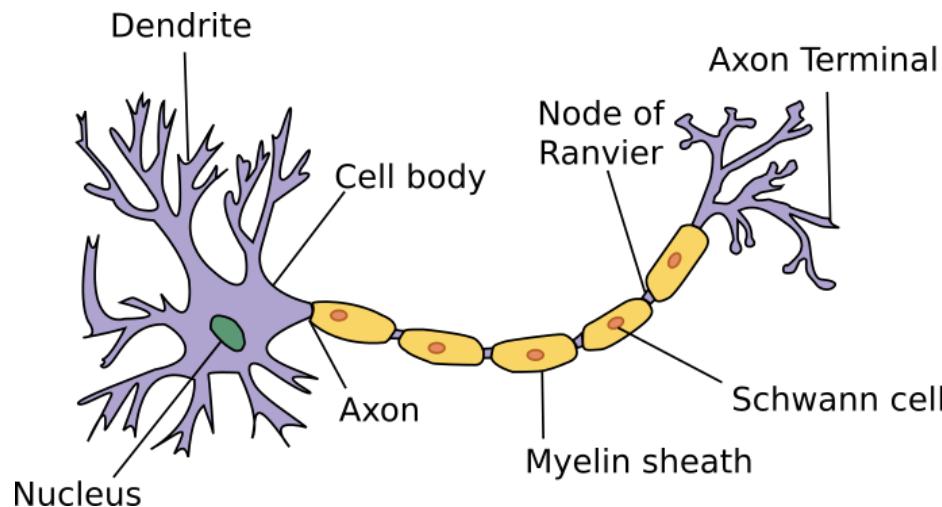
The Artificial Neuron



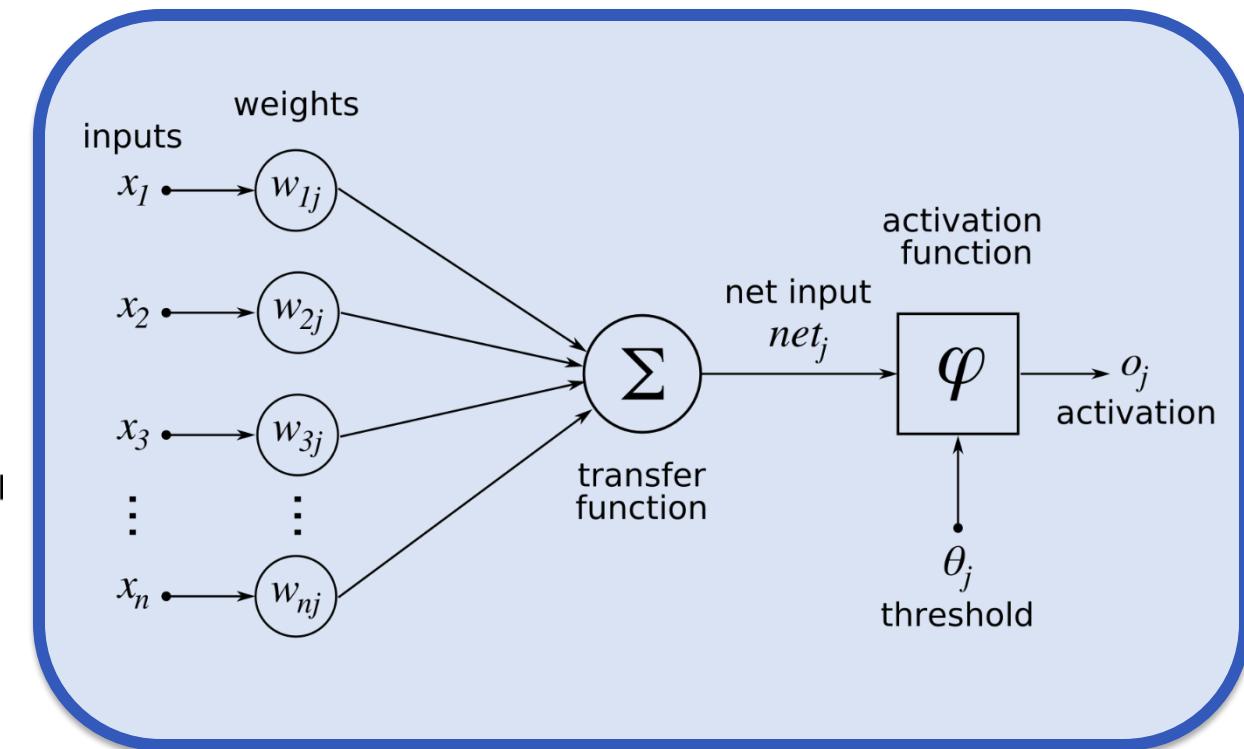
Warren
McCulloch (1898
- 1972)



Walter Pitts
(1924 - 1969)



Warren McCulloch and Walter Pitts, A Logical Calculus of Ideas Immanent in Nervous Activity, 1943, Bulletin of Mathematical Biophysics 5:115-133.



https://afit-r.github.io/ann_fundamentals

Multilayer Perceptron (MLP)

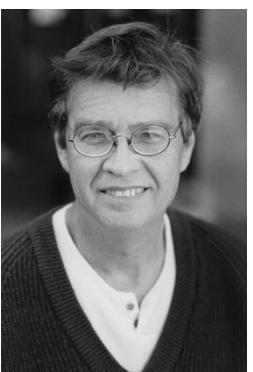
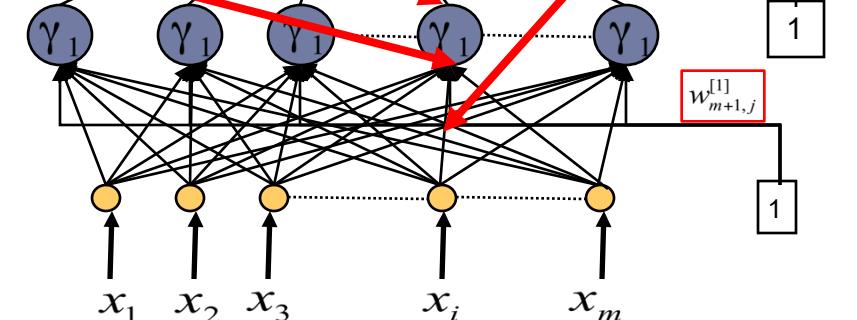
Output Layer

$$z = \sum_{j=1}^{\lambda} w_j^{[2]} a_j + w_{\lambda+1}^{[2]}$$

$$g_{\lambda}(x, w) = \gamma_2(z)$$

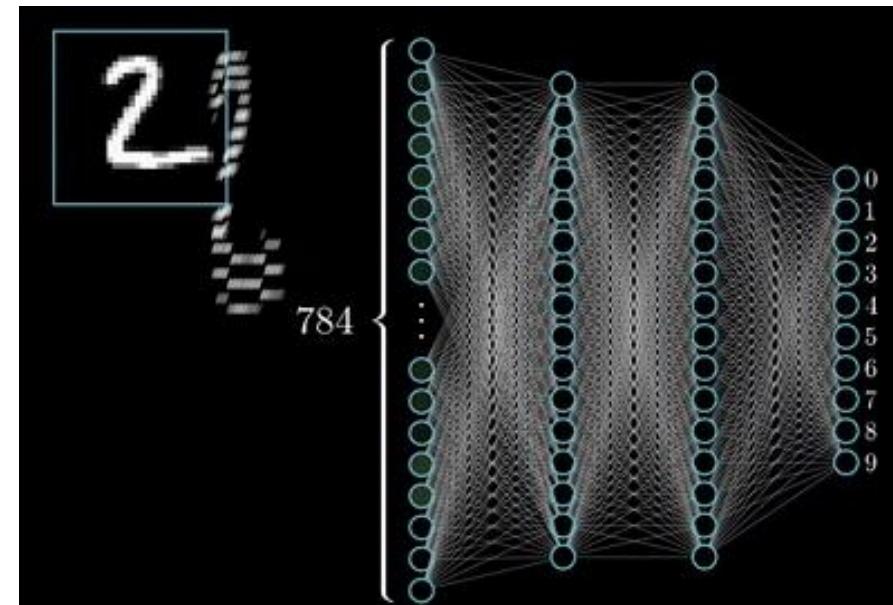
Hidden Layer

$$z_j = \sum_{i=1}^m w_{ij}^{[1]} x_i + w_{m+1,j}^{[1]}$$



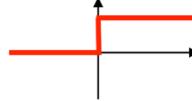
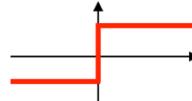
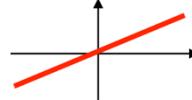
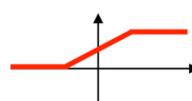
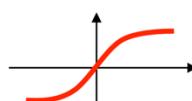
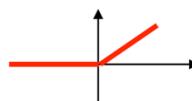
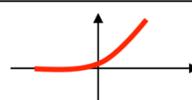
Rumelhart, David E., Geoffrey E. Hinton, and R. J. Williams. "Learning Internal Representations by Error Propagation", Parallel distributed processing: Explorations in the microstructure of cognition, Volume 1: Foundations. MIT Press, 1986

$$g_{\lambda}(x, w) = \gamma_2 \left(\sum_{j=1}^{\lambda} w_j^{[2]} \gamma_1 \left(\sum_{i=1}^m w_{ij}^{[1]} x_i + w_{m+1,j}^{[1]} \right) + w_{\lambda+1}^{[2]} \right)$$



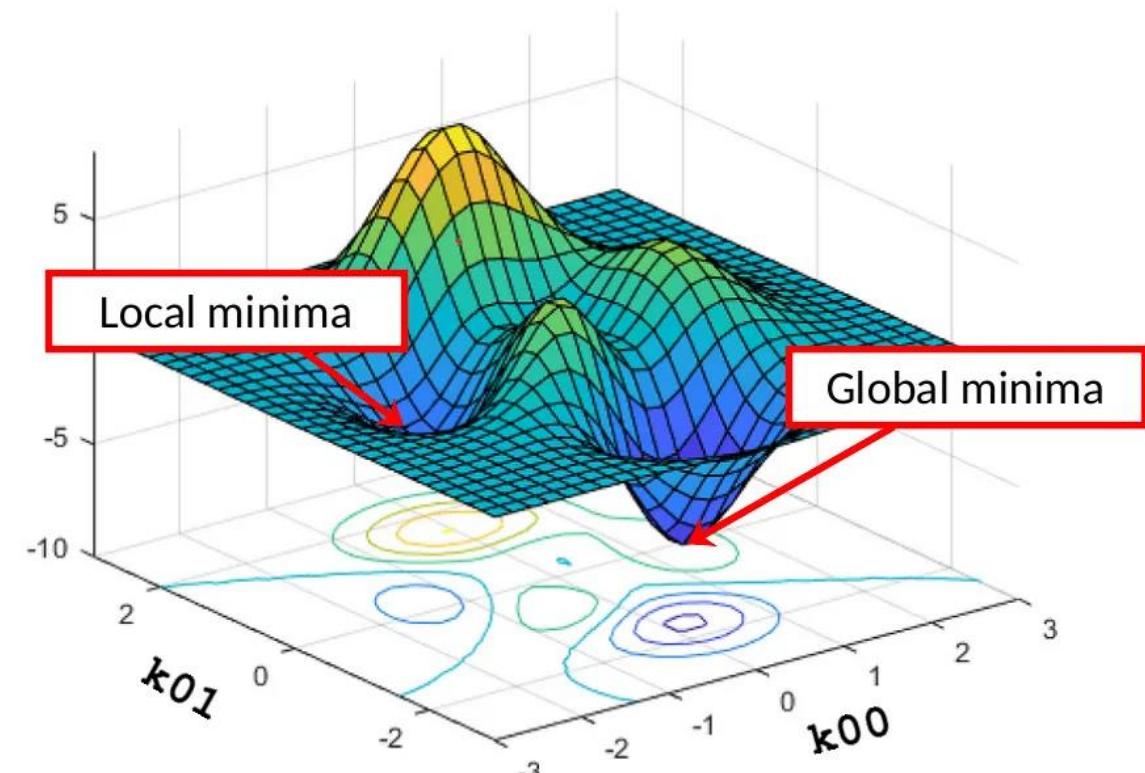
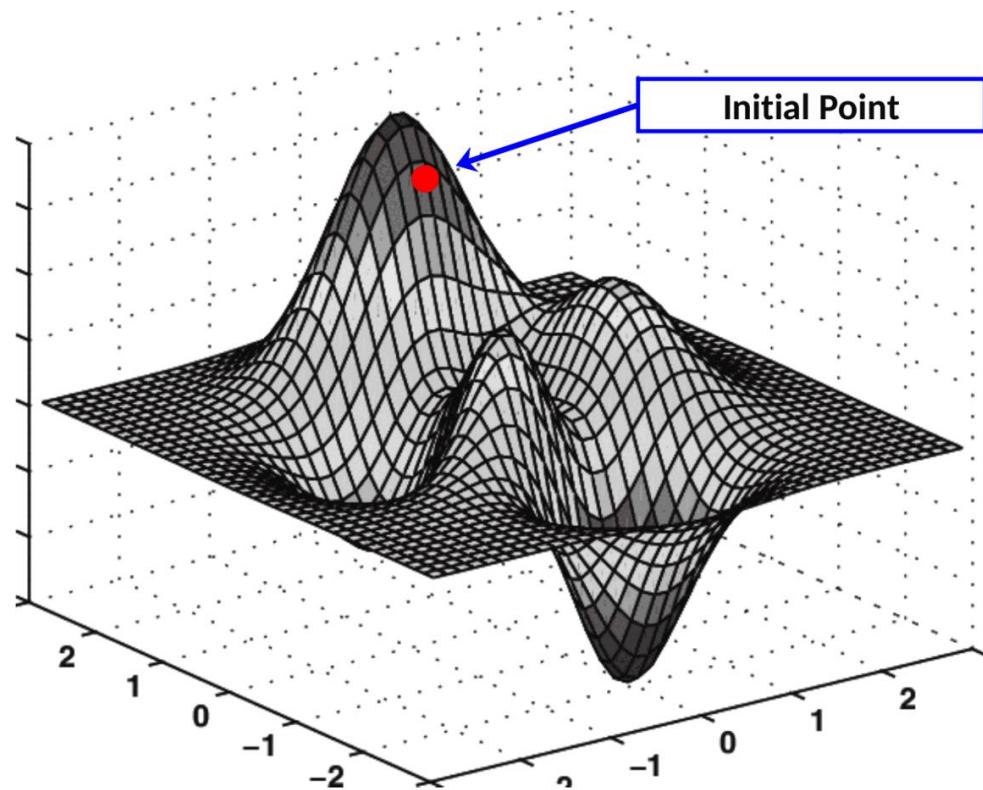
https://medium.com/@Suraj_Yadav/in-depth-knowledge-of-convolutional-neural-networks-b4bfff8145ab

Activation Functions

Activation function	Equation	Example	1D Graph
Unit step (Heaviside)	$\phi(z) = \begin{cases} 0, & z < 0, \\ 0.5, & z = 0, \\ 1, & z > 0, \end{cases}$	Perceptron variant	
Sign (Signum)	$\phi(z) = \begin{cases} -1, & z < 0, \\ 0, & z = 0, \\ 1, & z > 0, \end{cases}$	Perceptron variant	
Linear	$\phi(z) = z$	Adaline, linear regression	
Piece-wise linear	$\phi(z) = \begin{cases} 1, & z \geq \frac{1}{2}, \\ z + \frac{1}{2}, & -\frac{1}{2} < z < \frac{1}{2}, \\ 0, & z \leq -\frac{1}{2}, \end{cases}$	Support vector machine	
Logistic (sigmoid)	$\phi(z) = \frac{1}{1 + e^{-z}}$	Logistic regression, Multi-layer NN	
Hyperbolic tangent	$\phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$	Multi-layer Neural Networks	
Rectifier, ReLU (Rectified Linear Unit)	$\phi(z) = \max(0, z)$	Multi-layer Neural Networks	
Rectifier, softplus	$\phi(z) = \ln(1 + e^z)$	Multi-layer Neural Networks	

Delta Rule

$$w_j(t+1) = w_j(t) + \alpha(y - \hat{y})x_j$$



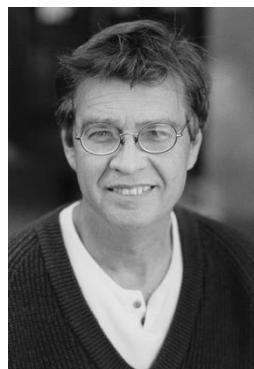
<https://pub.towardsai.net/deep-learning-from-scratch-in-modern-c-gradient-descent-670bc5889112>

Backpropagation

$$\mathbf{w}(t + 1) = \mathbf{w}(t) - \alpha(t) \nabla_{\mathbf{w}} l(h(\cdot, \mathbf{w}); z_t)$$

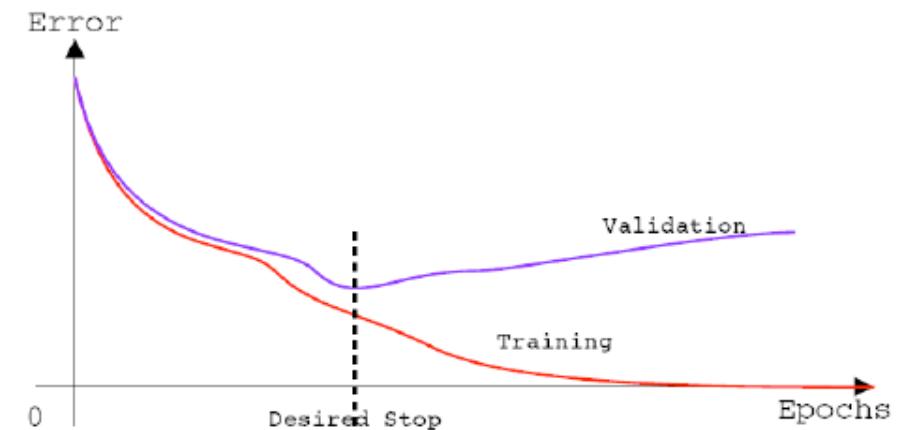
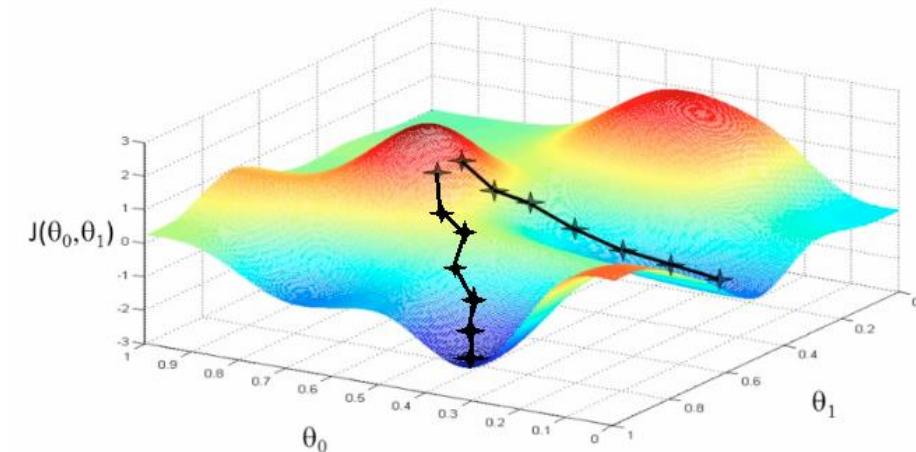


Paul Werbos.
1974 | In his Harvard PhD thesis, Paul Werbos describes training neural networks through backpropagation.

