

MAC 0460 / 5832
Introduction to Machine Learning

22 – Deep learning for Computer Vision (part 2)

IME/USP (05/07/2021)

Image classification

Bag of visual words (BOW)

Inspiration comes from **NLP**

Bag-of-words: widely used for document classification/retrieval

Vocabulary with N words

First idea: Document D represented by a vector f_D of size N :

$f_D[i] = \text{number of occurrences of word } v_i \text{ in the document}$

Improvement: term frequency-inverse document frequency (TF-IDF)

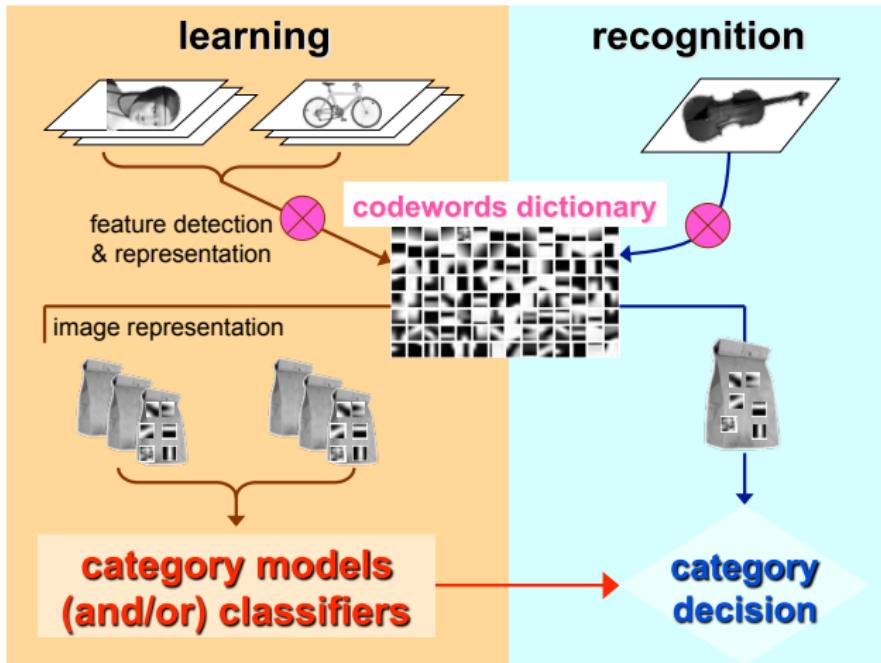
- normalization taking into account the frequency of the words
- $f_D[i]$ divided by the count of occurrences of v_i over all documents

Images (Computer Vision)

Bag of keypoints (Csurka, 2004)

1. **Create a dictionary** (visual patterns) from the set of training images
2. **Training:** compute the *bag-of-words* representation for each training image and train a classifier to discriminate object classes using the *bag-of-words* representation as features
3. **Classification:** compute the the *bag-of-words* representation of the image to be classified and apply the classifier

Bag-of-Words Model: Overview

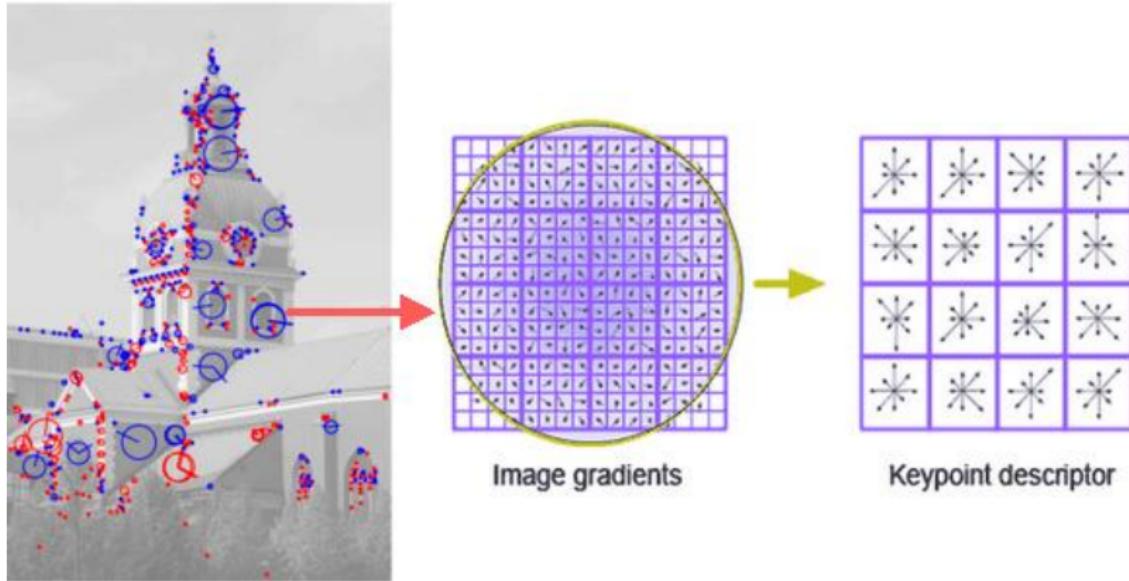


BoW – Dictionary creation

- Compute *keypoints* of the images
- Compute descriptors for each of the *keypoints*
- Csurka uses *Harris affine detector* for the detection of keypoints and SIFT for the description (feature vector of dimension 128)
- Descriptors are grouped using the K -means algorithm (parameter K is adjusted somehow)
- The dictionary consists of the set of centers of the K clusters (**visual words**)

Example of feature detection and description

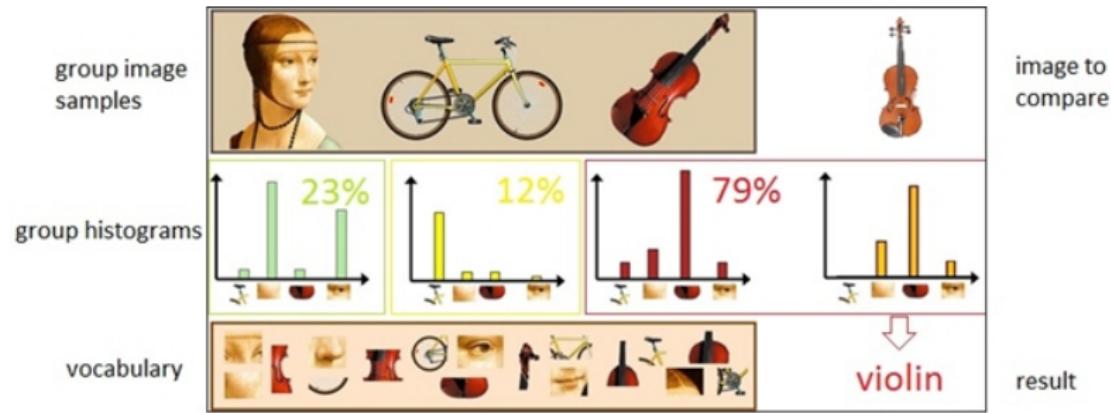
SIFT (Scale Invariant Feature Transform)



BoW – Training the classifier

- For each image, compute the set of descriptors
- Each descriptor is matched with the closest *visual word* in the dictionary
- Build a histogram of the *visual words* in the image (this is the *bag-of-words* of the image)
- The *bag-of-words* representation is the input feature used to train the classifier