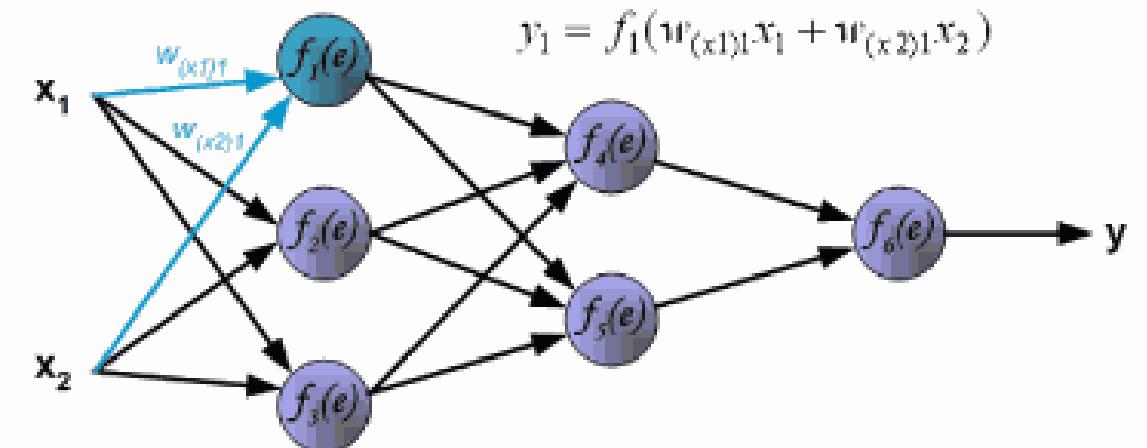
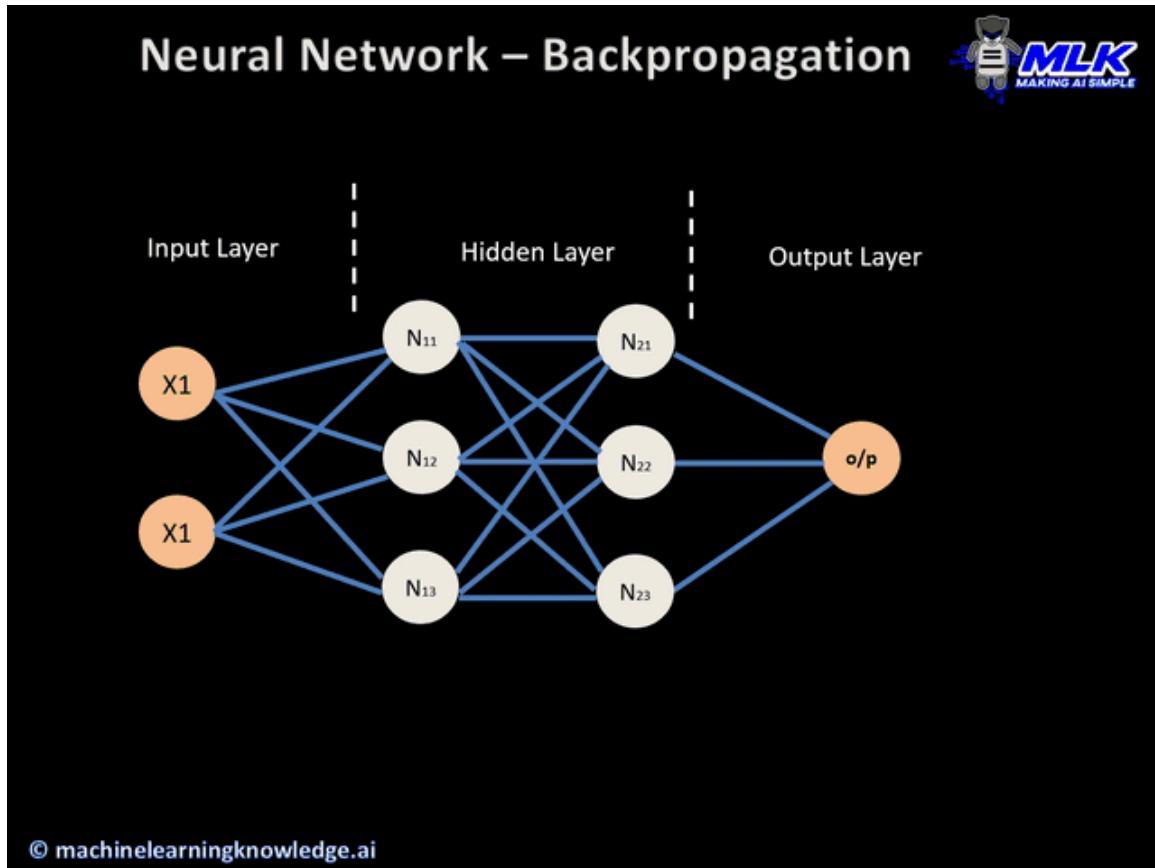


Rumelhart, David E., Geoffrey E. Hinton, and R. J. Williams. "Learning Internal Representations by Error Propagation", Parallel distributed processing: Explorations in the microstructure of cognition, Volume I: Foundations. MIT Press, 1986

http://www.holehouse.org/mlclass/01_02_Introduction_regression_analysis_and_gr.html



Backpropagation



<https://medium.com/deeper-deep-learning-tr/adim-adim-forward-and-back-propagation-cf4cd18276ee>

Backpropagation Algorithm

$$\delta(o, o) = \frac{\partial L}{\partial a_o} = \Phi'(a_o) \cdot \frac{\partial L}{\partial o}$$

$$\frac{\partial L}{\partial w_{(h_{r-1}, h_r)}} =$$

$$\frac{\partial L}{\partial o} \cdot \Phi'(a_o) \cdot \left[\sum_{[h_r, h_{r+1}, \dots, h_k, o] \in \mathcal{P}} \frac{\partial a_o}{\partial a_{h_k}} \prod_{i=r}^{k-1} \frac{\partial a_{h_{i+1}}}{\partial a_{h_i}} \right]$$

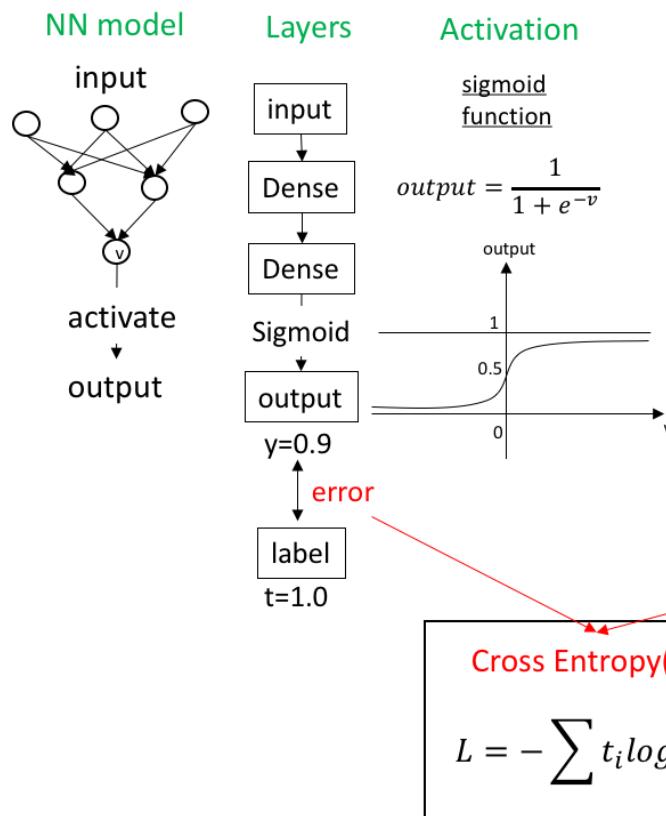
$$\underbrace{\frac{\partial a_{h_r}}{\partial w_{(h_{r-1}, h_r)}}}_{h_{r-1}}$$

Backpropagation computes $\delta(h_r, o) = \frac{\partial L}{\partial a_{h_r}}$

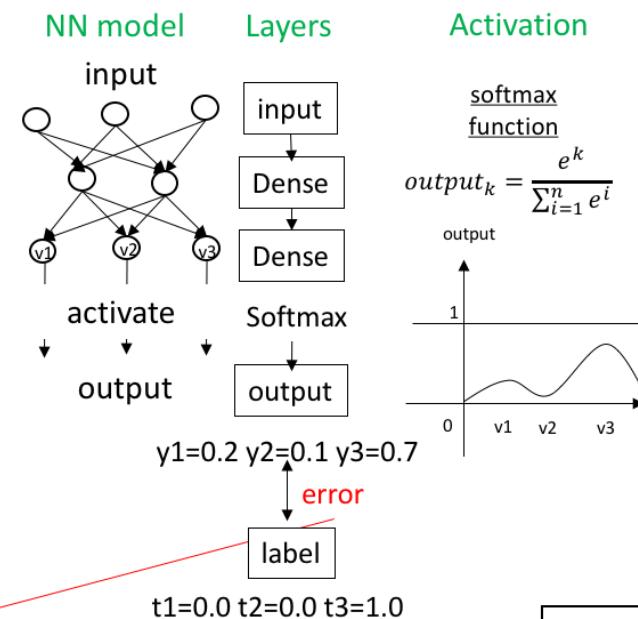
$$\delta(h_r, o) = \frac{\partial L}{\partial a_{h_r}} = \sum_{h: h_r \Rightarrow h} \widehat{\frac{\partial L}{\partial a_h}} \underbrace{\frac{\partial a_h}{\partial a_{h_r}}}_{\Phi'(a_{h_r})w_{(h_r, h)}} = \Phi'(a_{h_r}) \sum_{h: h_r \Rightarrow h} w_{(h_r, h)} \cdot \delta(h, o)$$

Loss Functions

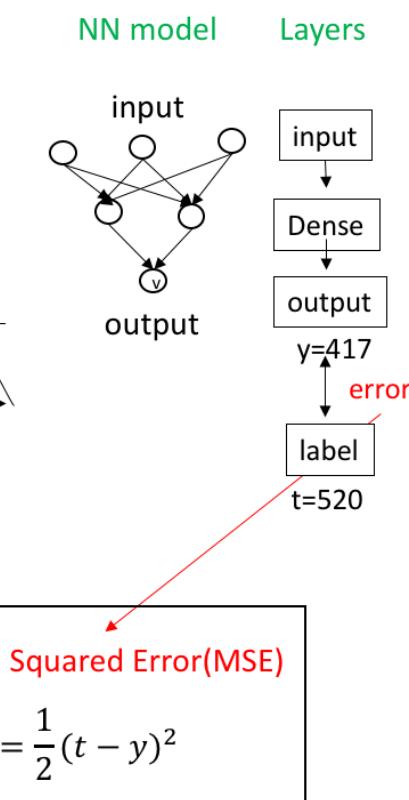
1. Binary Classification



2. Multiclass Classification



3. Regression



https://www.renom.jp/notebooks/tutorial/basic_algorithm/lossfunction/notebook.html

Universal Approximators

Theorem of Universal Approximation:

“Every bounded continuous function with bounded support can be approximated arbitrarily closely by a multi-layer perceptron by selecting enough but a finite number of hidden neurons with appropriate transfer function”



Cybenko., G. (1989) "Approximations by superpositions of sigmoidal functions",
Mathematics of Control, Signals, and Systems, 2 (4), 303-314

Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2:359--366, 1989

ANN Playground



<http://playground.tensorflow.org/>

Workshop 1 – MLP for Regression

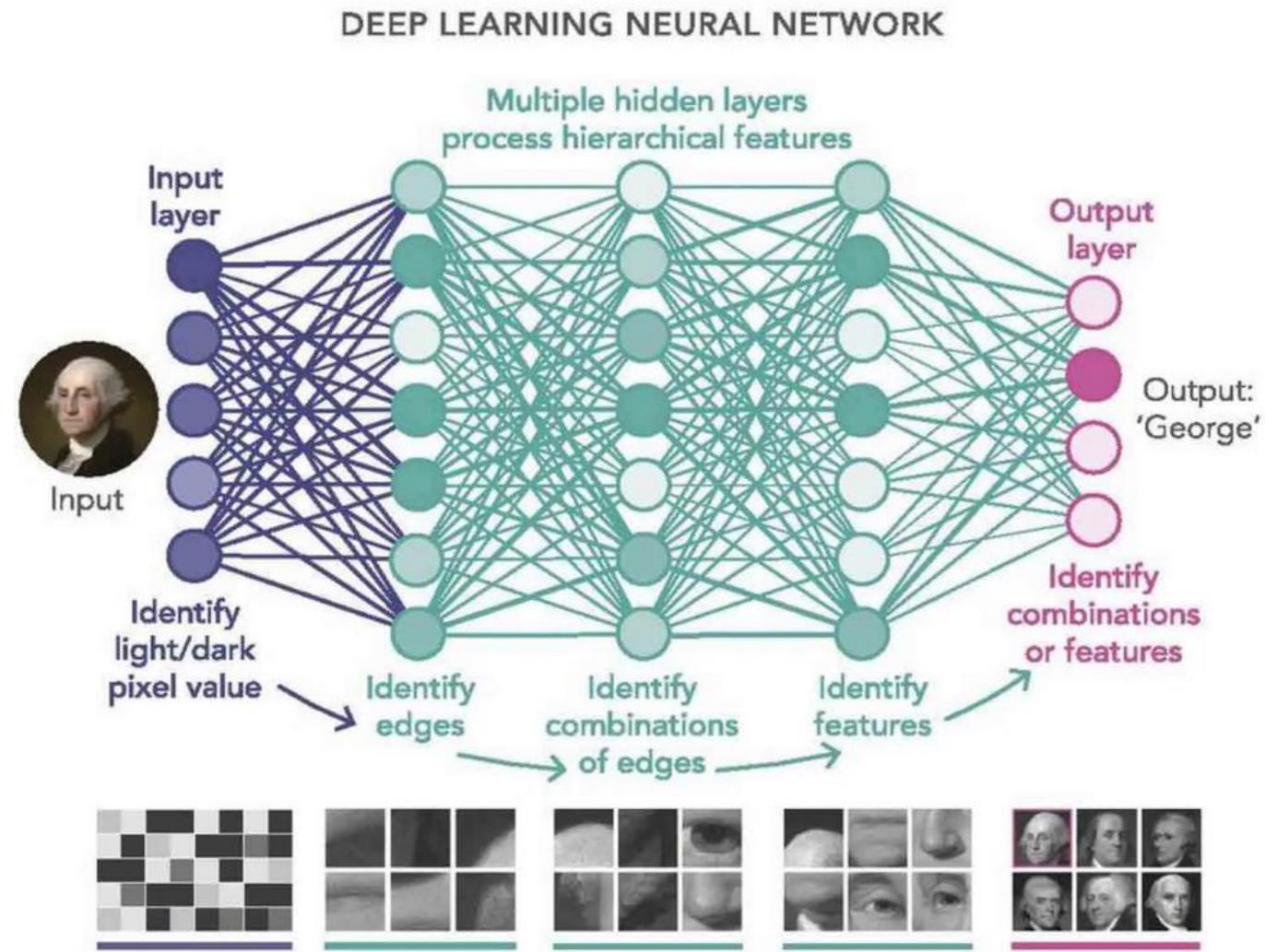


<https://colab.research.google.com/drive/1yW5Ia5hxcaDRNtPGAG3dDb4UQ2PeXvz2?usp=sharing>

Deep Learning and Convolutional Neural Networks **CNN**

¿What is Deep Learning?

- **Deep Learning (DL)** is a method based on deep neural networks that automatically learn complex structures and patterns from data through hierarchical and nonlinear representations.
- It is inspired by the architecture of the human brain.
- The model is composed of multiple intermediate layers.
- It enables the modeling of nonlinear relationships in large volumes of data.



M. Mitchell Warldrop, Artificial Neural Network; from <https://www.pnas.org>

Deep-Learning



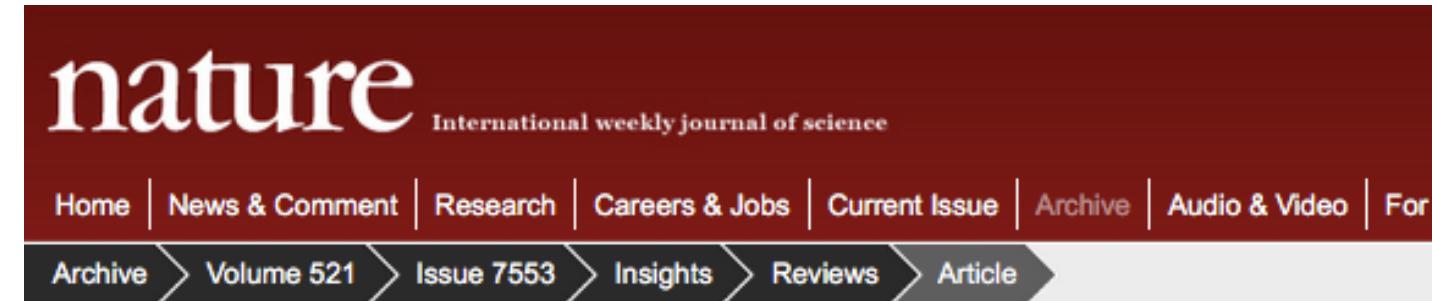
Yann LeCun



Yoshua Bengio



Geoffrey Hinton

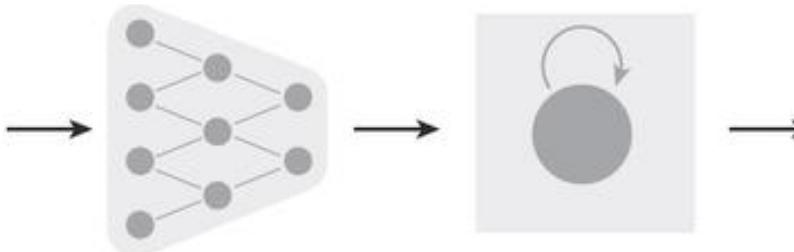


The screenshot shows the header of a Nature journal article. The title 'nature' is in large white letters, followed by 'International weekly journal of science'. Below the title is a navigation bar with links: Home, News & Comment, Research, Careers & Jobs, Current Issue, Archive, Audio & Video, and For Authors. A breadcrumb trail below the navigation bar shows the article's path: Archive > Volume 521 > Issue 7553 > Insights > Reviews > Article. The main heading 'ARTICLE PREVIEW' is centered above a link 'view full access options >'. The article title 'Deep learning' is in large bold black font. Below it are the authors' names: 'Yann LeCun, Yoshua Bengio & Geoffrey Hinton'. There are links for 'Affiliations' and 'Corresponding author'. The publication details 'Nature 521, 436–444 (28 May 2015) doi:10.1038/nature14539' are shown, along with the dates 'Received 25 February 2015 | Accepted 01 May 2015 | Published online 27 May 2015'. At the bottom are four buttons: 'Citation', 'Reprints', 'Rights & permissions', and 'Article metrics'.

High-level Translation from Images



Vision
Deep-CNN



Language
Generating RNN

A group of people
shopping at an outdoor
market.

There are many
vegetables at the
fruit stand.



A woman is throwing a **frisbee** in a park.



A **dog** is standing on a hardwood floor.



A **stop** sign is on a road with a
mountain in the background



A little **girl** sitting on a bed with a **teddy bear**.



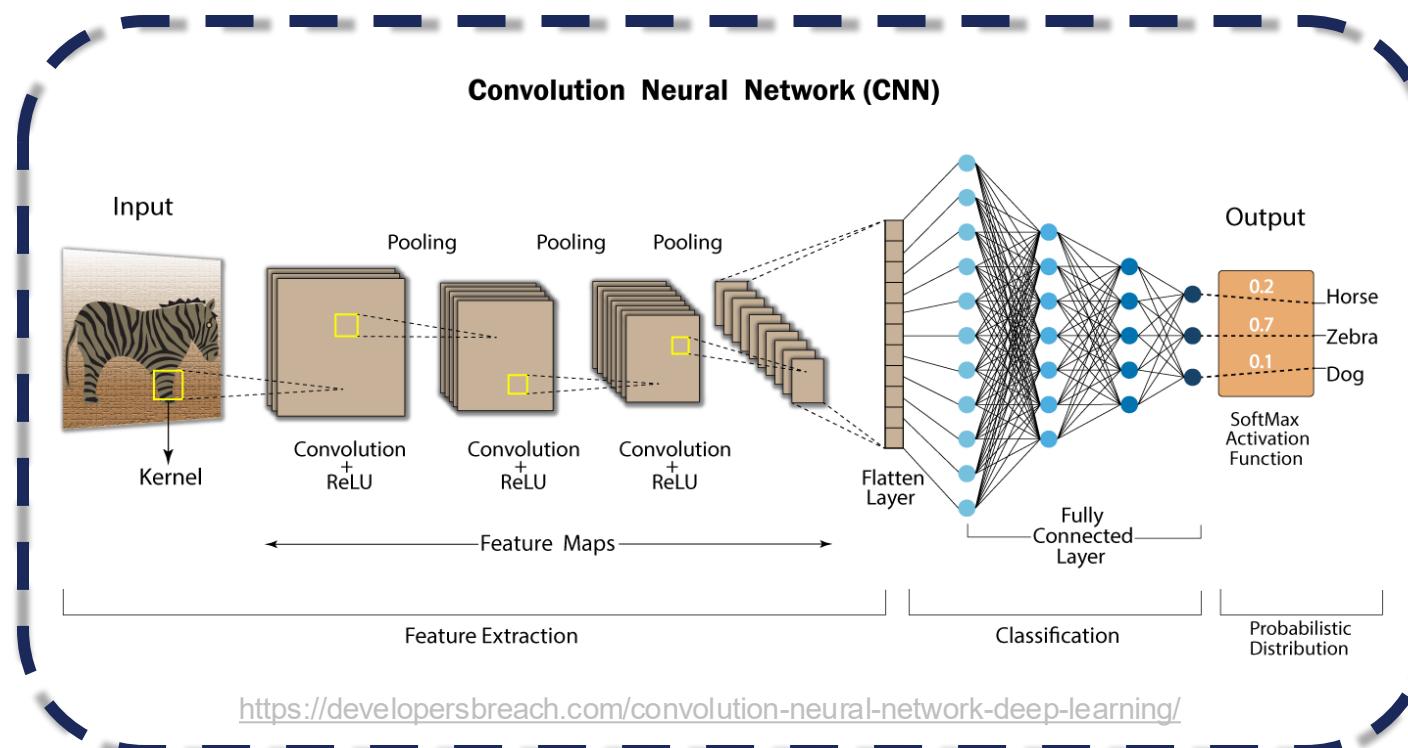
A group of **people** sitting on a boat in the water.



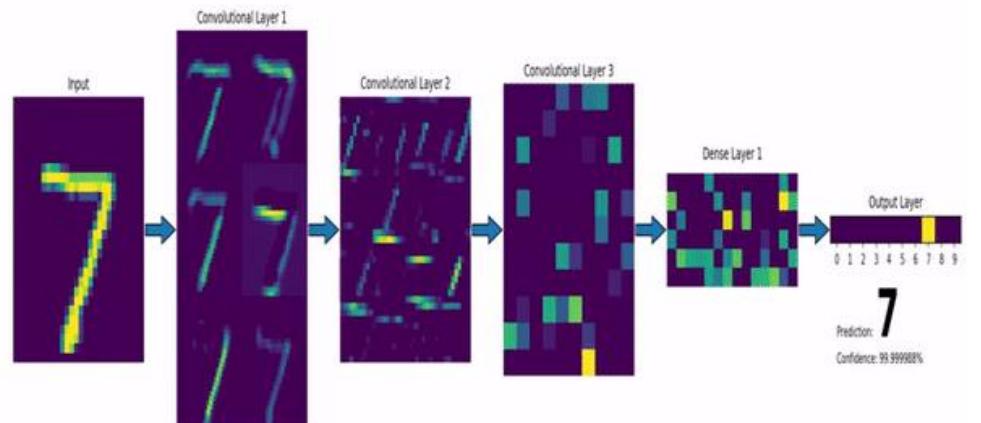
A **giraffe** standing in a forest with
trees in the background.

LeCunn, Bengio & Hinton
(2015). Deep learning.
Nature 521, pp. 436–444.
doi:10.1038/nature14539

Convolutional Neural Network



<https://developersbreach.com/convolution-neural-network-deep-learning/>



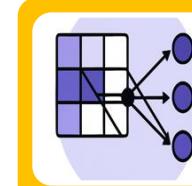
<https://www.louisbouchard.ai/densenet-explained/>

Convolutional Neural Network

A **Convolutional Neural Network (CNN)** is a type of neural network specialized in processing data with spatial structure. Its design enables the detection of visual patterns directly from the pixels, without the need to manually define features.



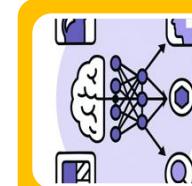
Each neuron connects only to a local region of the input (receptive field), capturing patterns such as edges or textures.



Layers are stacked hierarchically, learning from simple details to complex representations.

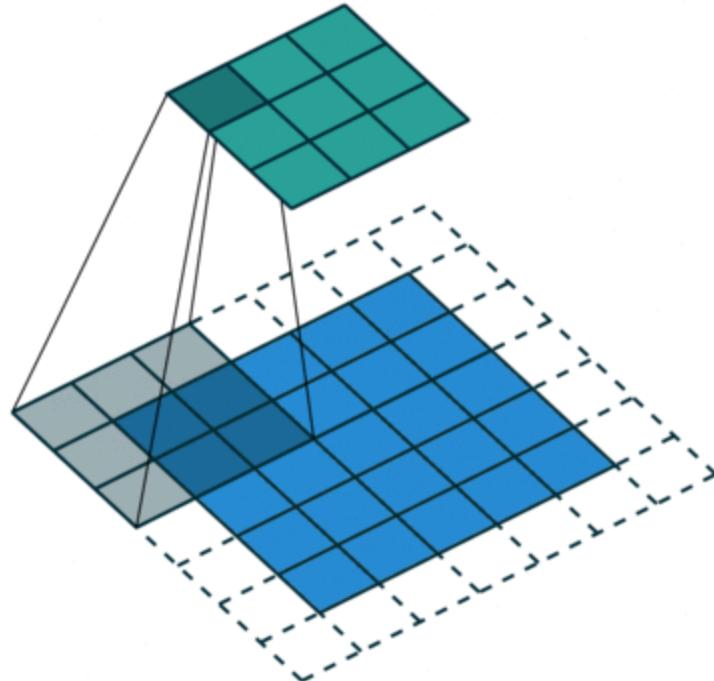


It uses **shared weights** and **filters** that slide across the image, drastically reducing the number of parameters.



It is highly effective in tasks such as **classification**, **segmentation**, and **object detection**.

Convolutional Operations



Kernel: A small filter that scans the image and extracts local features such as edges, textures, or patterns.

1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

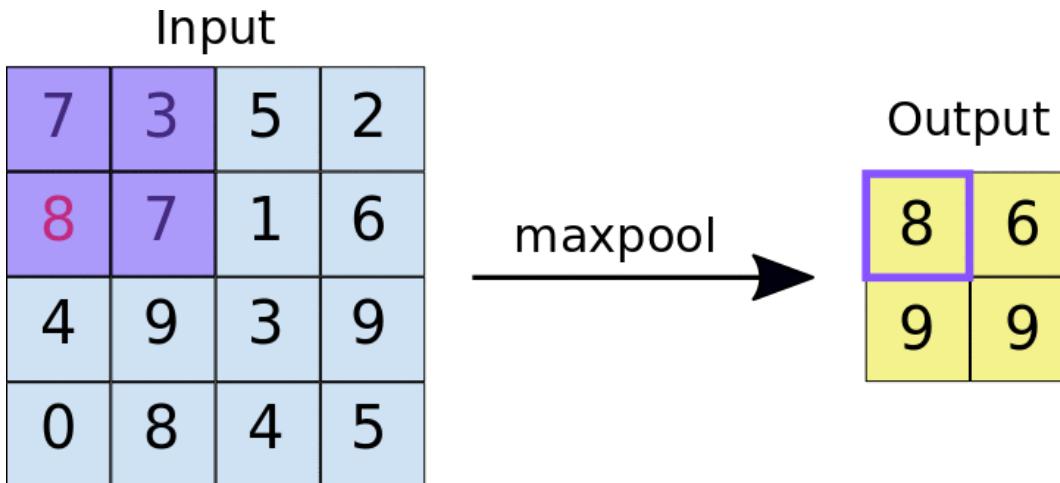
Image

4		

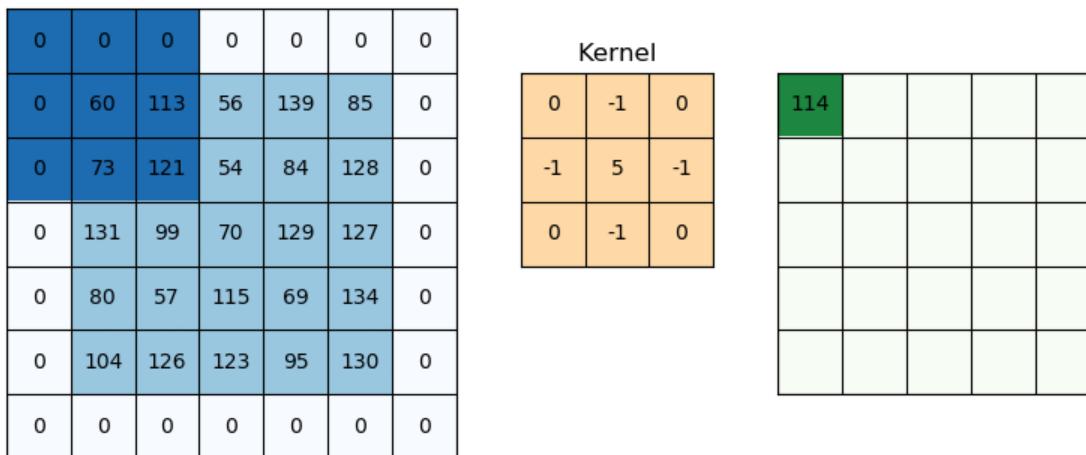
Convolved
Feature

Convolution: An operation that multiplies and sums the kernel values with local regions of the image to generate a feature map.

Convolutional Operations



Max-Pooling: A dimensionality reduction technique that keeps only the maximum value from local regions, preserving the most important features.

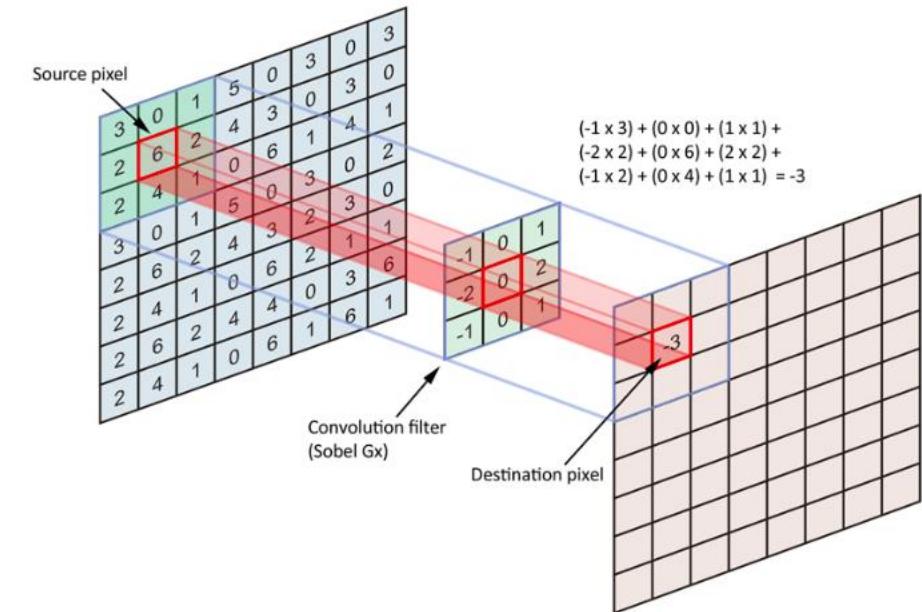


Padding: Adding borders (usually zeros) around the image to preserve the output size or prevent loss of information at the edges.

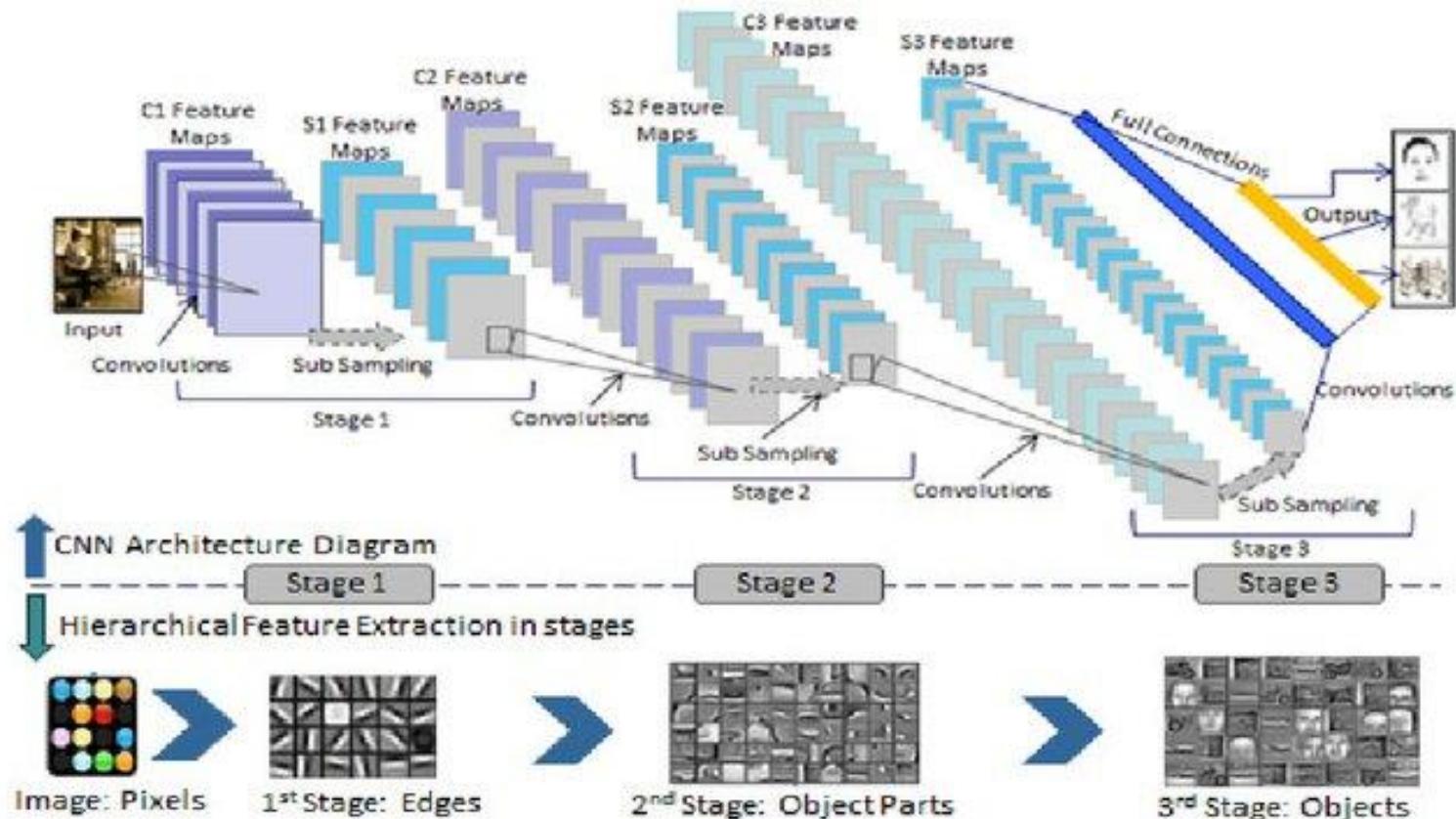
Automatic Feature Extraction with Convolutional Layers

Input Volume (+pad 1) (7x7x3)	Filter W0 (3x3x3)	Output Volume (3x3x2)
<code>x[:, :, 0]</code>	<code>w0[:, :, 0]</code>	<code>o[:, :, 0]</code>
<code>0 0 0 0 0 0 0</code>	<code>-1 0 1</code>	<code>2 3 3</code>
<code>0 0 0 1 0 2 0</code>	<code>0 0 1</code>	<code>3 7 3</code>
<code>0 1 0 2 0 1 0</code>	<code>1 -1 1</code>	<code>8 10 -3</code>
<code>0 1 0 2 2 0 0</code>	<code>w0[:, :, 1]</code>	<code>o[:, :, 1]</code>
<code>0 2 0 0 2 0 0</code>	<code>-1 0 1</code>	<code>-8 -8 -3</code>
<code>0 2 1 2 2 0 0</code>	<code>1 -1 1</code>	<code>1 -1 0</code>
<code>0 0 0 0 0 0 0</code>	<code>0 1 0</code>	<code>-3 1 0</code>
<code>x[:, :, 1]</code>	<code>w0[:, :, 2]</code>	<code>-3 -8 -5</code>
<code>0 0 0 0 0 0 0</code>	<code>-1 1 1</code>	
<code>0 2 1 2 1 1 0</code>	<code>1 1 0</code>	
<code>0 2 1 2 0 1 0</code>	<code>0 -1 0</code>	
<code>0 0 2 1 0 1 0</code>	<code>1 0 0</code>	
<code>0 1 2 2 2 2 0</code>		
<code>0 0 1 2 0 1 0</code>		
<code>0 0 0 0 0 0 0</code>		
<code>x[:, :, 2]</code>		
<code>0 0 0 0 0 0 0</code>		
<code>0 2 1 1 2 0 0</code>		
<code>0 1 0 0 1 0 0</code>		
<code>0 0 1 0 0 0 0</code>		
<code>0 1 0 2 1 0 0</code>		
<code>0 2 2 1 1 1 0</code>		
<code>0 0 0 0 0 0 0</code>		
	<code>Bias b0 (1x1x1)</code>	
	<code>b0[:, :, 0]</code>	
	<code>0</code>	

toggle movement

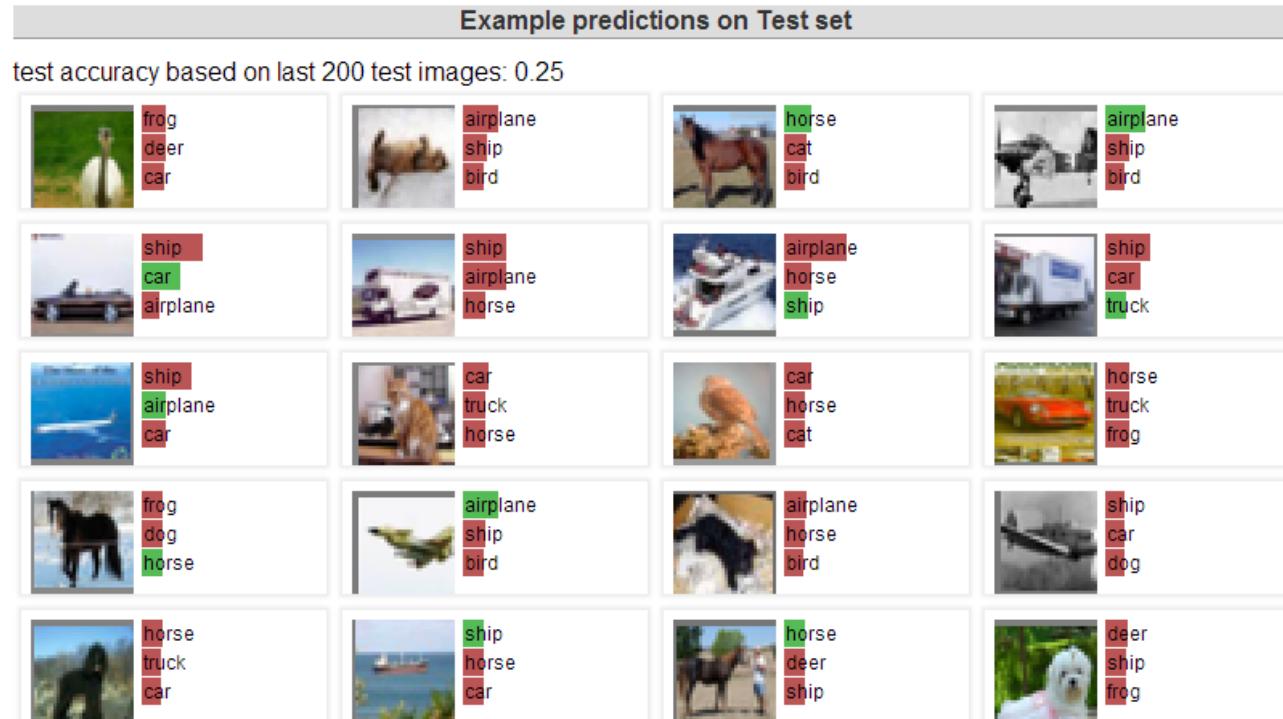
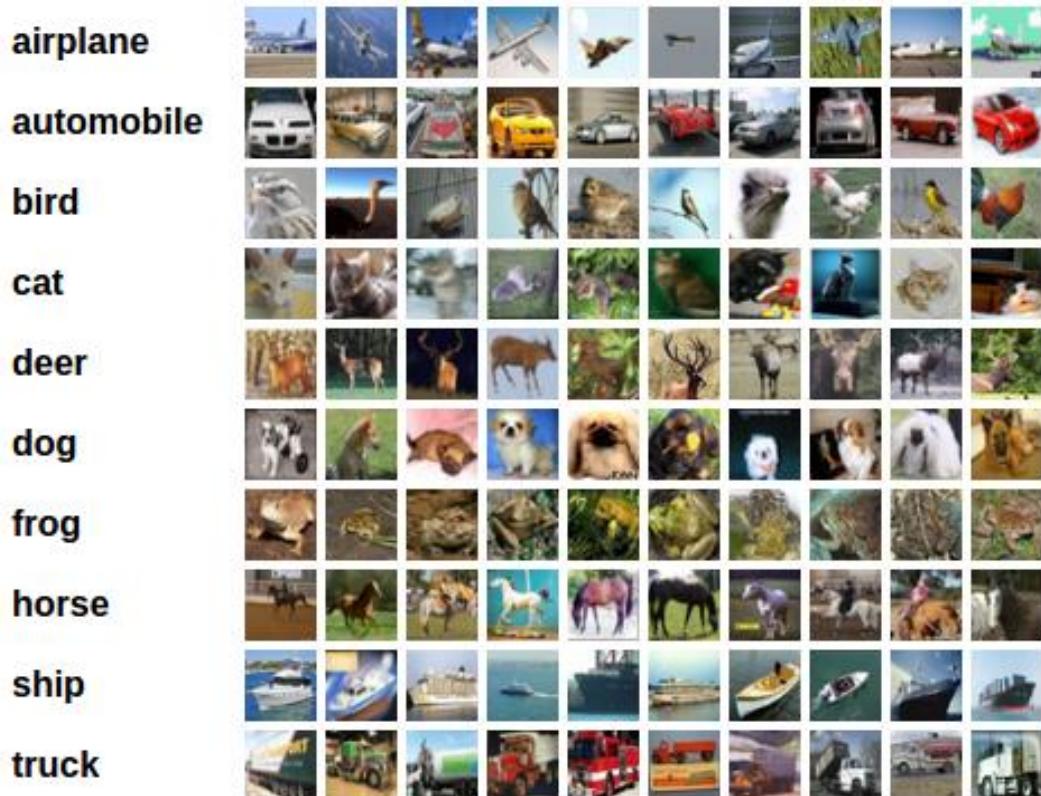