BoW – Training and classification process



(<http://vgg.fiiit.stuba.sk/wp-uploads/2015/02/BOW.jpg>)

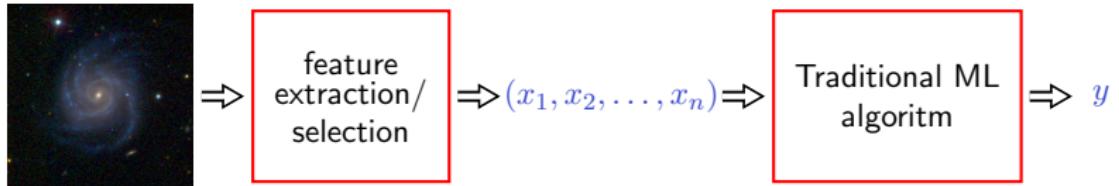
Computer Vision algorithms consisted of very complex processing pipelines

Examples we have seen:

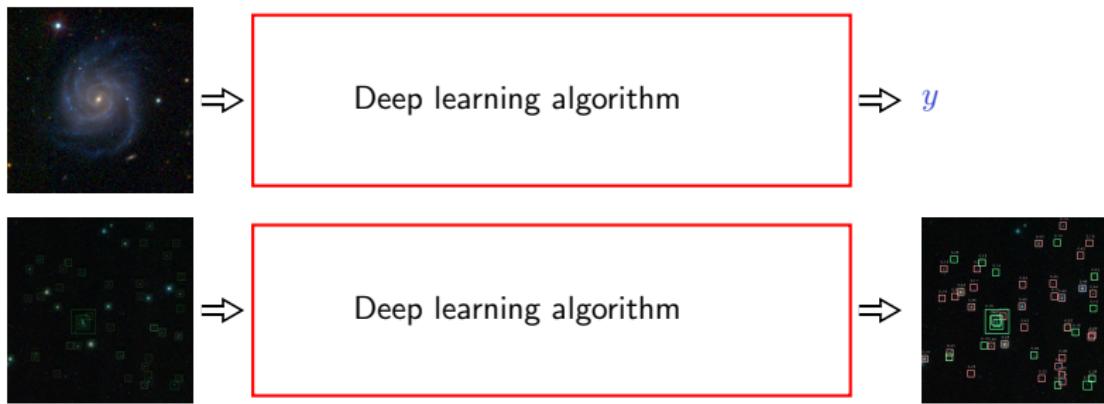
Object detection: Viola-Jones object detection framework

Image classification: Csurka's BoW model

Traditional ML



Deep learning: end-to-end processing



Neural network model for images

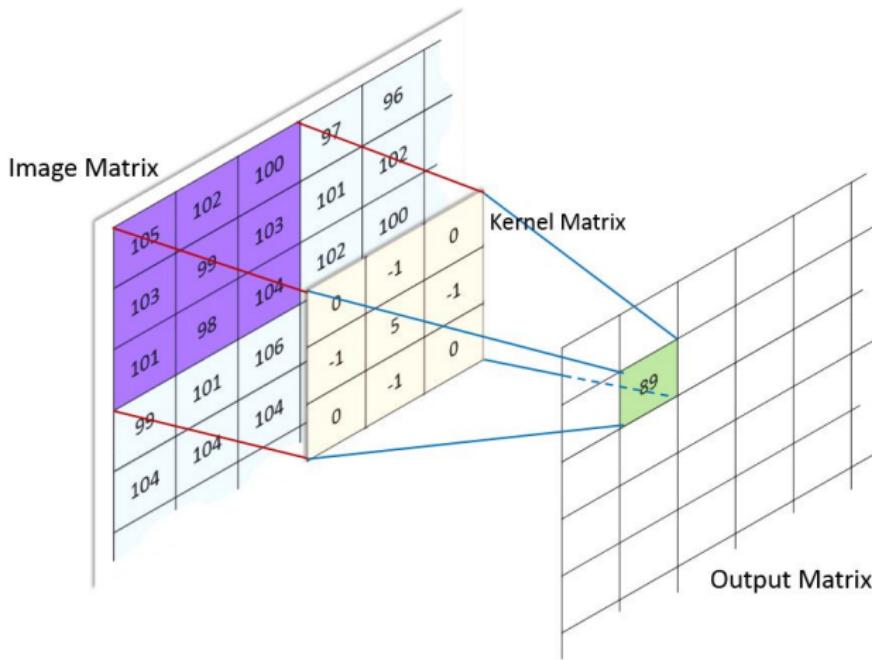
CNN (Convolutional Neural Networks)

They are based on the **convolution** operation

Local linear operation used in image transformations

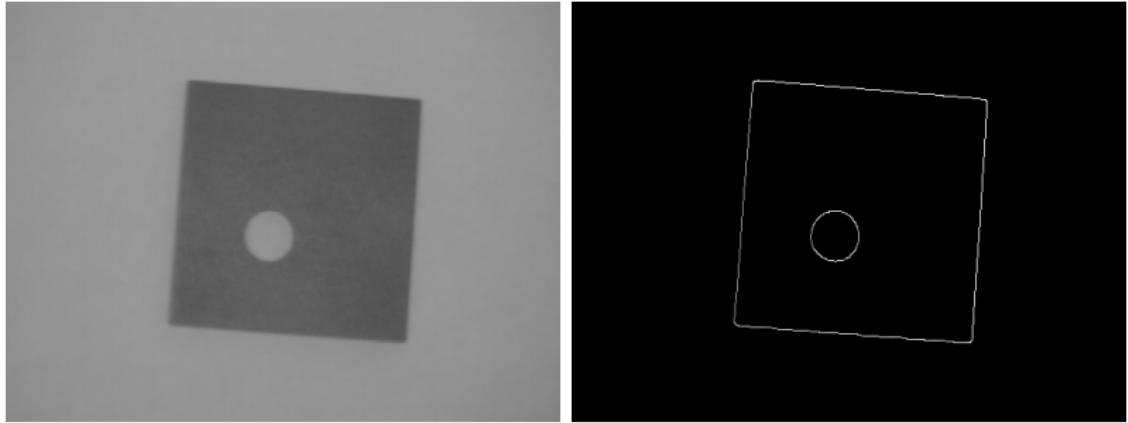
$$(f * g)[n] = \sum_{m=-M}^M f[n-m] g[m]$$

Convolution



(http://machinelearningguru.com/computer_vision/basics/convolution/image_convolution_1.html)

Examples of convolution kernels – edge detection



-1	-1	-1
-1	8	-1
-1	-1	-1

(<http://homepages.inf.ed.ac.uk/rbf/HIPR2/canny.htm>)

Examples of convolution kernels – edge detection



(<http://aishack.in/tutorials/image-convolution-examples/>)