Learning Hierarchy of Visual Features in CNN Architecture



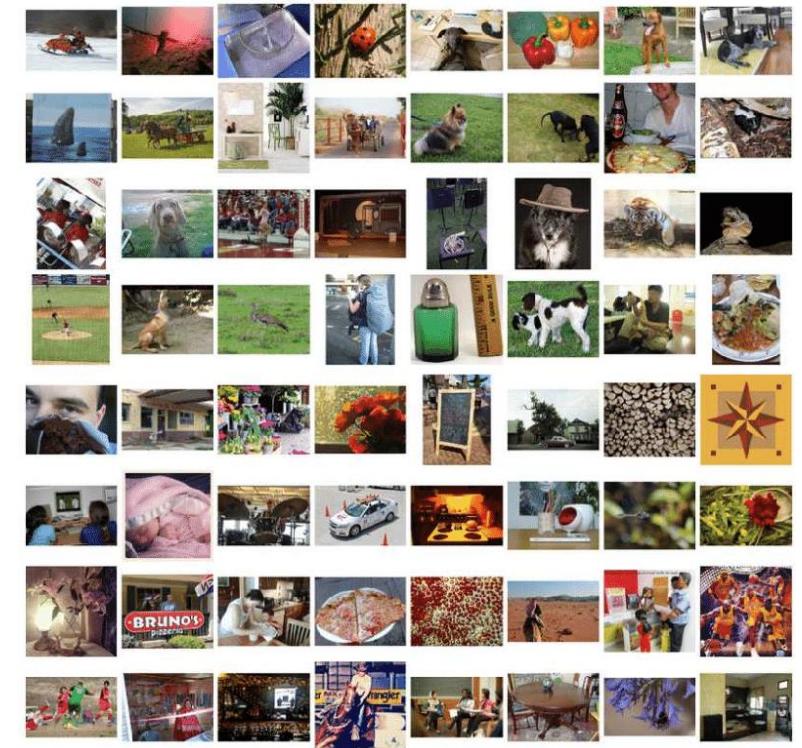
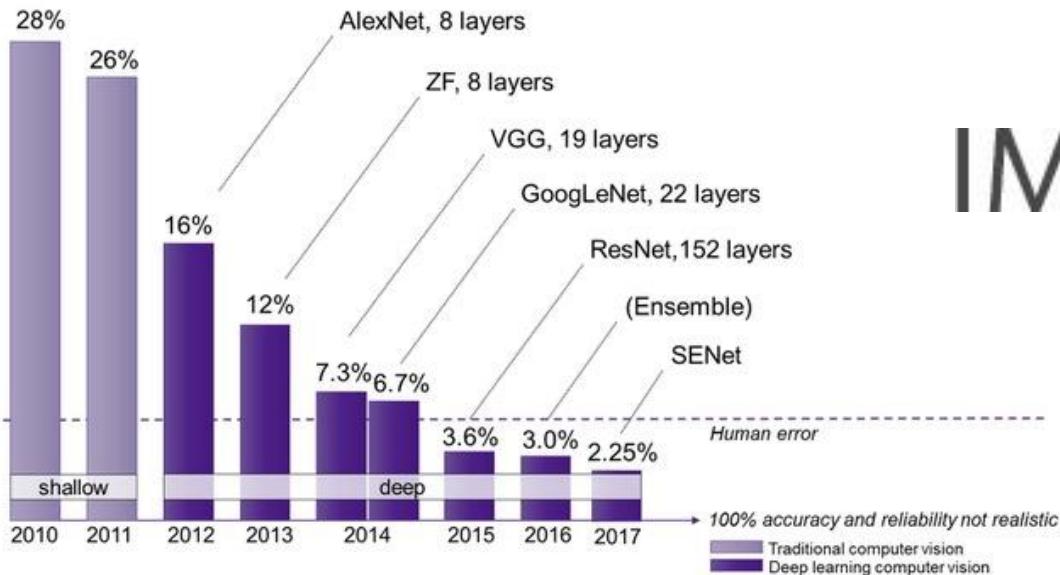
Katole et al. Hierarchical Deep Learning Architecture For 10K Objects Classification. Doi: 10.5121/csit.2015.51408

ConvNetJS CIFAR-10 demo



<http://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

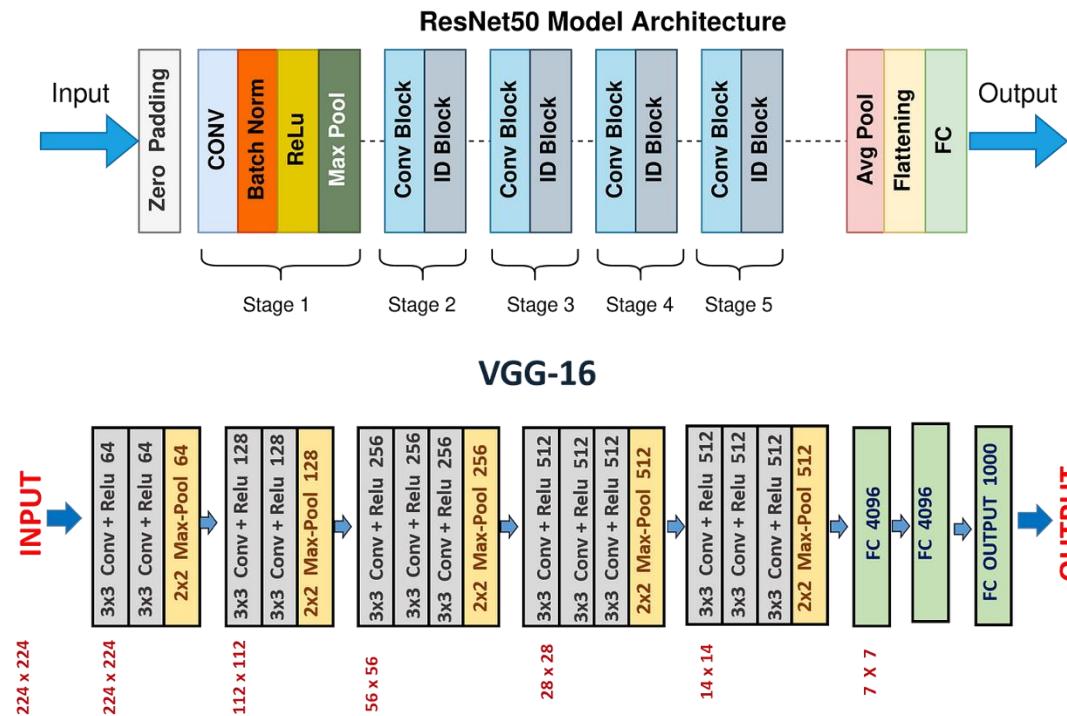
Large Scale Visual Recognition Challenge (ILSVRC)



<http://www.image-net.org/challenges/LSVRC/>

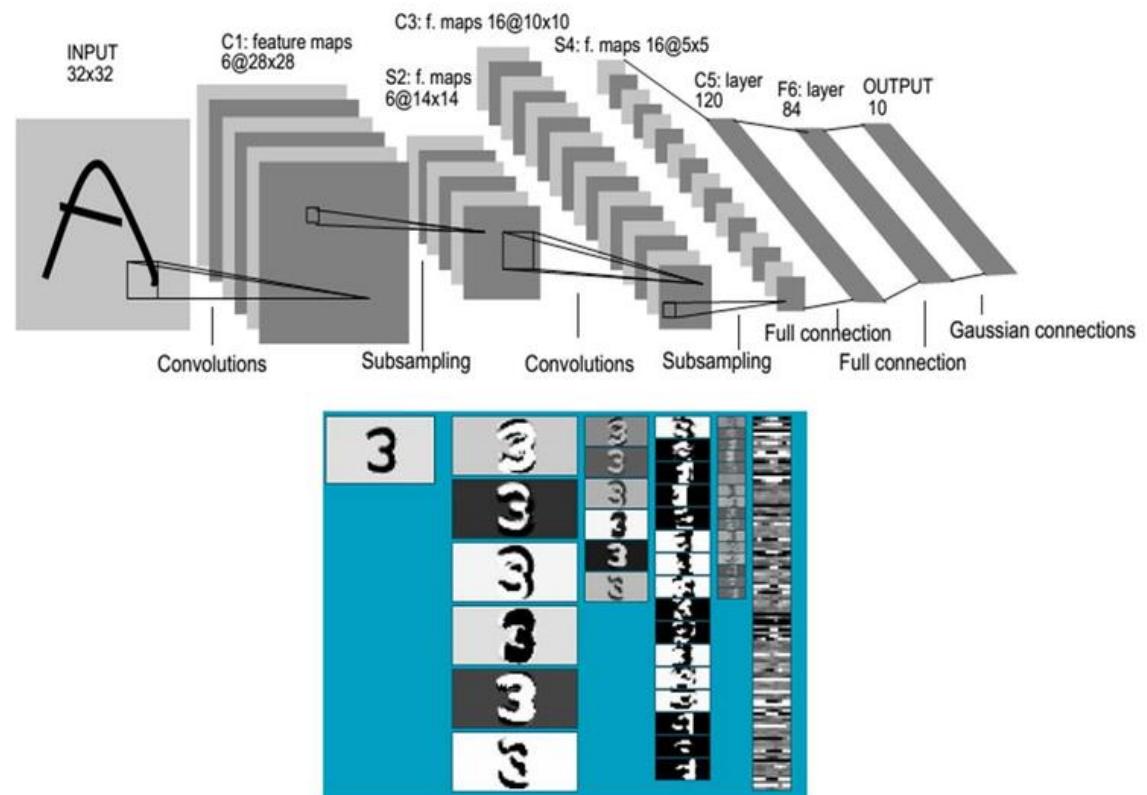
<https://gigaom.com/2014/08/22/with-enlitic-a-veteran-data-scientist-plans-to-fight-disease-using-deep-learning/>

Architectures of CNN

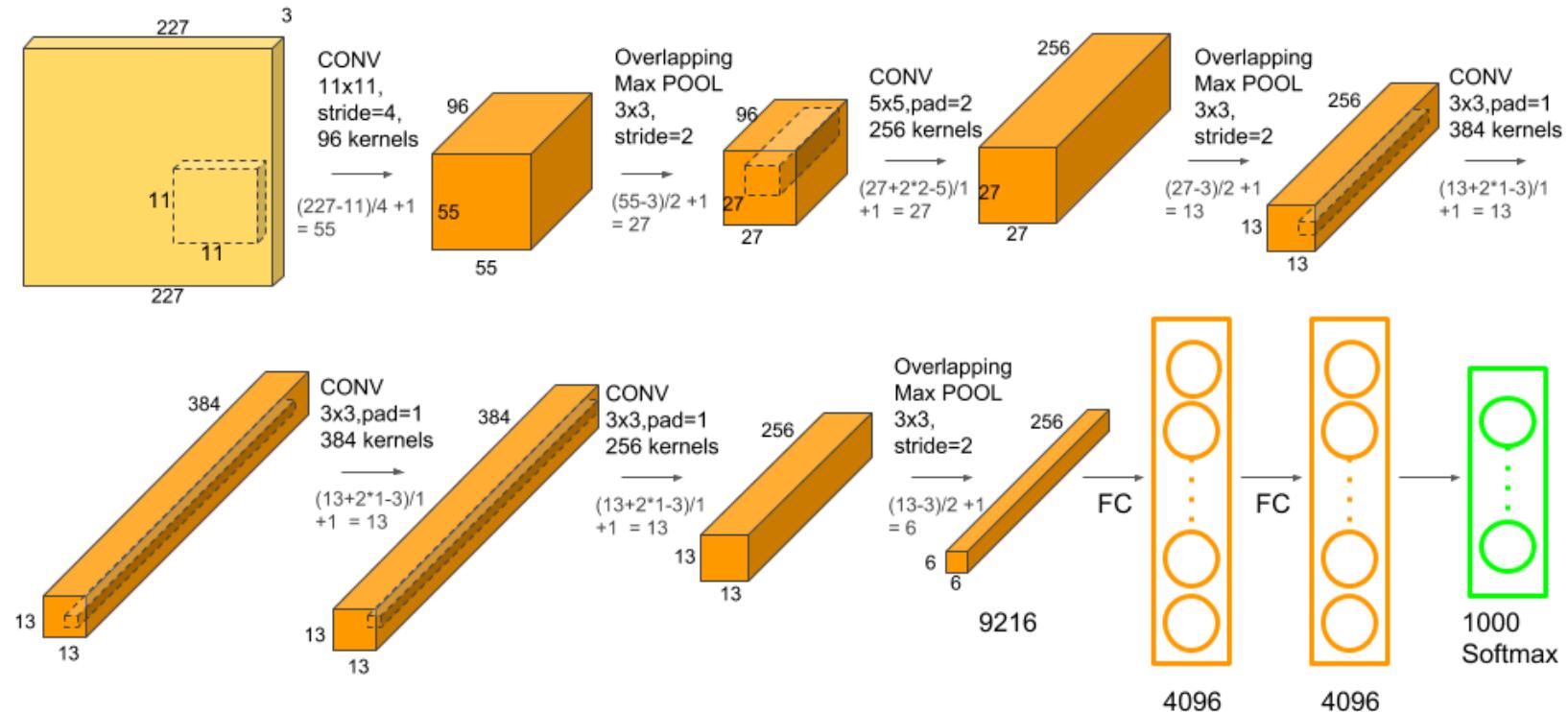


Estas arquitecturas CNN (VGG, ResNet, U-Net) aprenden representaciones visuales jerárquicas y se adaptan a tareas como clasificación y segmentación pixel a pixel.

LeNet

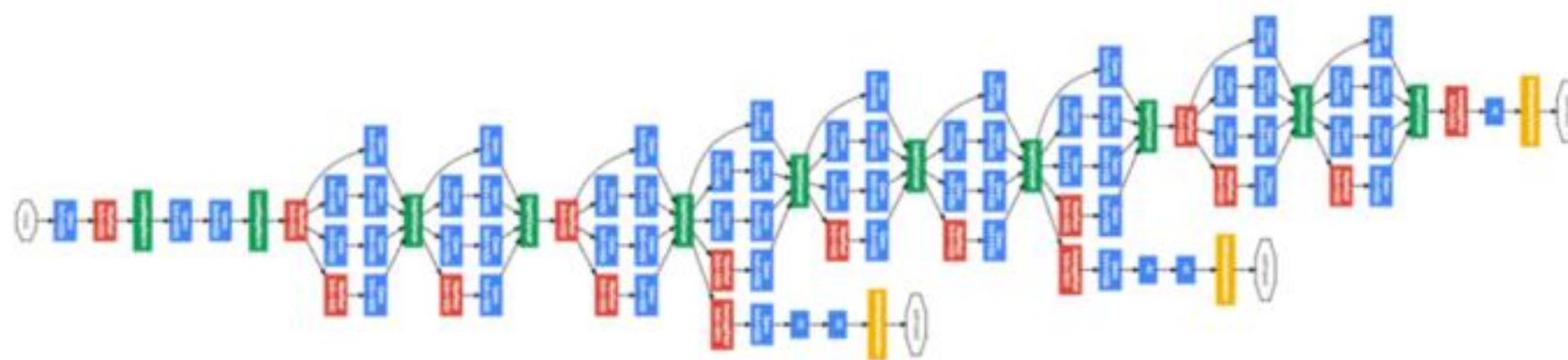


AlexNet

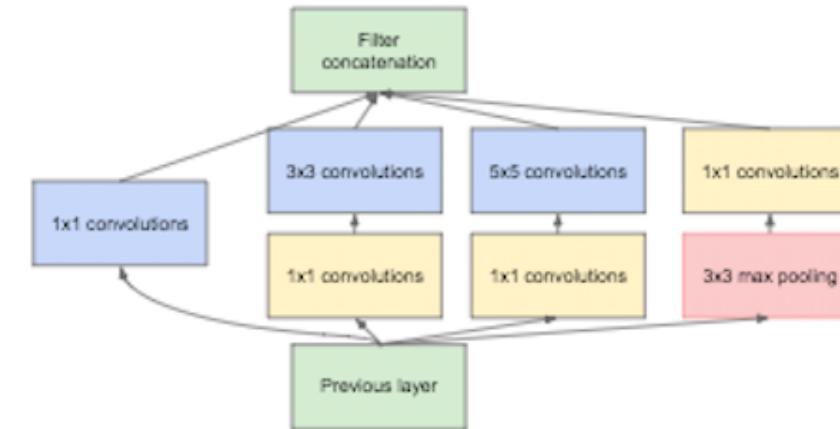


<https://medium.com/@smallfishbigsea/a-walk-through-of-alexnet-6cbd137a5637>

GoogLeNet



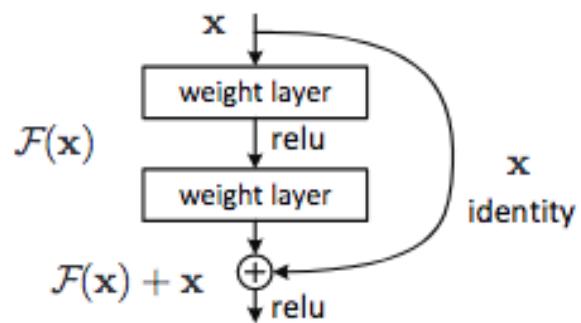
- **GoogLeNet:** developed by C. Szegedy et al. from Google. (2014).
- Introduced the **Inception module**, which reduced the number of parameters in the network (**sparsity**).



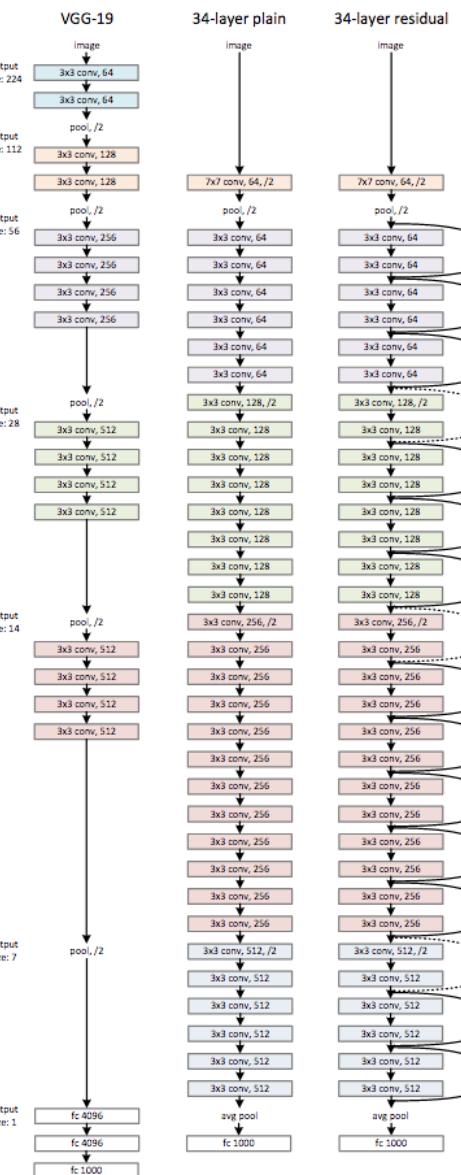
Módulo de Inception con reducción de la dimensionalidad

Residual Learning

- **Residual learning (ResNet)** was proposed by *K. He et al.* (2015).
 - It is used to enable learning at greater depth.

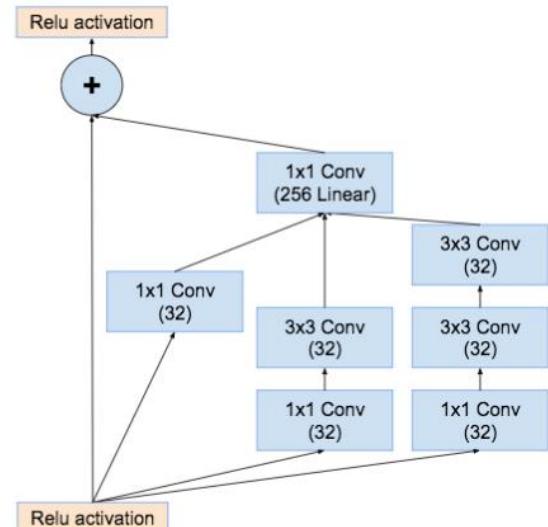


Bloque de aprendizaje residual

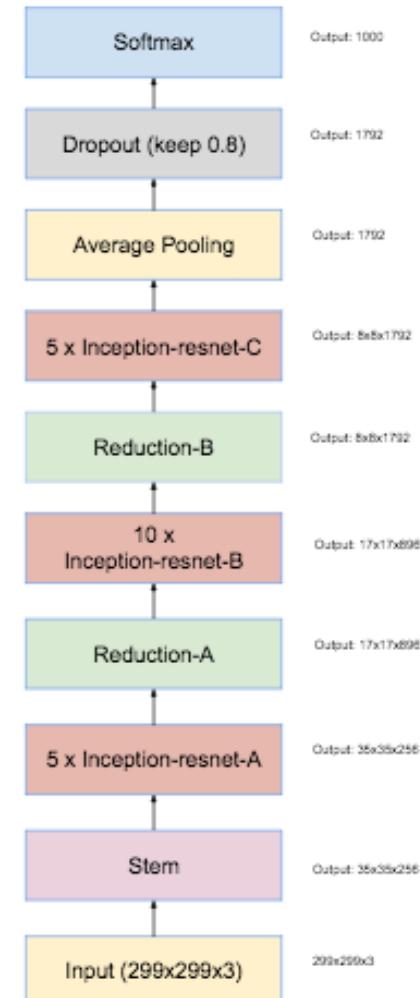


Inception-Resnet

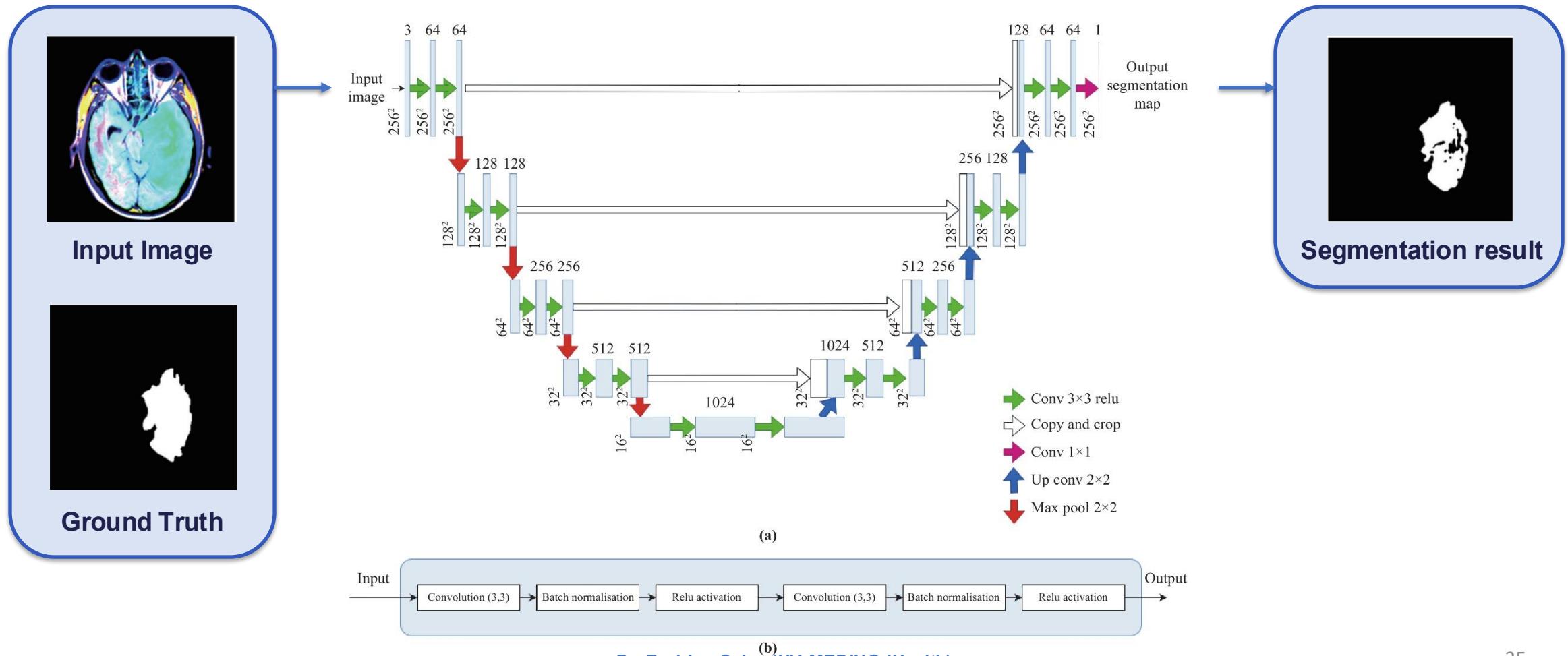
- Proposed by C. Szegedy 2016
- Combines the Inception module (Sparsity) with the ResNet module (Deep)



Módulo de Inception-ResNet-A



U-Net: Convolutional Networks for Biomedical Image Segmentation

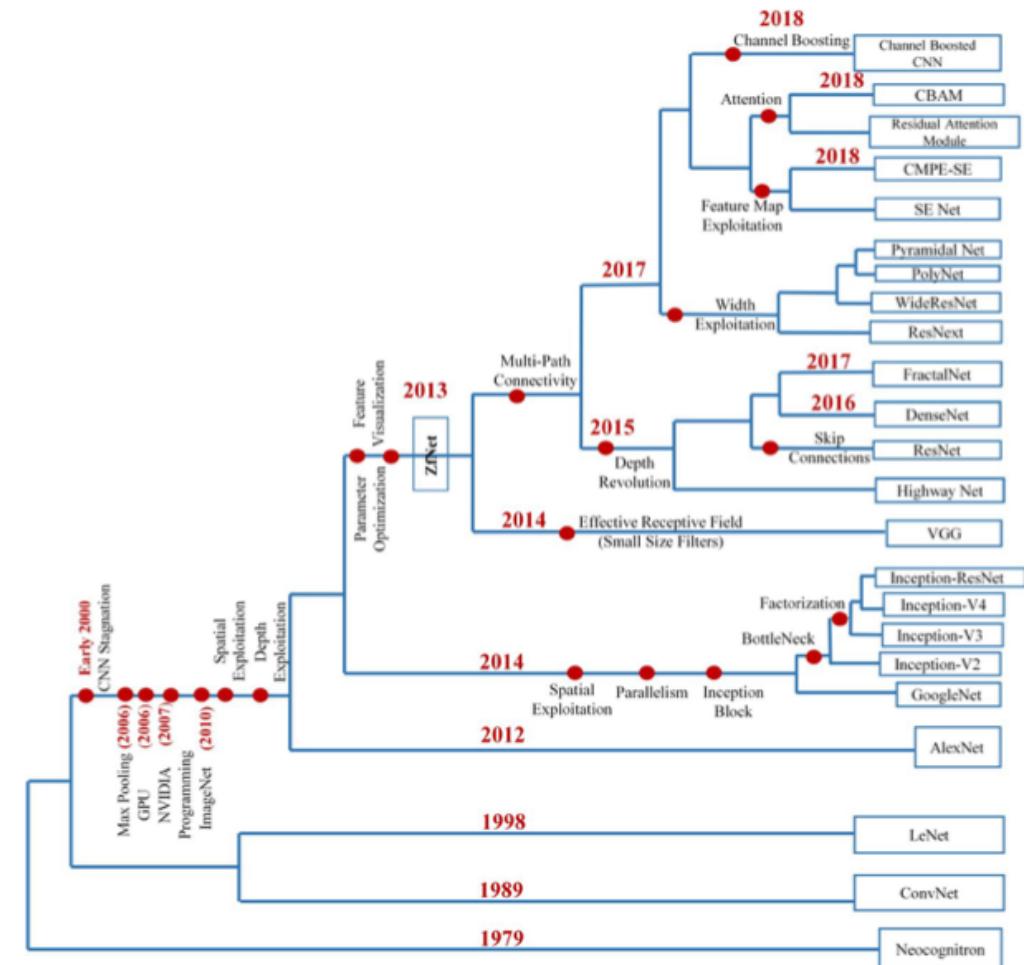
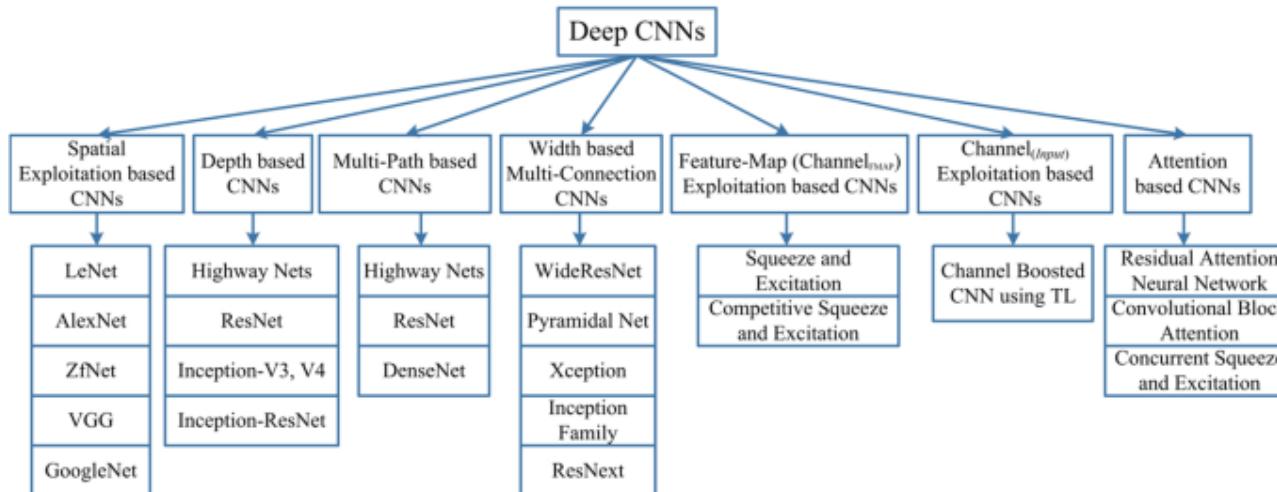


Architectures of Deep Convolutional Neural Networks

Artificial Intelligence Review (2020) 53:5455–5516
<https://doi.org/10.1007/s10462-020-09825-6>

A survey of the recent architectures of deep convolutional neural networks

Asifullah Khan^{1,2} · Anabia Sohail^{1,2} · Umme Zahoor¹ · Aqsa Saeed Qureshi¹



Benchmark Analysis of Deep Learning

IEEE Access

Multidisciplinary | Rapid Review | Open Access Journal

Received October 1, 2018, accepted October 17, 2018, date of publication October 24, 2018, date of current version November 19, 2018.

Digital Object Identifier 10.1109/ACCESS.2018.2877890

Benchmark Analysis of Representative Deep Neural Network Architectures

SIMONE BIANCO^{ID¹}, REMI CADENE², LUIGI CELONA^{ID¹}, AND PAOLO NAPOLITANO^{ID¹}

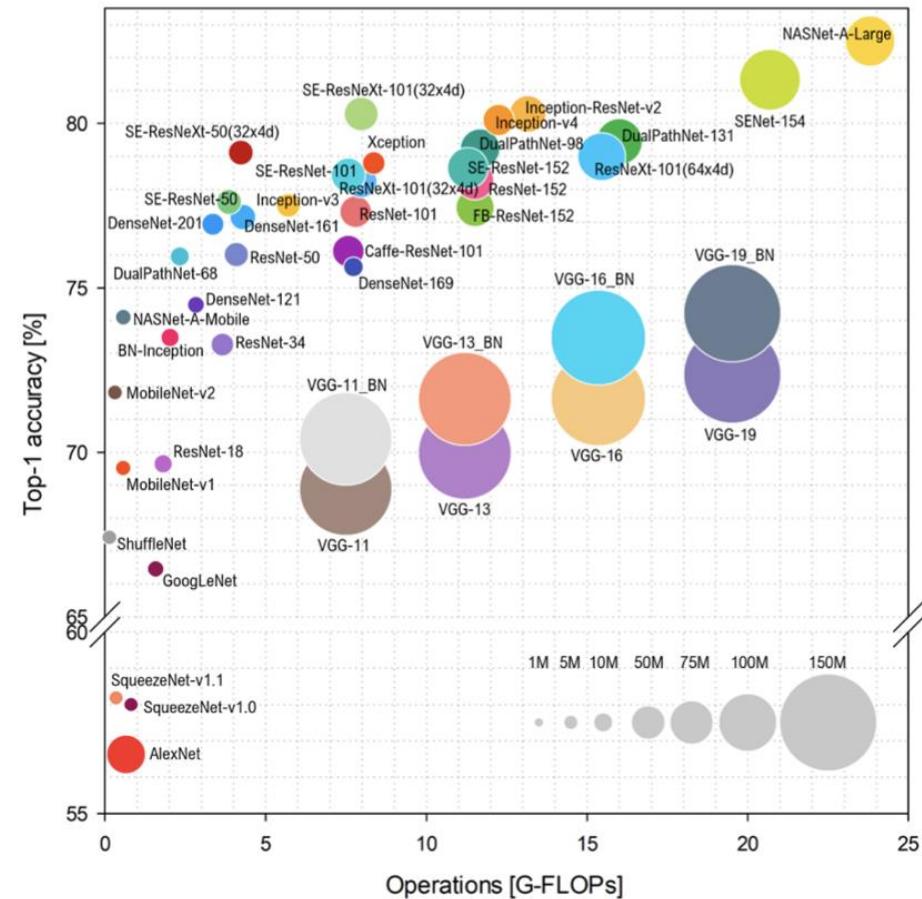
¹Department of Informatics, Systems and Communication, University of Milano-Bicocca, 20126 Milan, Italy

²LIP6, CNRS, Sorbonne Université, 75005 Paris, France

Corresponding author: Luigi Celona (luigi.celona@disco.unimib.it)

This paper presents an in-depth analysis of the majority of the deep neural networks (DNNs) proposed in the state of the art for image recognition.

For each DNN, multiple performance indices are observed, such as recognition accuracy, model complexity, computational complexity, memory usage, and inference time.



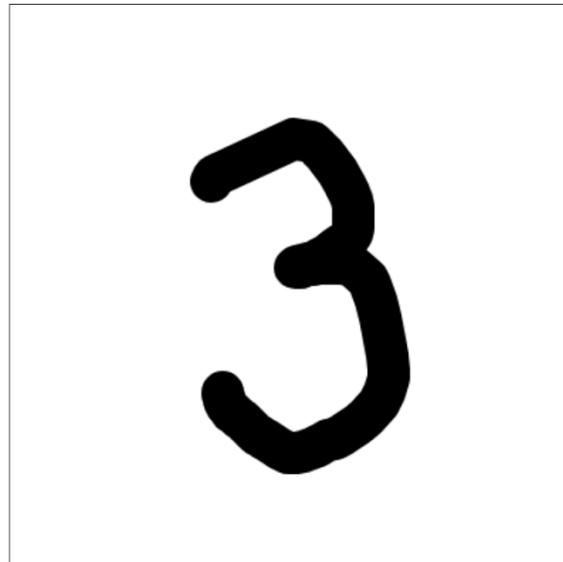
Workshop 2 – CNN applied for MNIST



<https://colab.research.google.com/drive/18Hr-Bz6W1FcgCe67vXwdKviPn8W8Ibj?usp=sharing>

MNIST Web Demo

MNIST Web Demo



checkbox: MNIST Preprocessing Thinner Black Thicker Black White Stroke Undo Stroke Reset

Image	0	1	2	3	4	5	6	7	8	9
	0.0000	0.0000	0.0000	1.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000

https://ufal.mff.cuni.cz/~courses/npfl129/2425/demos/mnist_web.html

Workshop 3 – CNN Applied to the Sign Language MNIST Dataset



<https://www.kaggle.com/code/rodsalasf/simposio-estadistica-taller-3-cnn>

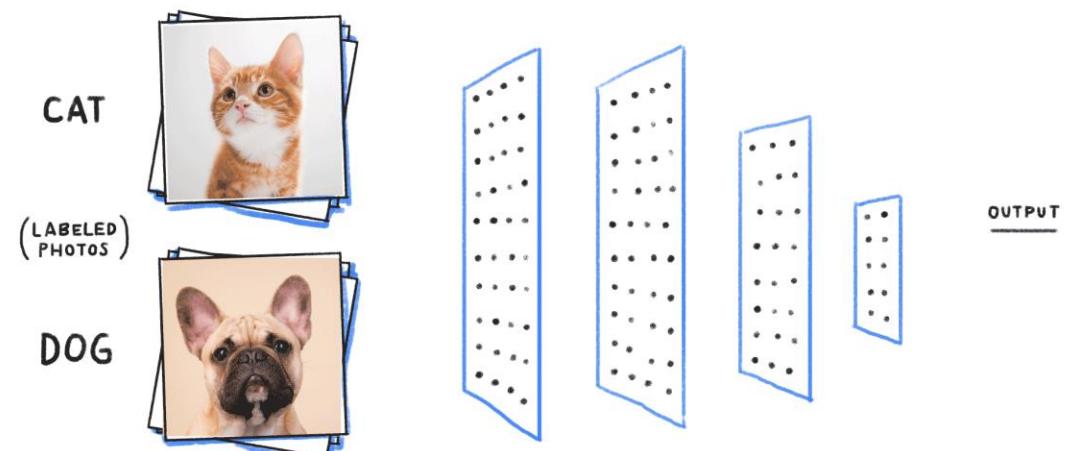
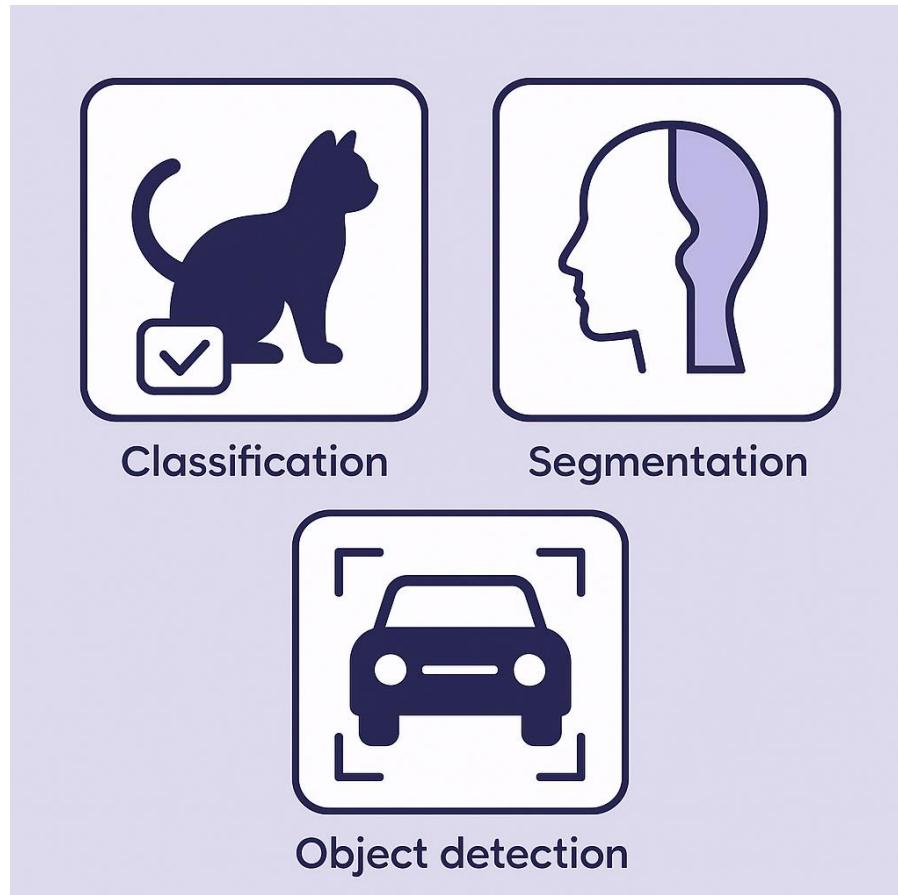
Dr. Ing. Rodrigo Salas (UV-MEDING-iHealth)