Smoothing



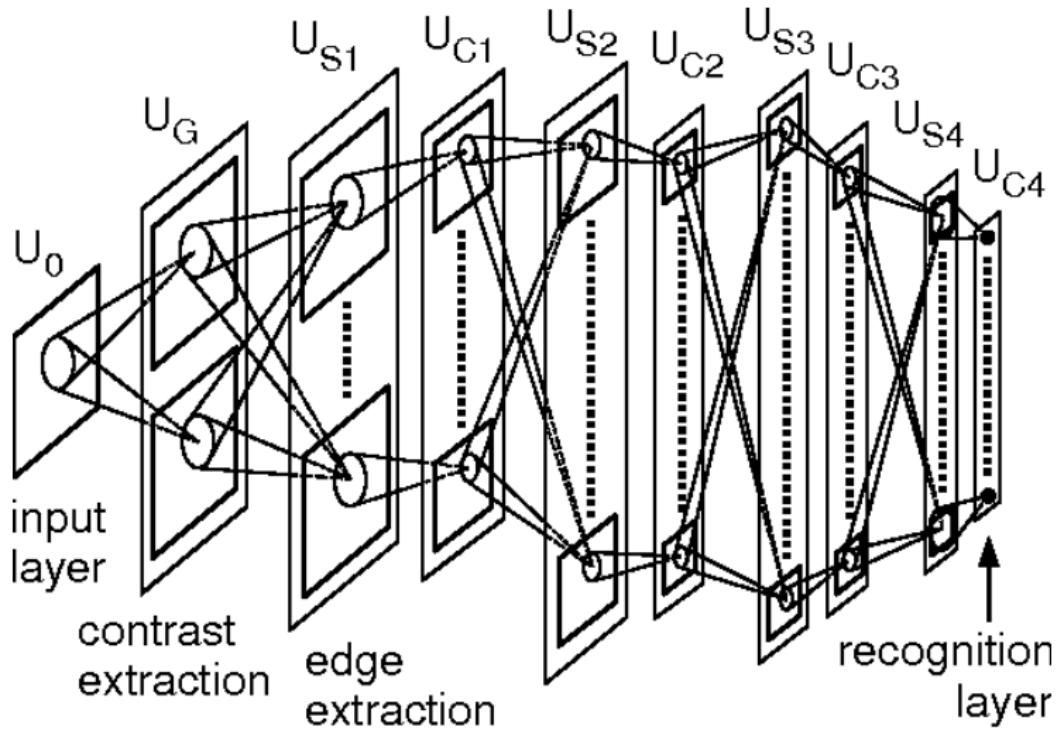
1/9	1/9	1/9
1/9	1/9	1/9
1/9	1/9	1/9

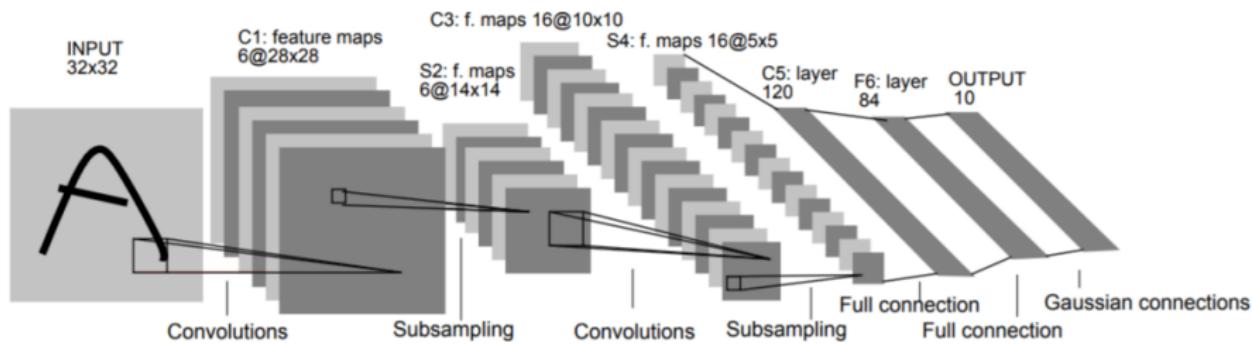
(<http://aishack.in/tutorials/image-convolution-examples/>)

Use of convolution in neural nets is not so recent:

- **Neocognitron** – Fukushima, 1990
Detailed explanation: <https://www.kiv.zcu.cz/studies/predmety/uir/NS/Neocognitron/en/index.html>
- **LeNet** - Yan Lecun, 1998, postal code recognition
<http://yann.lecun.com/exdb/lenet/index.html>

Neocognitron – Fukushima, 1990

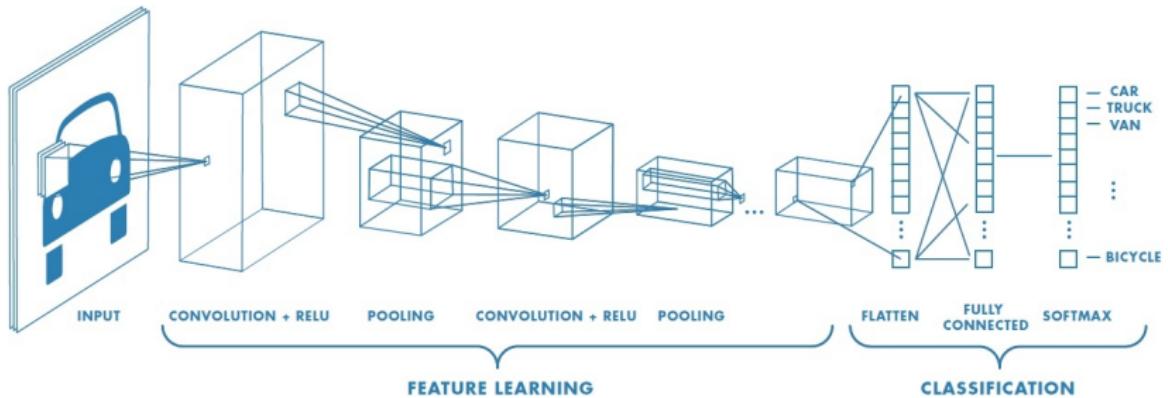




First CNN model, applied on postal code recognition

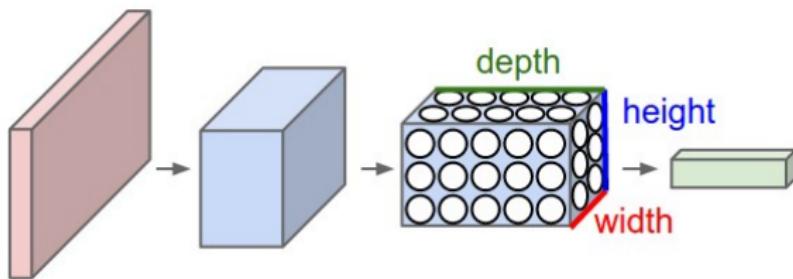
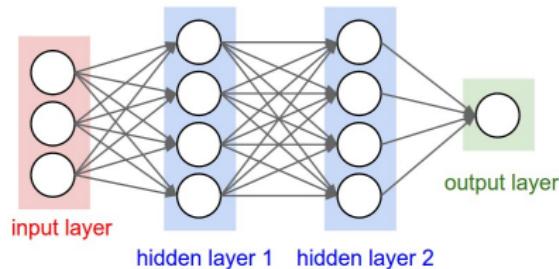
Basics of CNN

Typical architecture of Convnets



Source: https://iitmcvg.github.io/summer_school/DLSession3/

Graphical representation of ConvNet structure



```
INPUT -> [[CONV -> RELU]*N -> POOL?] *M -> [FC -> RELU]*K -> FC
```

Convolution on multiple channels

0	0	0	0	0	0	0	...
0	156	155	156	158	158	...	
0	153	154	157	159	159	...	
0	149	151	155	158	159	...	
0	146	146	149	153	158	...	
0	145	143	143	148	158	...	
...

Input Channel #1 (Red)

0	0	0	0	0	0	0	...
0	167	166	167	169	169	...	
0	164	165	168	170	170	...	
0	160	162	166	169	170	...	
0	156	156	159	163	168	...	
0	155	153	153	158	168	...	
...

Input Channel #2 (Green)

0	0	0	0	0	0	0	...
0	163	162	163	165	165	...	
0	160	161	164	166	166	...	
0	156	158	162	165	166	...	
0	155	155	158	162	167	...	
0	154	152	152	157	167	...	
...

Input Channel #3 (Blue)

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1



308

+

1	0	0
1	-1	-1
1	0	-1

Kernel Channel #2



-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



164

+

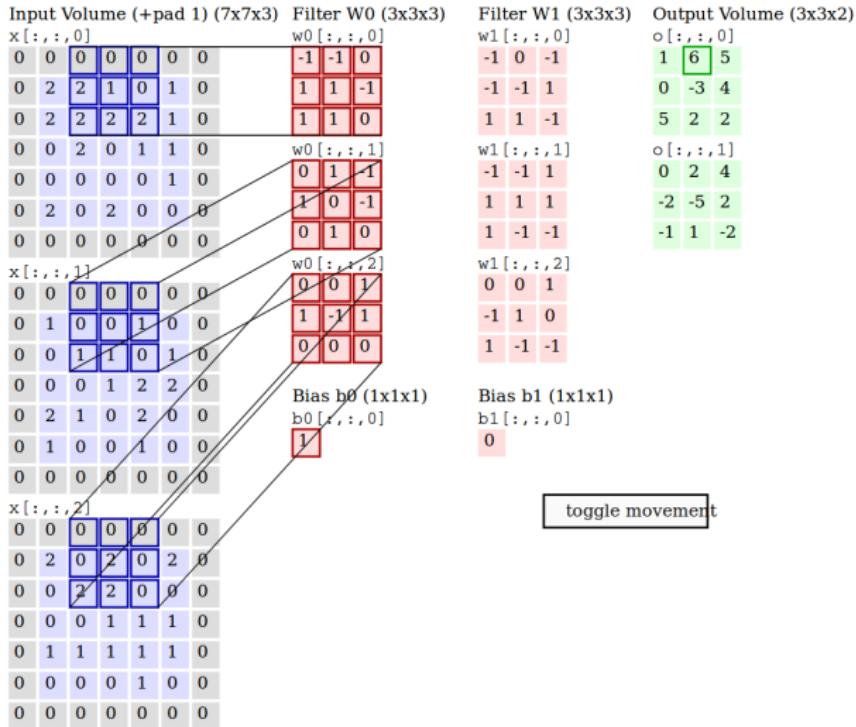
+ 1 = -25

Bias = 1

-25					...
					...
					...
					...
...

(http://machinelearningguru.com/computer_vision/basics/convolution/convolution_layer.html)

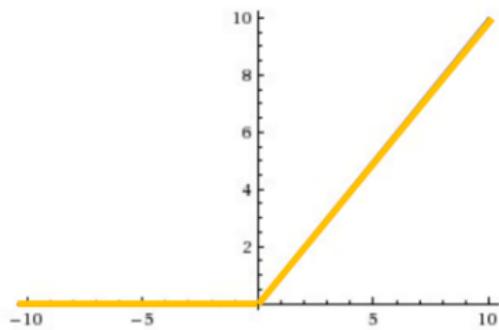
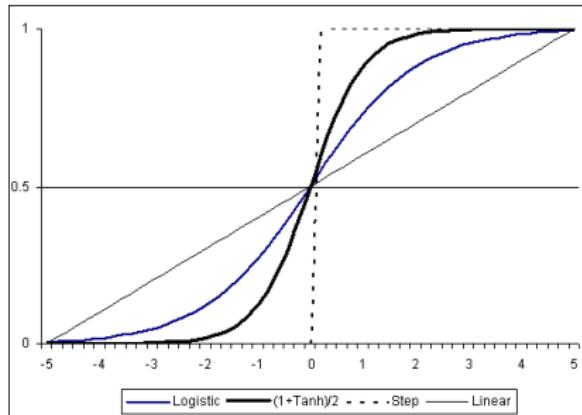
Multiple filters



Concepts: padding, stride, kernel size, number of filters/maps

Activation functions

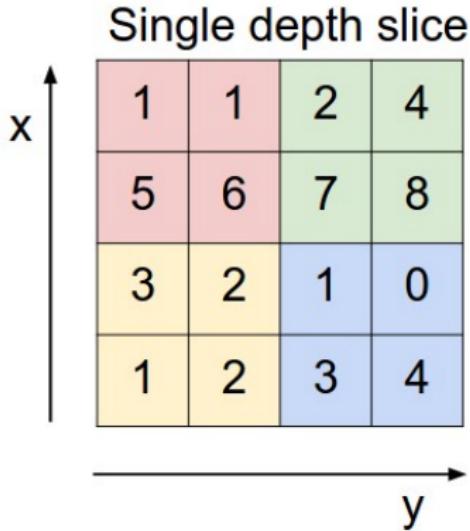
ReLU



$$F(x) = \max(0, x)$$

- via strides – skip pixels by steps of size d
- via pooling – aggregate $d \times d$ pixels