<https://www.kaggle.com/code/madz2000/cnn-using-keras-100-accuracy>

Deep Learning Applied for Image Processing

Deep Learning applied for Image Processing



Deep Learning Tasks in Computer Vision

Classification



Cat

Semantic
Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

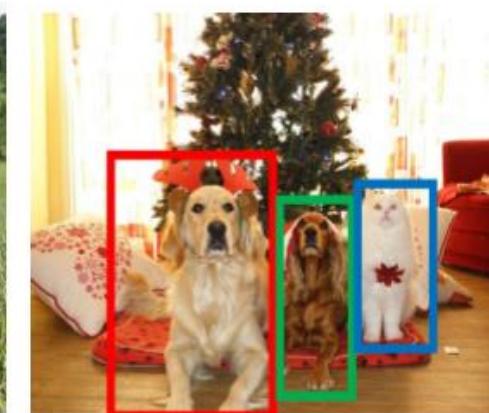
Classification
+ Localization



CAT

Single Object

Object
Detection



DOG, DOG, CAT

Multiple Object

Instance
Segmentation

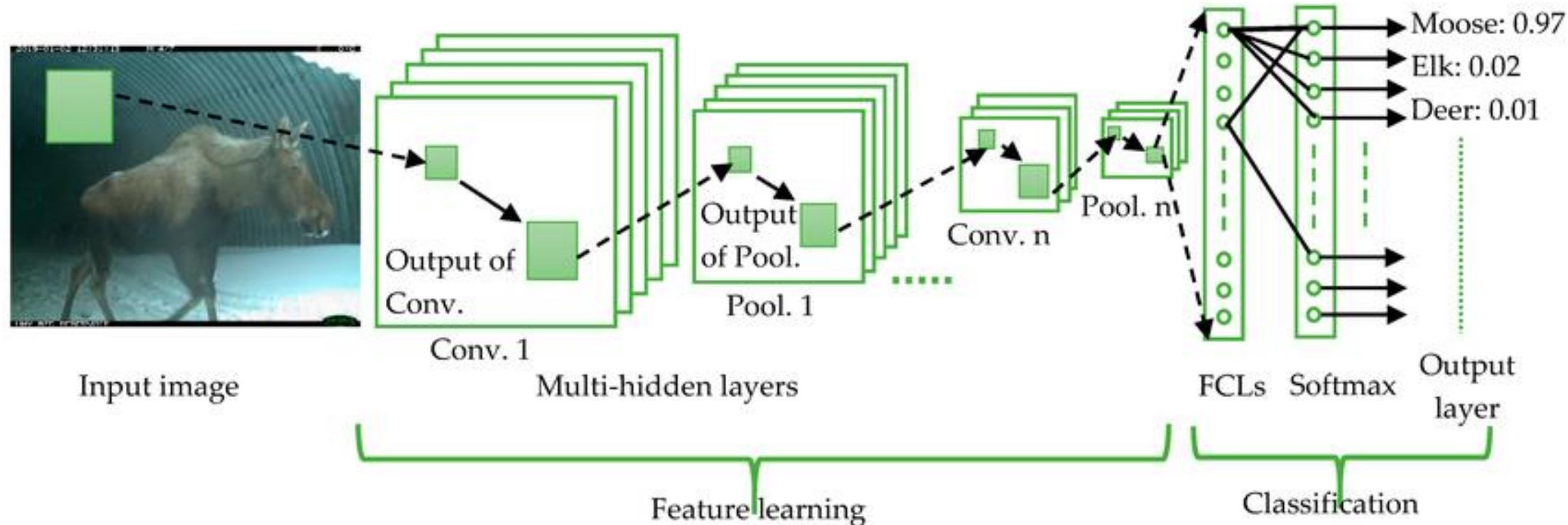


DOG, DOG, CAT

This image is CC0 public domain

Comparison of semantic segmentation, classification and localization, object detection and instance segmentation (Li, Johnson and Yeung, 2017)

Image Classification

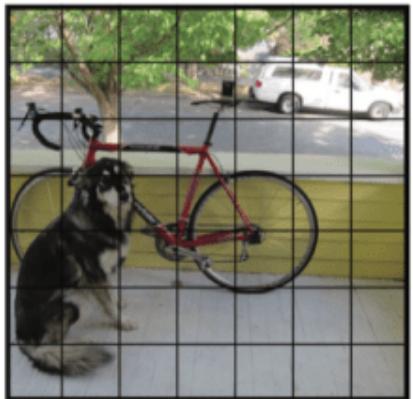


The network extracts **visual features** from the image through convolution and pooling layers. These representations are then processed in fully connected layers to assign probabilities to each class. The model learns to identify patterns that allow it to automatically classify the content of the image.

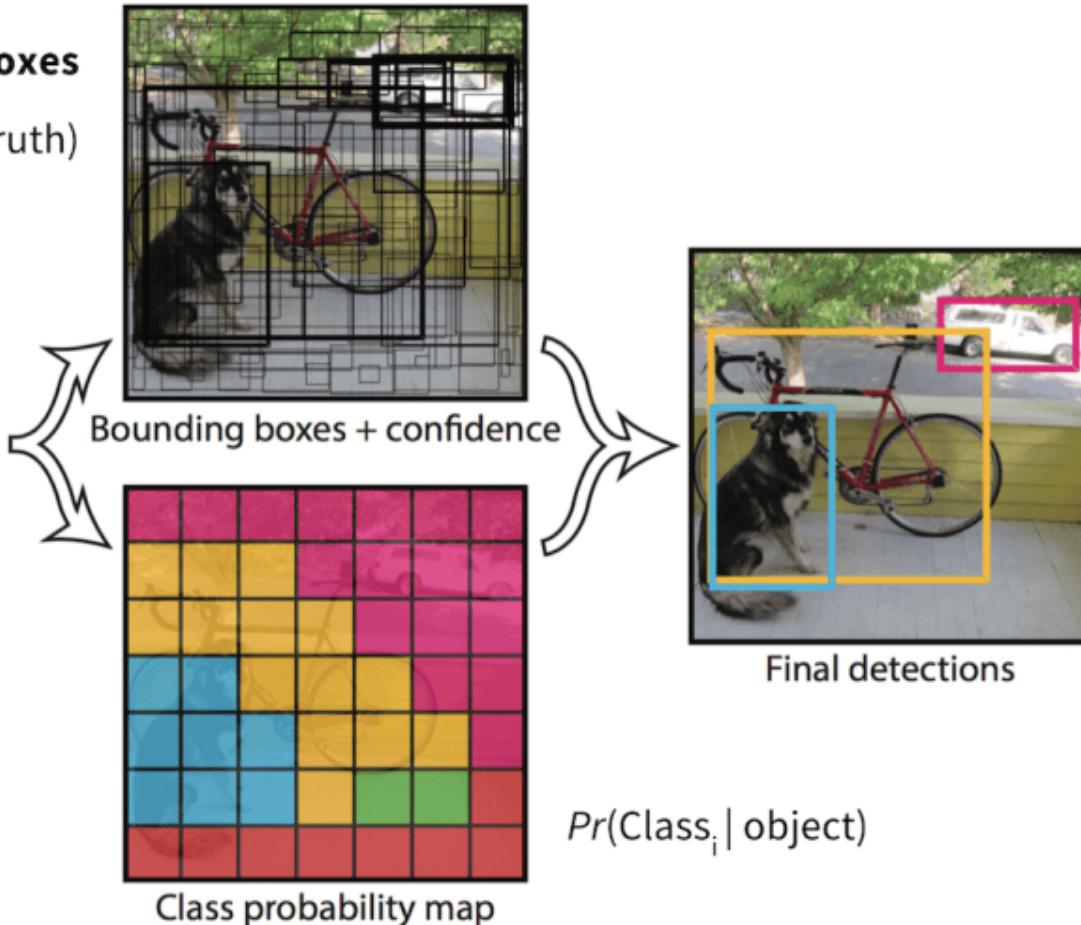
Object Detection

$S \times S \times B$ bounding boxes

confidence = $Pr(\text{object}) \times \text{IoU}(\text{pred, truth})$

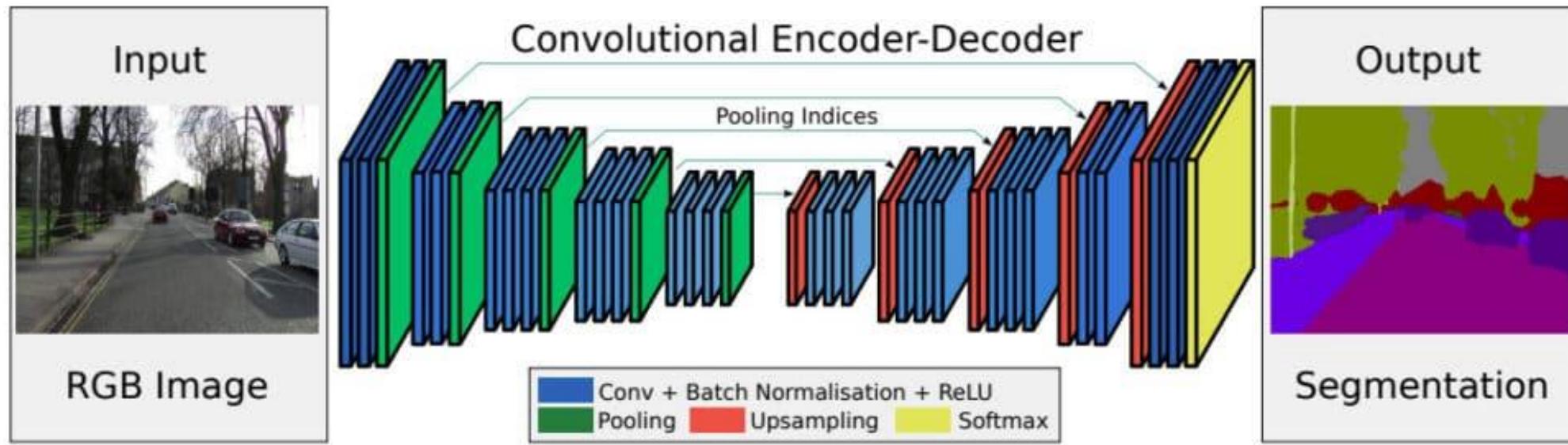


$S \times S$ grid on input



Object detection combines localization and classification: the model divides the image into a grid, predicts bounding boxes, and assigns a class probability to each one. Using metrics such as **IoU** and **confidence scores**, multiple objects can be identified and located within a single image.

Semantic Segmentation



Semantic segmentation assigns a label to every pixel in the image, identifying which class each region belongs to (e.g., street, car, person).

This process is performed using **encoder–decoder networks**, where the encoder extracts features and the decoder reconstructs a class map with the same resolution as the original image.

U-NET Applied for Semantic Segmentation



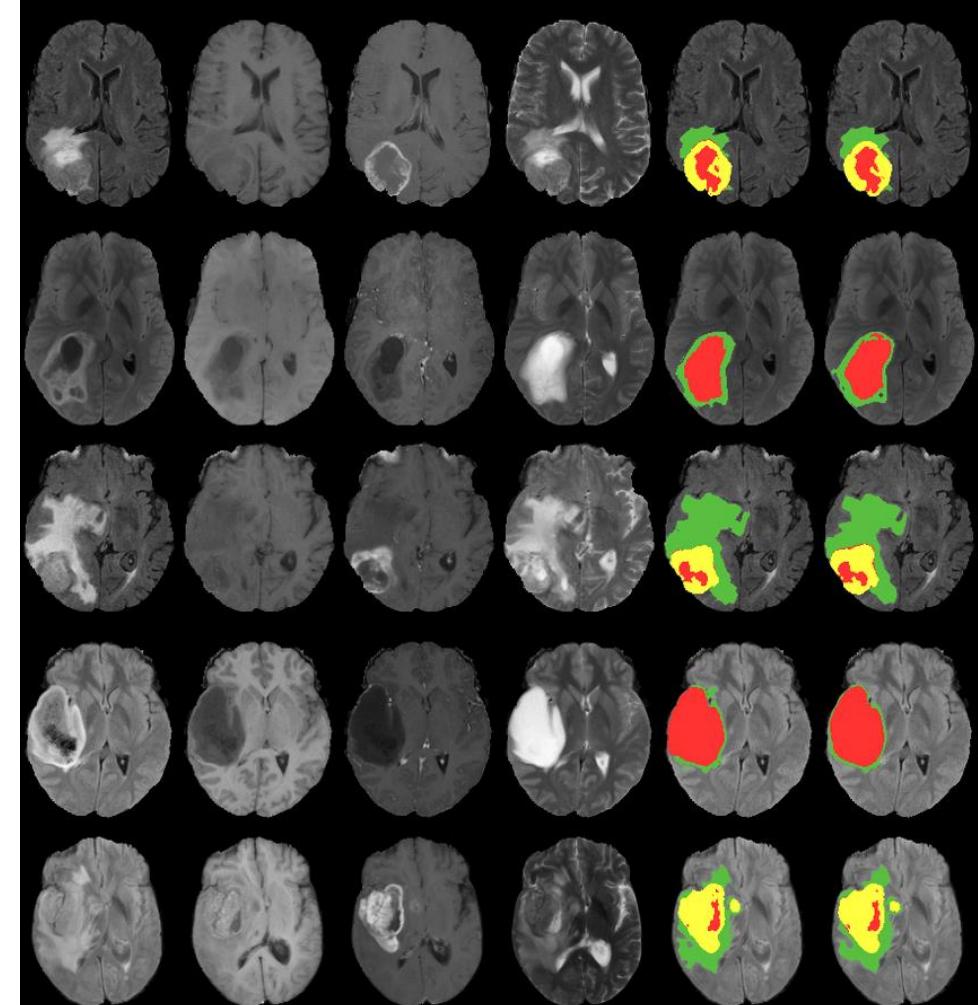
MSc. Student
Gabriel Guerra

Material desarrollado con el apoyo del estudiante del
Magister en Ciencias e Ingeniería para la Salud,
Universidad de Valparaíso
Sr. Gabriel Guerra

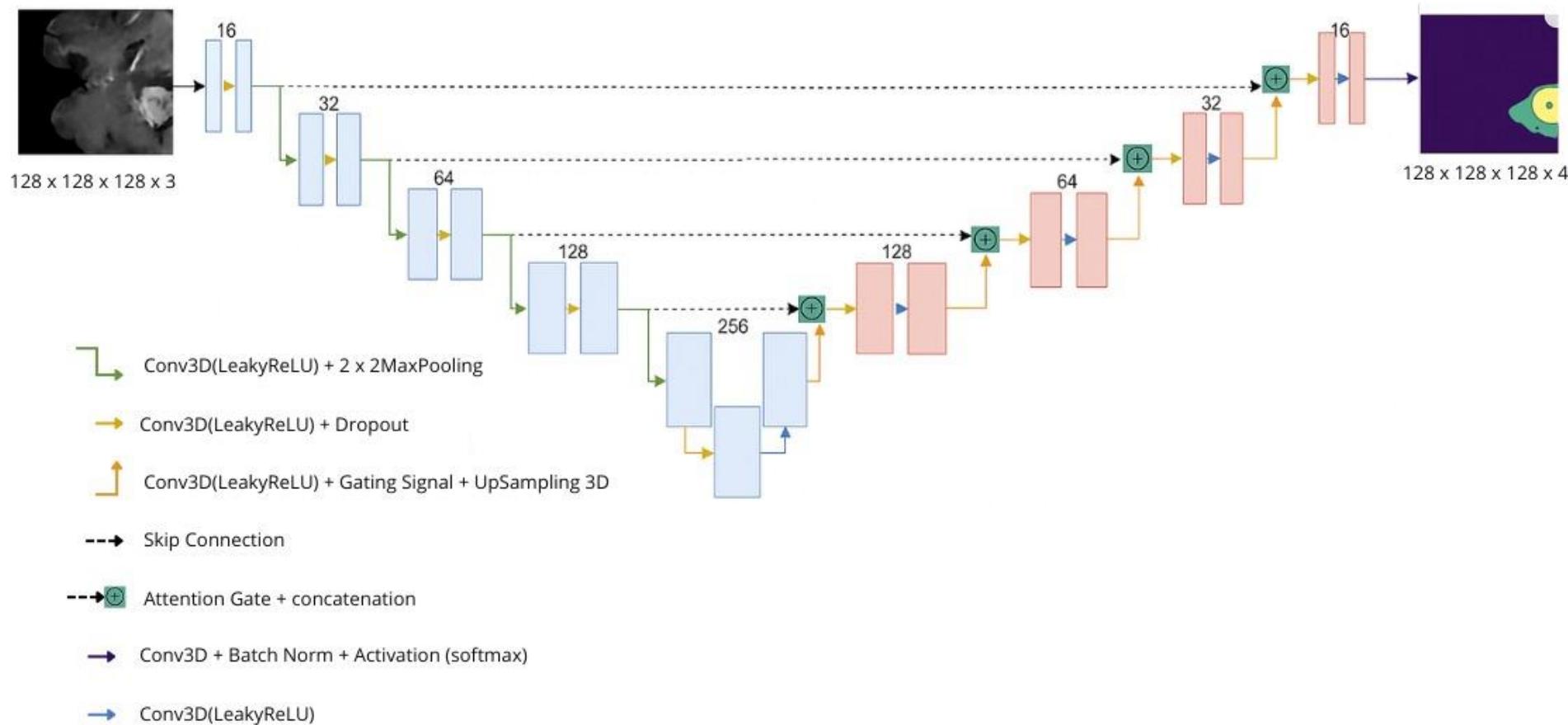
Segmentación de tumores cerebrales



BraTS es un conjunto de datos público que contiene resonancias magnéticas cerebrales multimodales (T1, T1c, T2, FLAIR) de pacientes con gliomas. Incluye segmentaciones manuales de tres regiones tumorales: tumor completo, tumor realizado y tumor no realizado. Su objetivo es evaluar algoritmos de segmentación automática en contextos clínicos realistas.

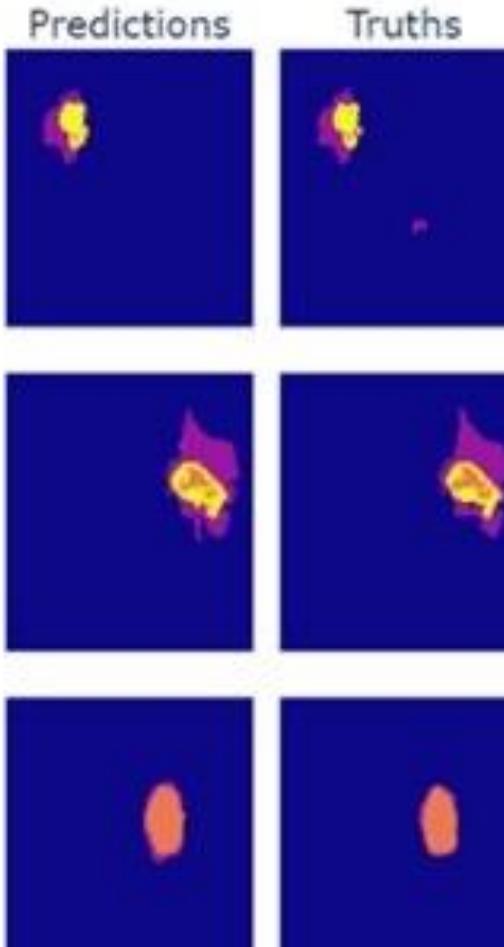


Attention U-net

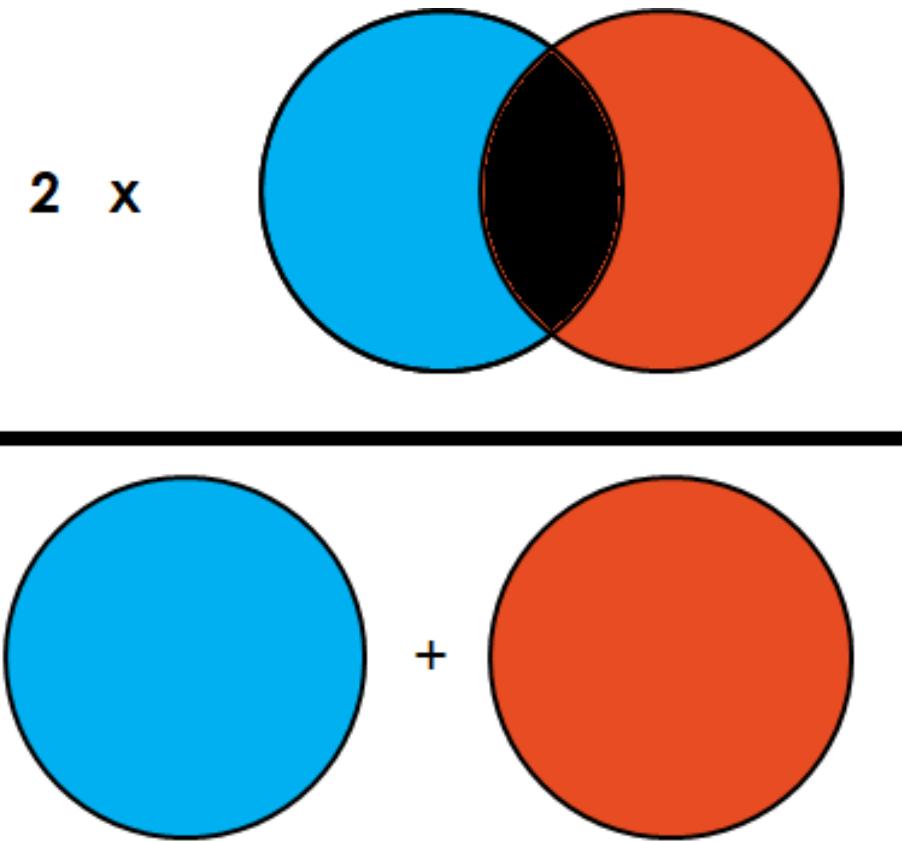


Gitonga, M. M. (2023). *Multiclass MRI Brain Tumor Segmentation using 3D Attention-based U-Net*. arXiv preprint arXiv:2305.06203. <https://arxiv.org/abs/2305.06203>

DICE Coefficient



$$\frac{2 * |X \cap Y|}{|X| + |Y|}$$



Workshop 4 – Tumor Semantic Segmentation



<https://colab.research.google.com/drive/1F5xrndrCGmRqiyZbsMDymTnrrHXWjcV8?usp=sharing>

The Black Box Problem

The Black Box Problem of Artificial Intelligence



The Black Box Problem in Artificial Intelligence



Lack of Transparency

No clear explanation of how predictions are made.