Max pooling

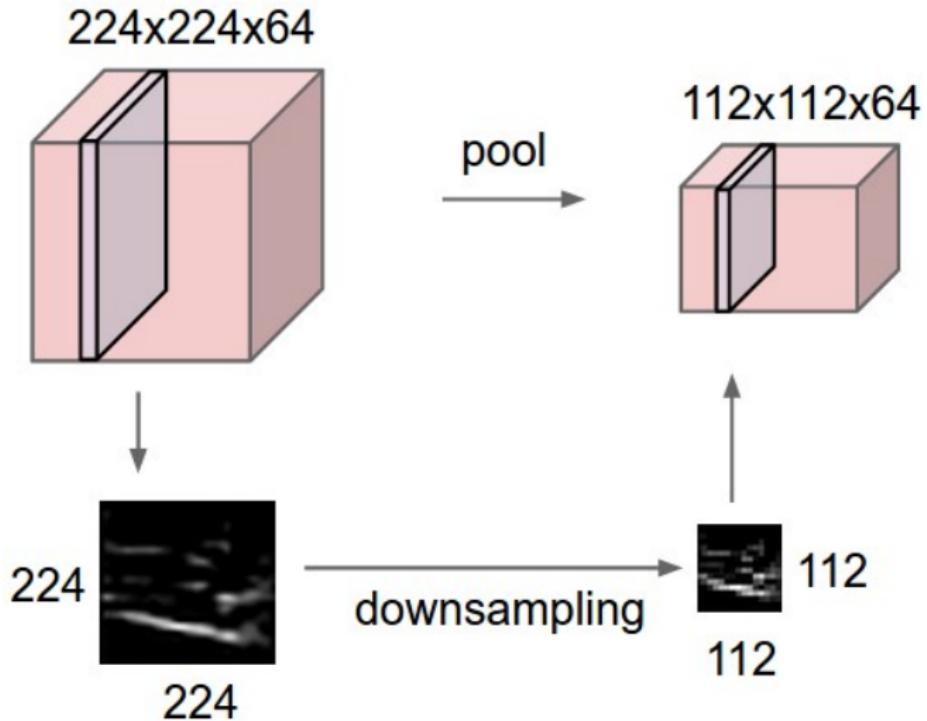


max pool with 2x2 filters
and stride 2

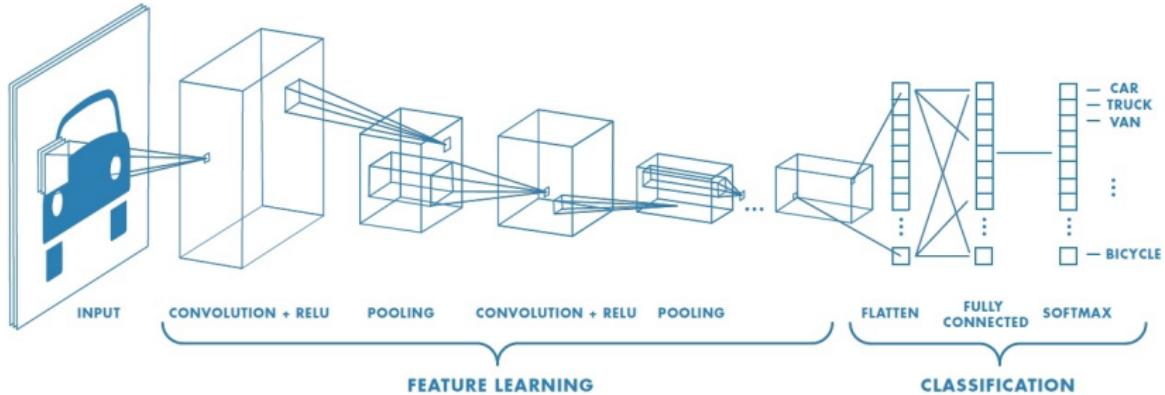
A 2x2 output matrix resulting from max pooling with 2x2 filters and stride 2. The matrix has two columns labeled x and two rows labeled y. The values are:

6	8
3	4

Pooling



Typical architecture of Convnets



Source: https://iitmcvg.github.io/summer_school/DLSession3/

- Convolutional layers: feature extraction
- Fully connected layers: classification

Training of CNN

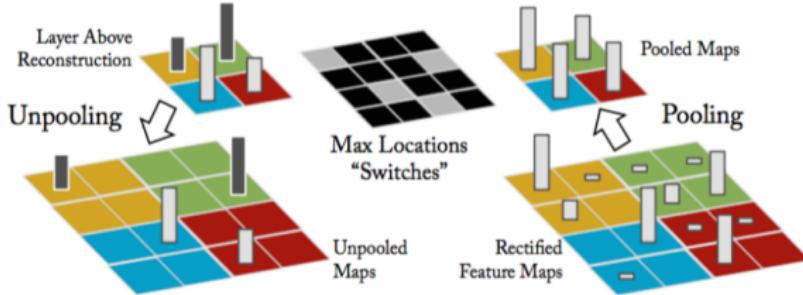
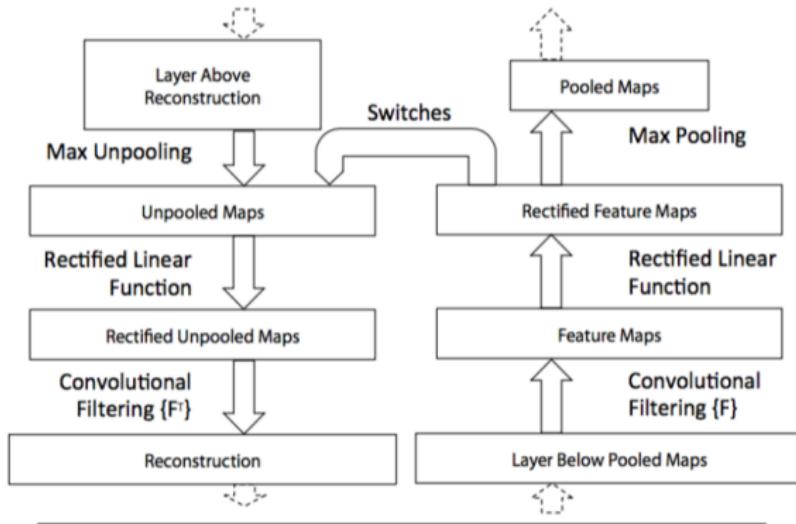
It is done using the backpropagation algorithm
(gradient descent)

Convolutional layers act as feature extractors. How ?

Zeiler et al., 2013, proposed a visualizaton technique

At a given convolutional layer, take the node with largest activation

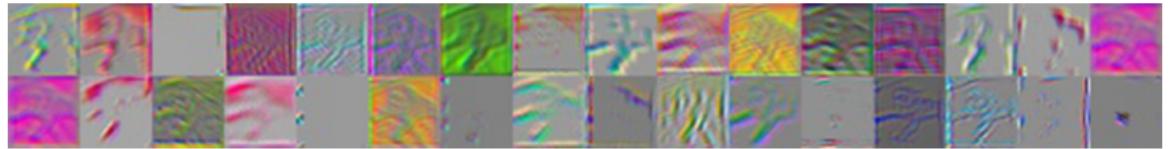
Reconstruct the input image by backpropagating from that node



(Zeiler et al., 2013)



Images reconstructed from neurons 1 to 32 of layer 1, for one input image

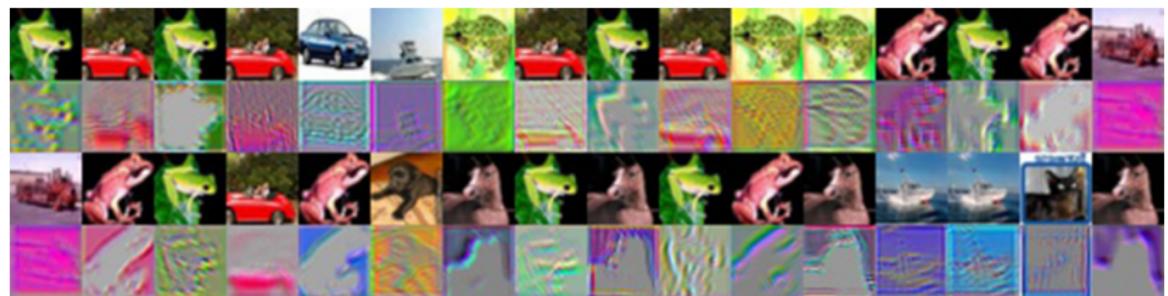


(<http://kvfrans.com/visualizing-features-from-a-convolutional-neural-network/>)

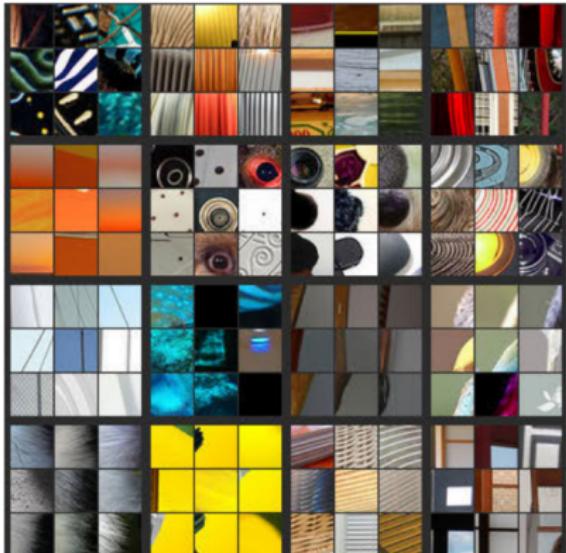
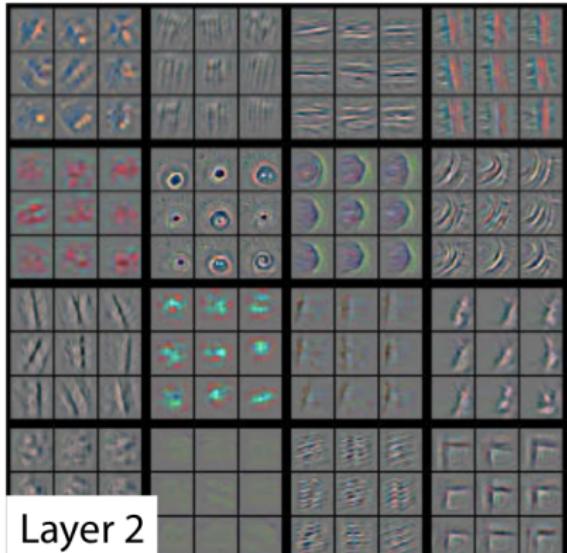
Images reconstructed from neuron 7 of layer 1, for multiple input images



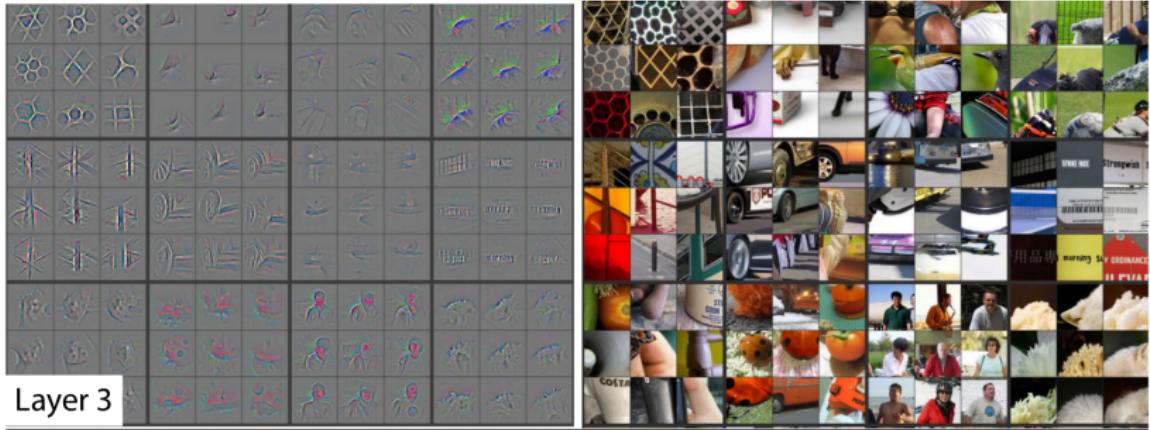
For each of the 32 neurons, reconstruction of the images that most activated the neuron



(Zeiler et al.) 16 neurons in layer 2; for each neuron, reconstruction of the 9 images that correspond to strongest activation of the node



(Zeiler et al.) 12 neurons in layer 3; for each neuron, reconstruction of the 9 images that correspond to strongest activation of the node



Large neural networks require large amount of training data

Transfer learning

It works as a nice weight initialization

Transfer learning refers to using whatever has been learned in a certain domain to tackle similar problems in another domain

Why? Because training data is a scarce resource

Common practice: train a network with images of ImageNet, and then fine-tune the weights with data of the target domain

Example: galaxy image classification

Dataset used: EFIGI (<https://www.astromatic.net/projects/efigi>)

1. Only spiral and elliptic
2. Extract *features* for each image
 - 2.1 Using Morfometryka hand-engineered features (Fabrício Ferrari, FURG)
 - 2.2 Using a pre-trained CNN
3. Separate data into training and validation sets
4. Logistic regression classifier trained on training data
5. Evaluation on the validation data

In the next slides:

- label above the image is the true class
(blue="elliptic", orange="spiral")
- first column: examples of wrongly classified images
- columns 2 to 5: most similar images to the ones at column 1,
according to feature similarity

Features morfométricos

Classificação errada



Features morfométricos

Classificação errada

4 vizinhos mais próximos



Features convolucionais

Classificação errada



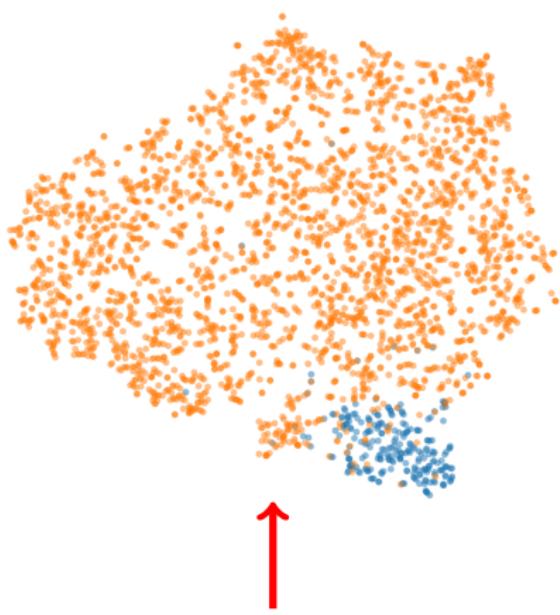
Astronomia ao meio-dia (N.S.T. Hirata)

Features convolucionais

Classificação errada



Projeções das *features*



Features morfométricos



The previous example did not use fine tuning

Main CNN architectures used in ILSVRC

- **AlexNet**, 2012: winner of the ImageNet Large Scale Visual Recognition Challenge (ILSVRC), 60M network parameters

(Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. NIPS 2012, pages 1097–1105)

- **VGG-11, 16 e 19**, 2014: 8, 13 e 16 convolutional layers, VGG-19 138M network parameters

(Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for largescale image recognition. CoRR, abs/1409.1556, 2014)

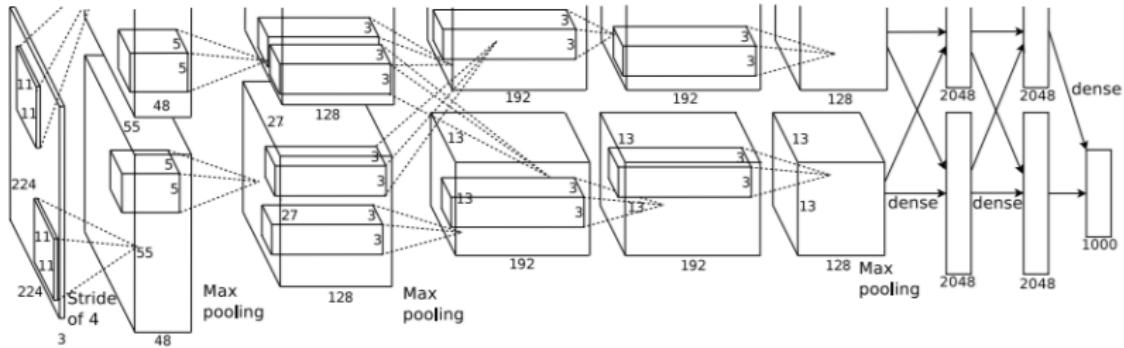
- **GoogleLeNet (Inception)**, 2014: winner of ILSVRC 2014, inception layers, 7M network parameters