Limited Trust

Difficult for medical professionals to rely on outputs they cannot understand.



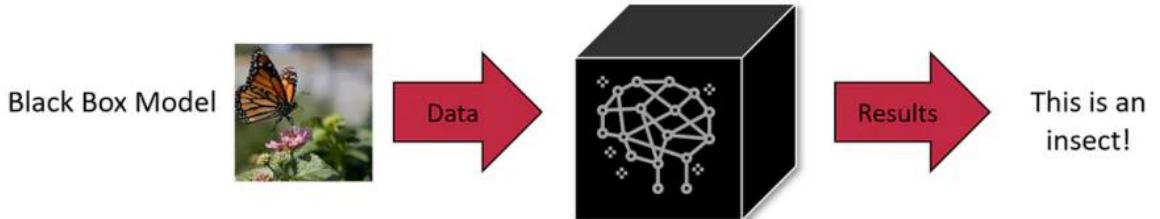
Bias and Hidden Errors

Potentially dangerous if incorrect decisions go undetected.

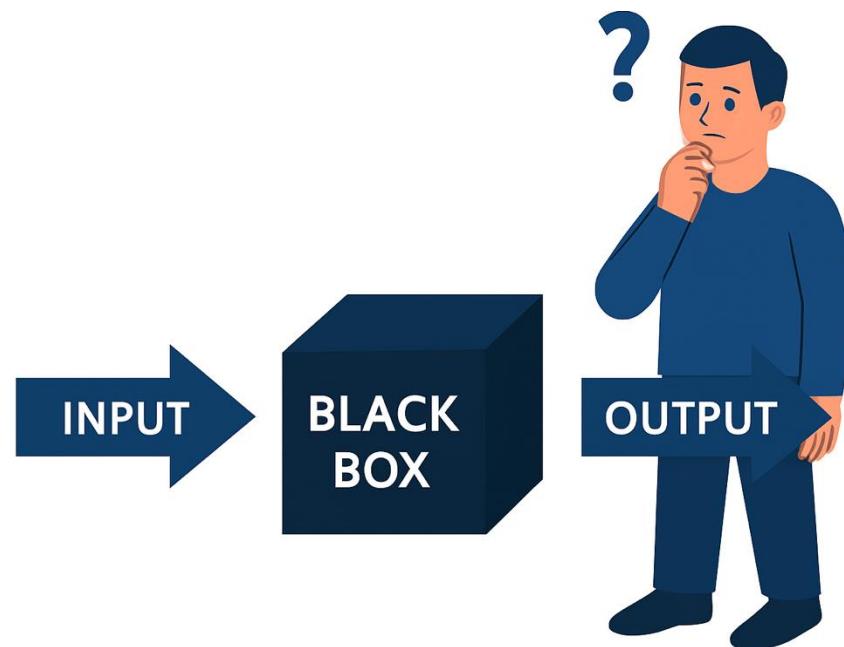


Ethical Concerns

Users need understandable reasoning behind critical decisions.



<https://www.unite.ai/the-black-box-problem-in-langs-challenges-and-emerging-solutions/>



Huskies vs Wolf

Explain the Prediction



Predicted: Wolf
True: Wolf



Predicted: Husky
True: Husky



Predicted: Husky
True: Husky



Predicted: Wolf
True: Wolf



Predicted: Wolf
True: Wolf



Predicted: Wolf
True: Wolf



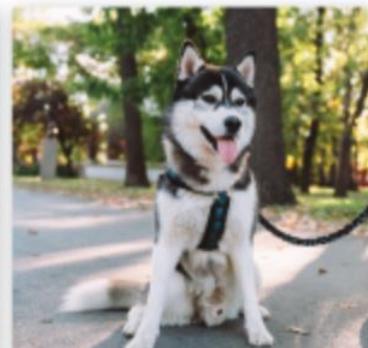
Predicted: Husky
True: Wolf



Predicted: Wolf
True: Wolf



Predicted: Wolf
True: Husky



Predicted: Husky
True: Husky

Besse, Philippe & Castets-Renard, Céline & Garivier, Aurélien & Loubes, Jean-Michel. (2018). Can Everyday AI Be Ethical? Machine Learning Algorithm Fairness (English version). 10.13140/RG.2.2.22973.31207.

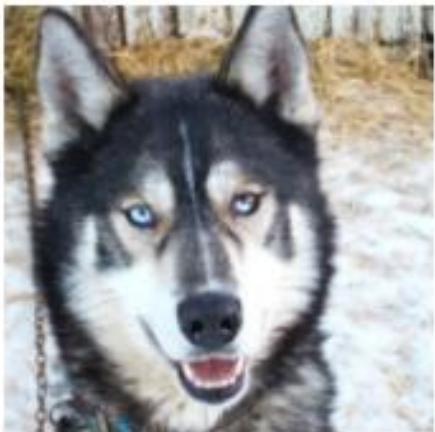
Huskies vs Wolf

“Why Should I Trust You?” Explaining the Predictions of Any Classifier

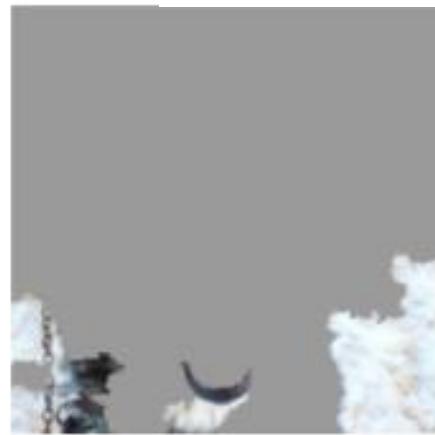
Marco Tulio Ribeiro
University of Washington
Seattle, WA 98105, USA
marcotcr@cs.uw.edu

Sameer Singh
University of Washington
Seattle, WA 98105, USA
sameer@cs.uw.edu

Carlos Guestrin
University of Washington
Seattle, WA 98105, USA
guestrin@cs.uw.edu



(a) Husky classified as wolf



(b) Explanation

Figure 11: Raw data and explanation of a bad model’s prediction in the “Husky vs Wolf” task.

	Before	After
Trusted the bad model	10 out of 27	3 out of 27
Snow as a potential feature	12 out of 27	25 out of 27

Table 2: “Husky vs Wolf” experiment results.



(a) The image of the Arctic Wolf in the background of snow is correctly classified.



(b) The image of the Husky in the background of snow is correctly classified.



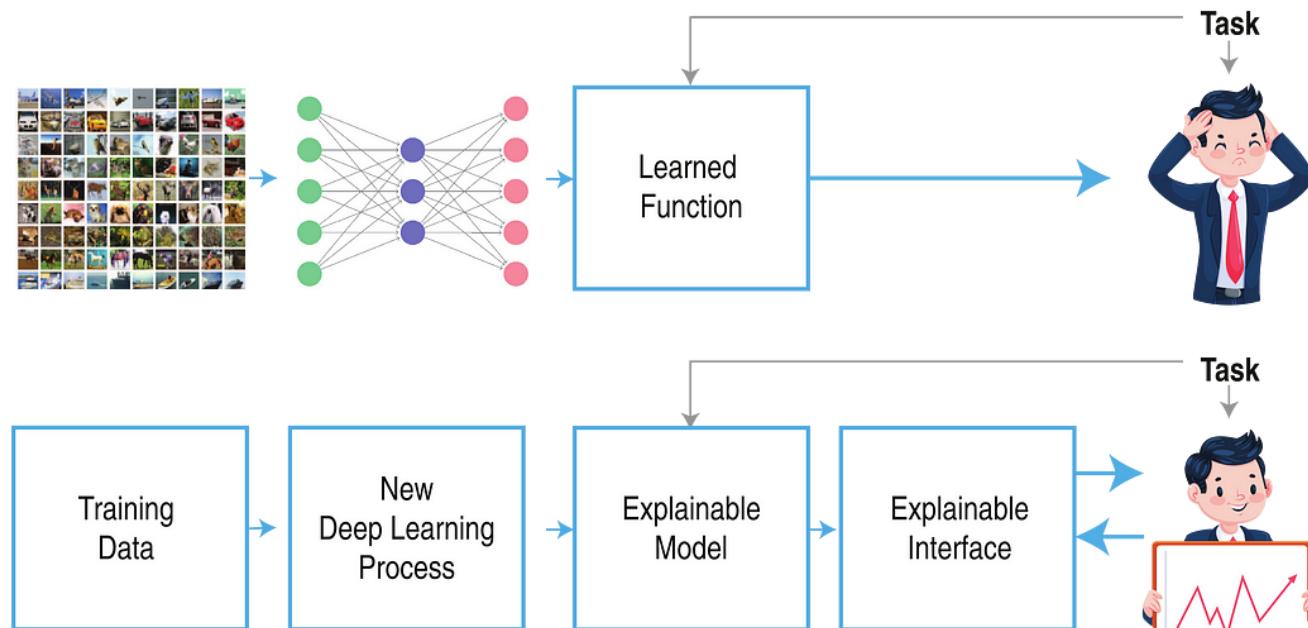
(c) The image of the Husky in the background of snow is correctly classified.

¡Open the Black Box!



eXplainable Artificial Intelligence **XAI**

What is Explainable Artificial Intelligence



Benefits of Explainable AI

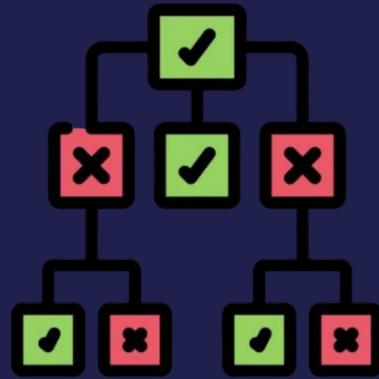
- Enables the deployment of trustworthy and understandable AI models.
- Enables continuous evaluation for faster AI outcomes.
- Reduces risk, bias, and auditing costs in AI systems.

Explainable Artificial Intelligence(XAI) is a set of processes and methods that allows human users to comprehend and trust the results and output created by AI algorithms.

<https://medium.com/deeppviz/what-is-xai-explainable-ai-and-visualization-part-10-da41c981c5fa>

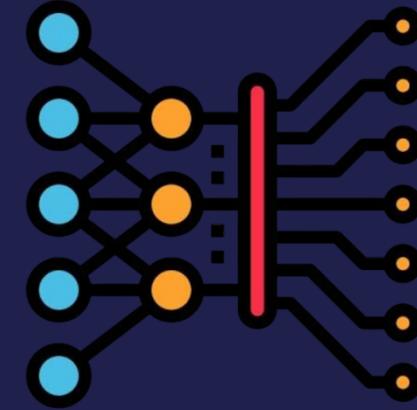
Interpretability vs Explainability

INTERPRETABILITY



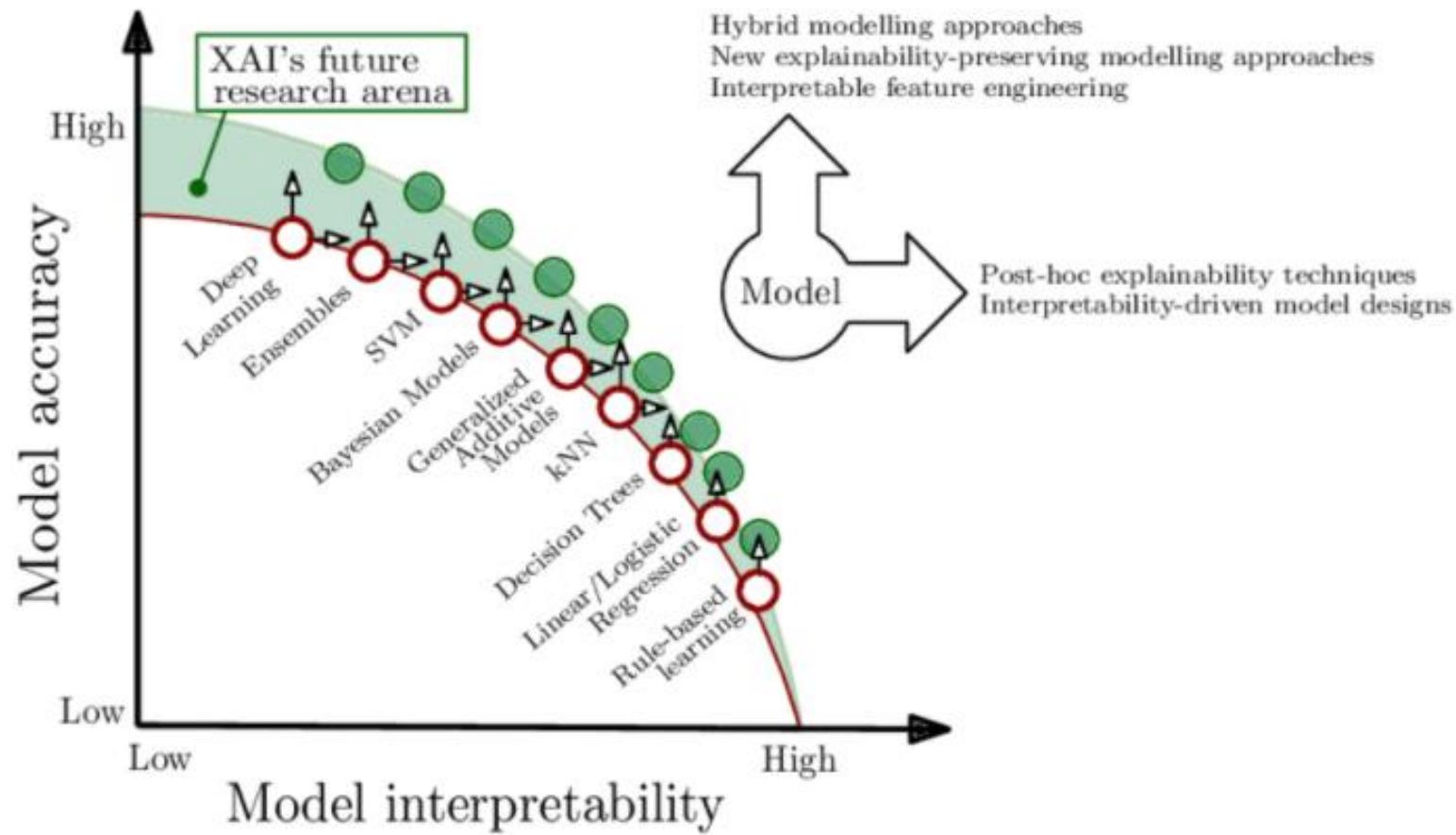
Describes *how* the model works internally, from the **developer's perspective**.

EXPLAINABILITY



Explains *why* the model made a decision, from the **user's perspective**.

Open the Black Box



<https://doi.org/10.1016/j.colec.2024.101629>

Post-Hoc vs. Ante-Hoc in XAI

Refers to models designed to be explainable from the start.

ANTE-HOC

Their internal structure allows for understanding without additional techniques.

Limited to specific models



Inherent model transparency

Applicable to any model



External explanation generation

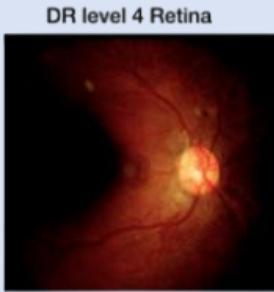
Applies explanatory methods after training complex models.

POST-HOC

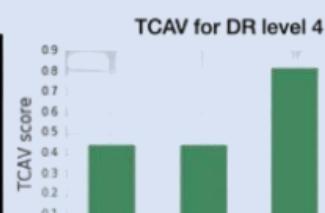
Aims to justify or interpret decisions made by black-box models.

Post-Hoc Explainability Approaches

<https://doi.org/10.48550/arXiv.1711.11279>

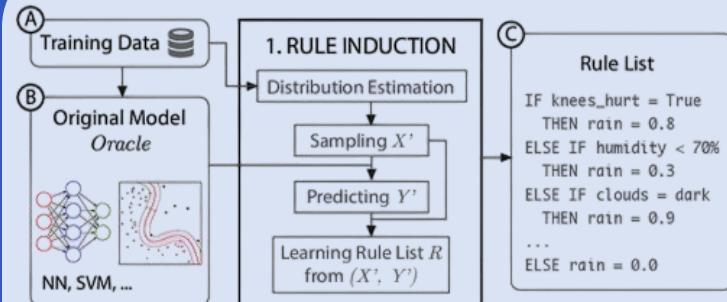


DR level 4 Retina



Numeric Explanations

<http://dx.doi.org/10.1109/TVCG.2018.2864812>

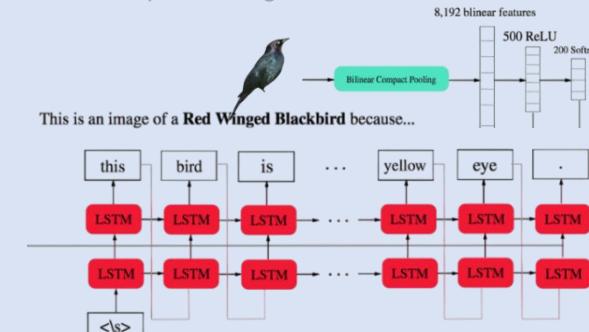


Rule-Based Explanations

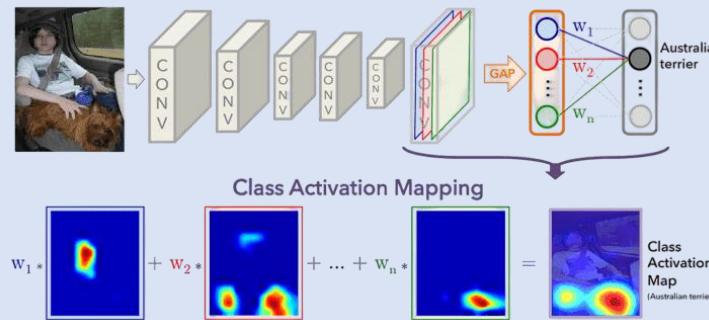
<https://doi.org/10.3390/make3030032>



This is an image of a **Red Winged Blackbird** because...



Textual Explanations



Visual Explanations

<https://doi.org/10.48550/arXiv.1512.04150>

Q: Is this a healthy meal? Textual Justification Visual Pointing



A: No

...because it is a hot dog with a lot of toppings.



A: Yes

...because it contains a variety of vegetables on the table.