(C. Szegedy, Wei Liu, Yangqing Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. CVPR 2015, pages 1–9)

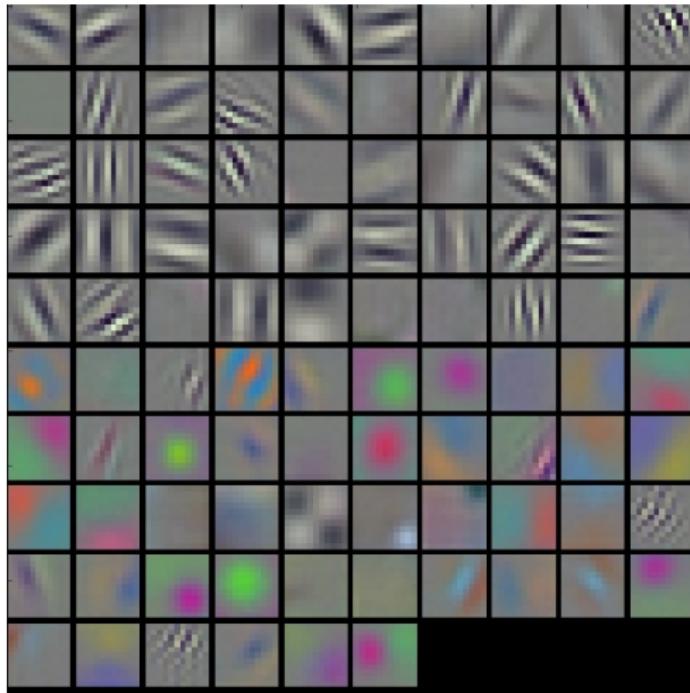
- **Residual Network (ResNet)**, 2015: winner of ILSVRC 2015, 25.5M network parameters, residual block, vanishing gradient

(K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. CVPR 2016, pages 770–778)

AlexNet (winner ILSVRC 2012)

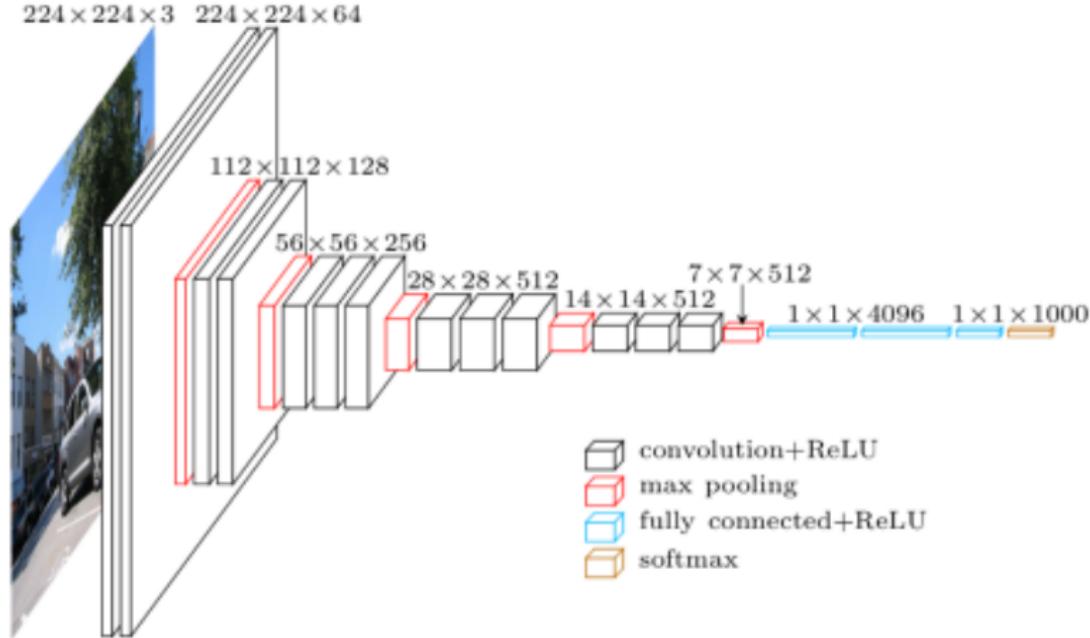


- succeeded on training a large CNN
- 60 million parameters
- used 2 GPUs
- proposed ReLU
- used dropout
- used data augmentation (rotation, scale, crops, etc)
- indication that number of layers is important



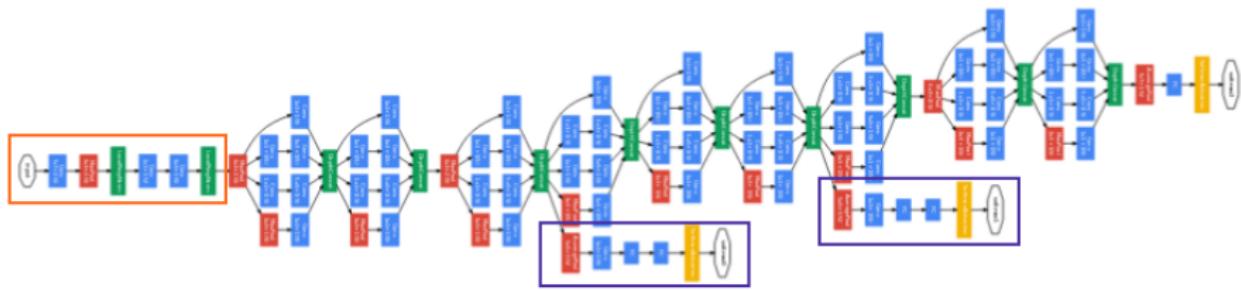
Images that maximize the response of the filters (AlexNet)

VGGNet (2nd place ILSVRC 2014)

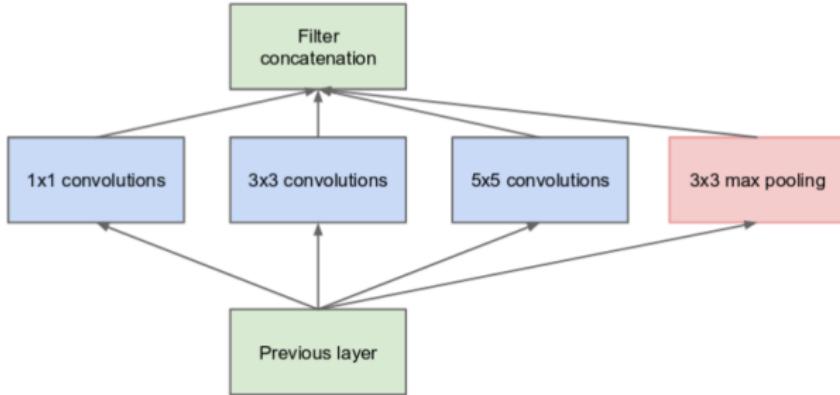


- multiple layers (ex., 16 e 19)
- only 3x3 convolutions
- stride 1 for the convolutions (AlexNet used stride ≥ 1)
- first layers with few filters and gradually increasing as we advance through the layers
- 144 millions of parameters