Mixed Explanations

<https://doi.org/10.3390/make3040048>

Visual Explainability in Deep Learning



MSc. Student
Esthefanía Astargo

Material desarrollado con el apoyo de la estudiante
del Magister en Ciencias e Ingeniería para la Salud,
Universidad de Valparaíso
Sr. Esthefanía Astargo

XAI in Images

	Saliency Map Visualization	 <p>Simonyan et al., 2013</p>
Image	Shapley Value Importance	<p>Ghorbani et al., 2020</p>
	Concept Attribution	<p>Graziani et al., 2020</p>

<https://spectra.mathpix.com/article/2021.09.00007/demystify-post-hoc-explainability>

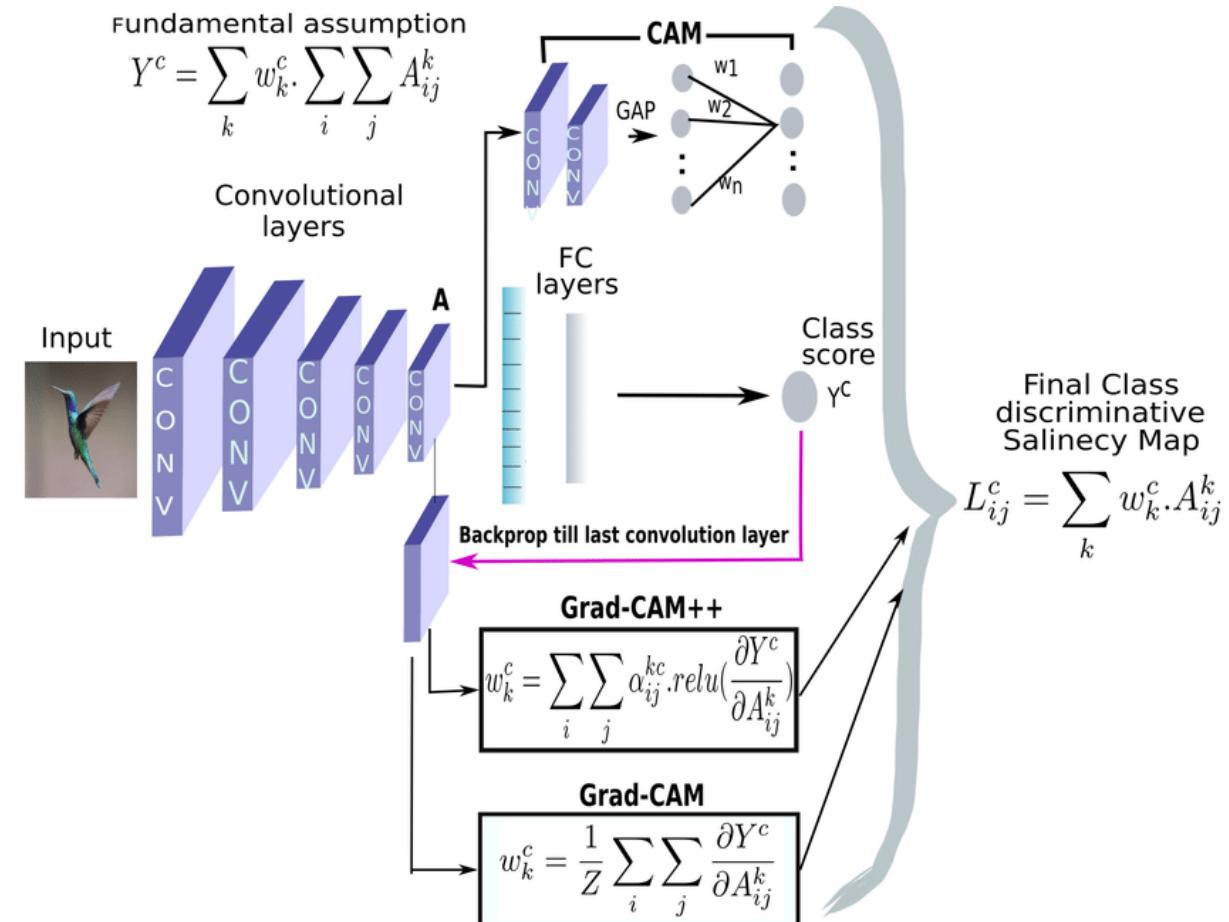
Grad-CAM: Gradient-weighted Class Activation Mapping

Grad-CAM es un método visual post-hoc y específico de modelo, diseñado para redes convolucionales profundas (CNNs).

Se basa en los gradientes del modelo para identificar qué regiones del input activaron más una clase específica.

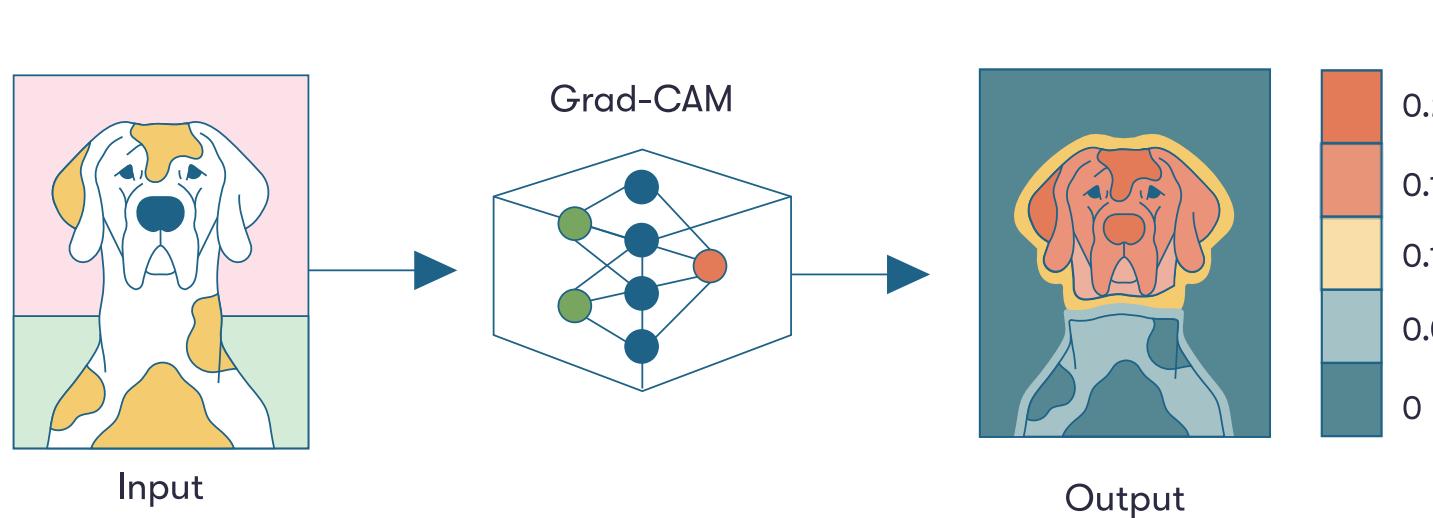
Utiliza los mapas de activación de la última capa convolucional, conservando la información espacial relevante.

Calcula un mapa ponderado de relevancia combinando gradientes y activaciones, destacando zonas clave.



<https://doi.org/10.48550/arXiv.1710.11063>

Grad-CAM: Gradient-weighted Class Activation Mapping

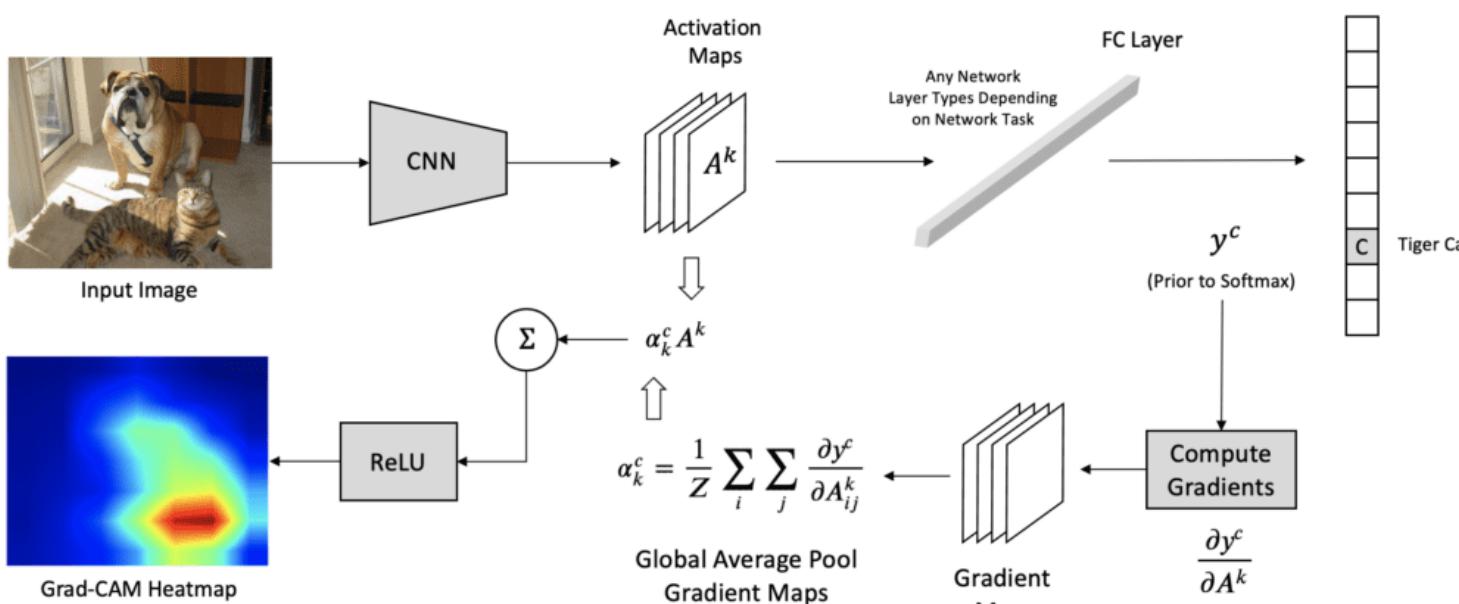


<https://courses.mnnlearn.com/en/courses/trustworthy-ai/preview/explainability/types-of-explainable-ai/>

Grad-CAM aplicado a clasificación de imágenes

- El modelo recibe como entrada la imagen de un perro y realiza una predicción de clase.
- Grad-CAM propaga hacia atrás los gradientes desde la salida hasta la última capa convolucional.
- A partir de los mapas de activación y gradientes, se genera un mapa de calor.
- El heatmap superpuesto indica qué regiones del perro fueron más relevantes para la predicción.

Grad-CAM paso a paso



https://xai-tutorials-readthedocs-io.translate.goog/en/latest/_model_specific_xai/Grad-CAM.html?_x_tr_sl=en&_x_tr_t=es&_x_tr_h=es&_x_tr_pto=t

Paso 1: Forward Pass

- Obtener los mapas de características de la última capa convolucional y las salidas sin procesar antes de la función softmax.
- Estos mapas de características se denominan A^k , donde k es un mapa de características específico dentro de una capa convolucional.

Paso 2: Selección de la clase objetivo

- Se elige la clase c que se desea explicar (normalmente la clase predicha con la puntuación más alta) y se toma su valor de activación antes de la función softmax.
- y^c representa la activación asociada a la clase c antes de aplicar la función softmax.

Grad-CAM paso a paso

Paso 3: Calcular los gradientes

- Se calculan los gradientes de y^c con respecto a los mapas de características A^k , es decir:

$$\frac{\partial y^c}{\partial A^k}$$

Paso 4: Calcular el mapa Grad-CAM

- Se calcula un peso por cada mapa (media global de gradientes):

$$\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^k}$$

- Luego se combinan todos los mapas con sus pesos:

$$L_{Grad-CAM}^c = \text{ReLU}\left(\sum_k \alpha_k^c A^k\right)$$

- La función ReLU conserva solo las regiones que contribuyen positivamente.

https://xai--tutorials-readthedocs-io.translate.goog/en/latest/_model_specific_xai/Grad-CAM.html?_x_tr_sl=en&_x_tr_tl=es&_x_tr_hl=es&_x_tr_pto=tc

Workshop 5 – Using Grad-CAM for Explainability in Image Classification

Google
colab



<https://colab.research.google.com/drive/1qlAtUqSFkIVSK2LWjMBzvVM4nGAC1CvO?usp=sharing>

LIME: Local Interpretable Model-agnostic Explanations

Construye un modelo sustituto que aproxima al modelo complejo solo en torno a una predicción específica.

Genera variaciones del dato de entrada y observa cómo el modelo original responde a estos cambios.

Utiliza esas observaciones para entrenar un modelo simple que imita al modelo complejo localmente.

Este modelo sustituto revela qué características influyen más en esa predicción puntual.

Puede aplicarse a texto, imágenes, o datos tabulares, sin requerir acceso interno al modelo.

Matemáticamente, los modelos sustitutos se pueden expresar de la siguiente manera:

$$\xi(x) = \arg \min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

f es el modelo complejo original.

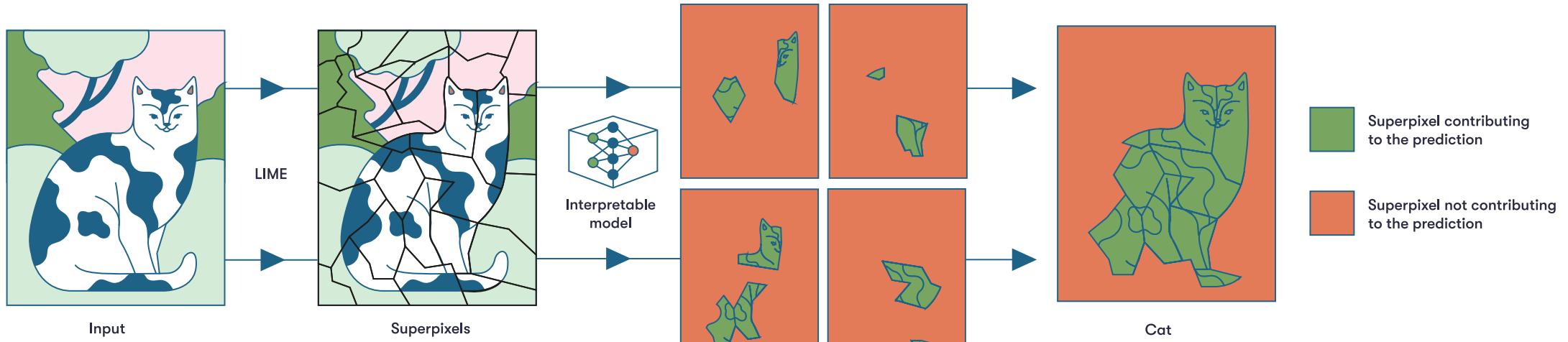
g es el modelo interpretable local.

$\mathcal{L}(f, g, \pi_x)$ mide cuán bien g aproxima a f cerca de x .

π_x define la vecindad local ponderando las instancias perturbadas.

$\Omega(g)$ penaliza la complejidad de g .

LIME: Local Interpretable Model-agnostic Explanations



LIME aplicado a clasificación de imágenes

Divide la imagen en superpíxeles, que agrupan regiones visuales coherentes y comprensibles.

Se generan múltiples versiones de la imagen desactivando distintas combinaciones de superpíxeles

El modelo original predice sobre cada versión modificada y LIME observa cómo cambian las salidas.

Luego, entrena un modelo interpretable para estimar qué superpíxeles son claves en la decisión.

Ejemplo: los superpíxeles verdes explican por qué el modelo clasificó la imagen como gato.

Workshop 6 - Uso de LIME para Explicabilidad en Clasificación de Imágenes

¡Escanéame!



[https://colab.research.google.com/drive/1pWRDK32SPg
BA3OD83_8qkTwsuBBkqqaz?usp=sharing](https://colab.research.google.com/drive/1pWRDK32SPgBA3OD83_8qkTwsuBBkqqaz?usp=sharing)

Material del Curso y la Conferencia



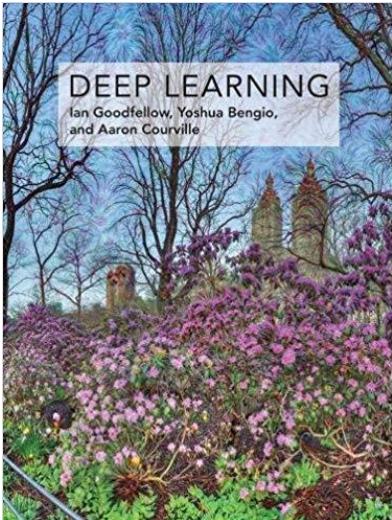
Github



<https://github.com/rodsalasf/>

Dr. Rodrigo Salas Fuentes
rodrigo.salas@uv.cl

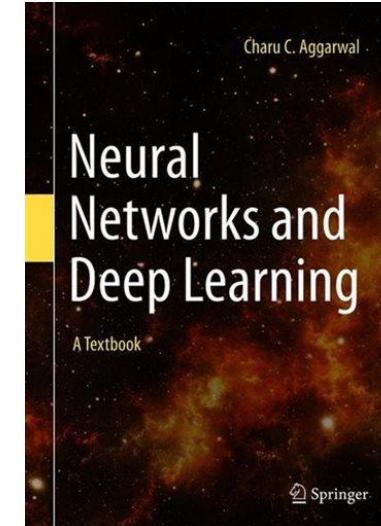
Deep Learning Books



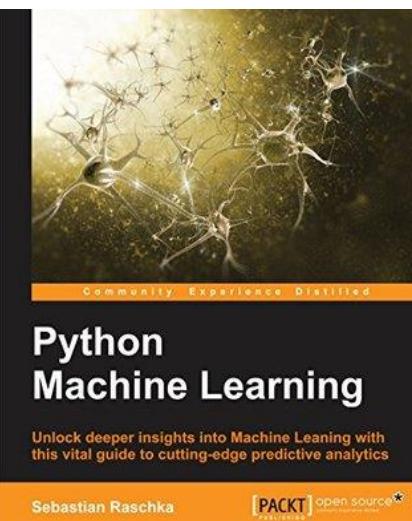
Ian Goodfellow



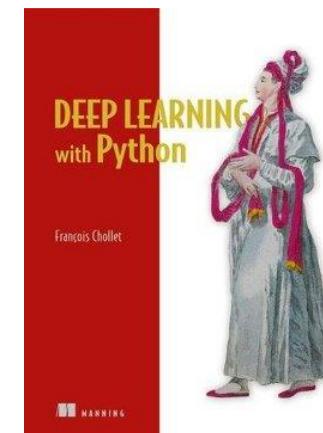
Yoshua Bengio **Aaron Courville**



Charu Aggarwal



Sebastian Raschka

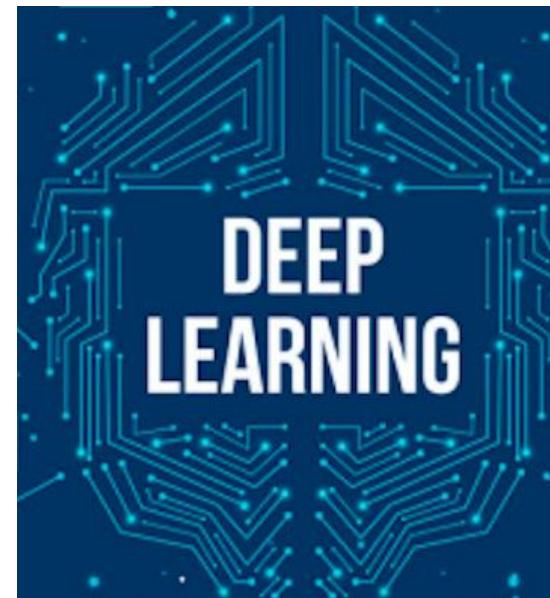


Francois Chollet

Toolbox for Deep Learning



PYTORCH
Deep Learning with PyTorch



Caffe

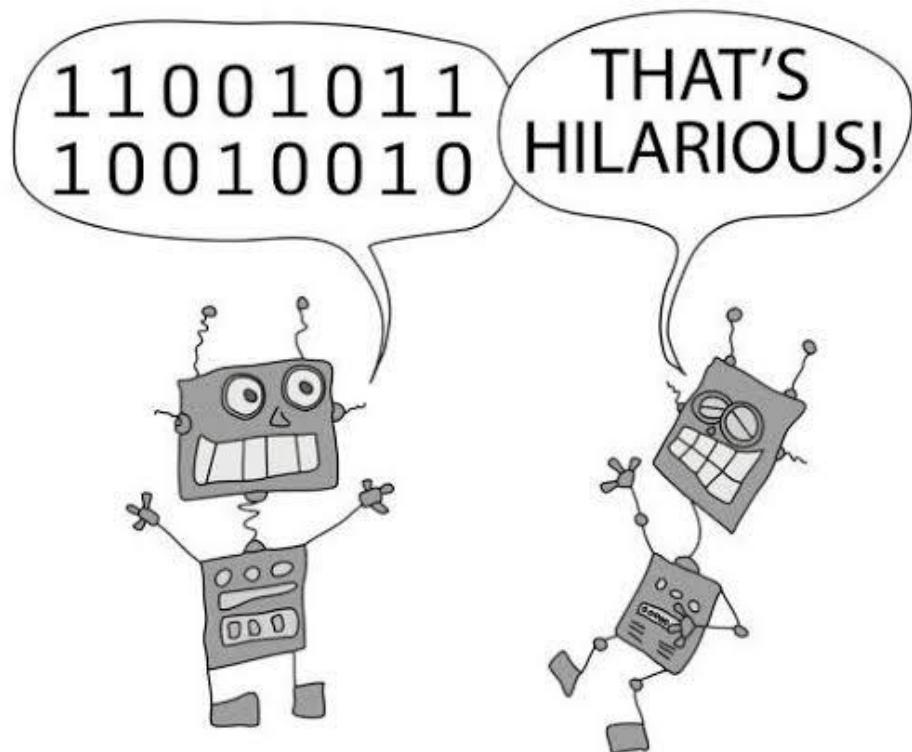


theano



dmlc
mxnet



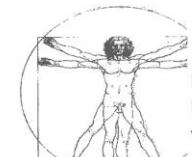


**Dr. Rodrigo Salas
Fuentes**
rodrigo.salas@uv.cl



Thanks for your Attention

-  LinkedIn: linkedin.com/in/rodrigo-salas-fuentes/
-  Google Scholar: scholar.google.com/citations?user=ZaqDIPcAAAAJ
-  ORCID: orcid.org/0000-0002-0350-6811
-  Email: rodrigo.salas@uv.cl



Ingeniería Civil Biomédica
Facultad de Ingeniería
Universidad de Valparaíso
www.biomedica.uv.cl