Inception - GoogleLeNet (winner ILSVRC 2014)

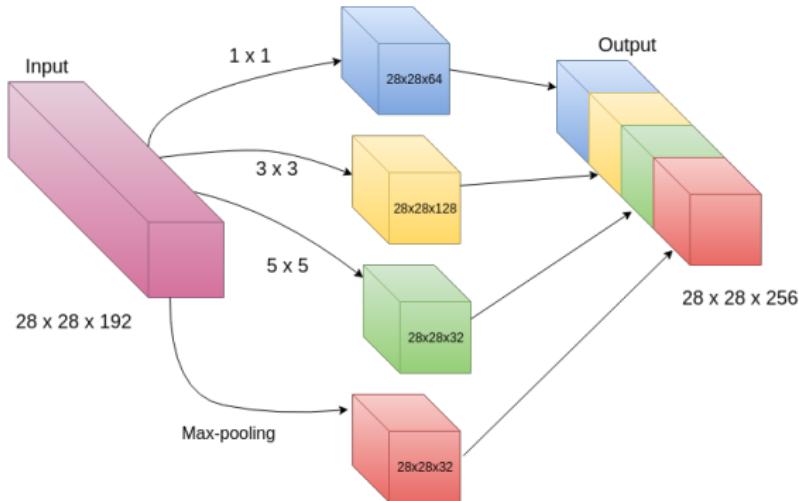


Inception – naïve version



(a) Inception module, naïve version

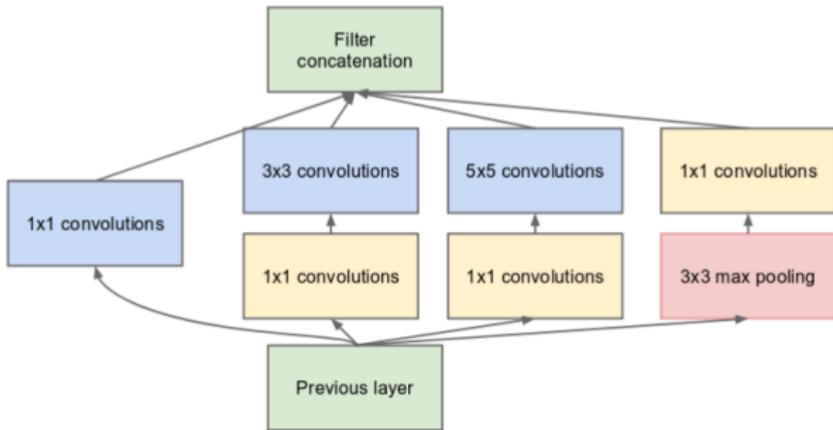
Inception – naïve version example



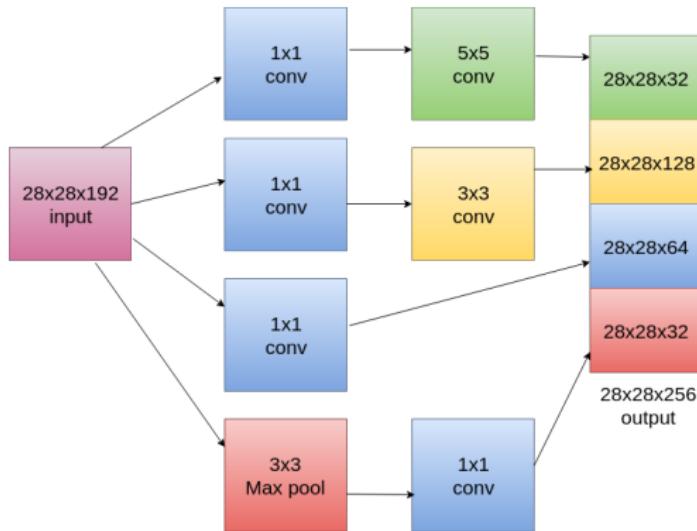
Camada de entrada = $28 \times 28 \times 192$

Convolução 5×5 : com 32 filtros, camada de saída = $28 \times 28 \times 32$. Para cada ponto da camada de saída, o número de multiplicações necessárias é $5 \times 5 \times 192$. Como são $28 \times 28 \times 32$ pontos na camada de saída, o total de multiplicações necessárias é $5 \times 5 \times 192 \times 28 \times 28 \times 32 = 120$ milhões.

Inception – effective version

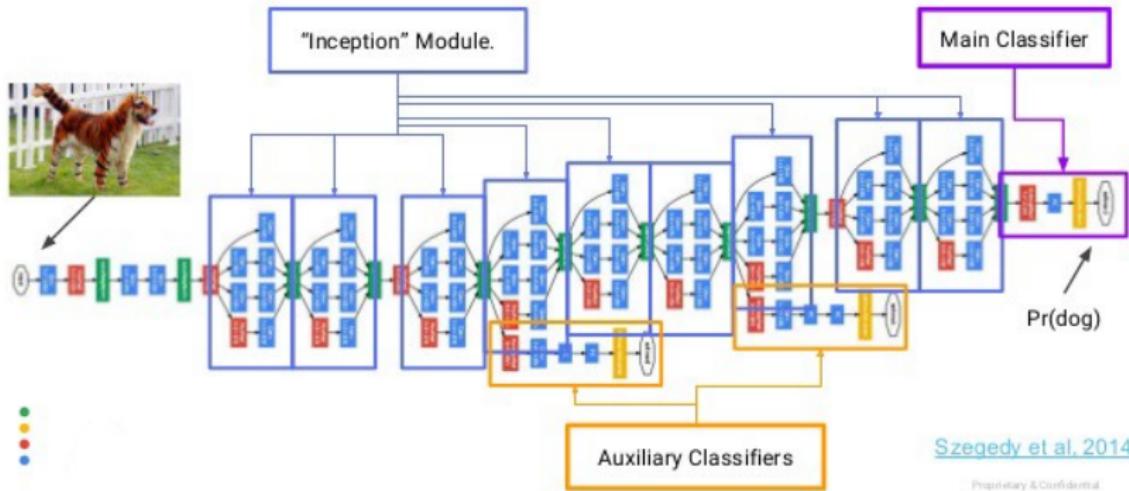


(b) Inception module with dimension reductions



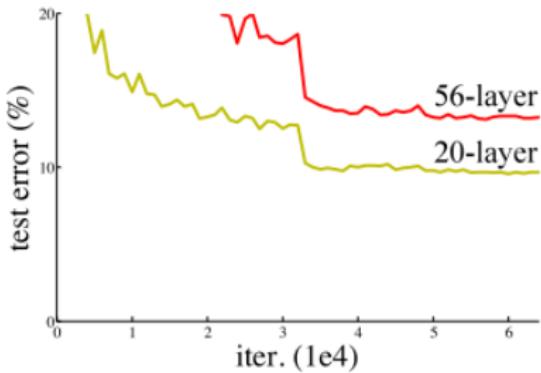
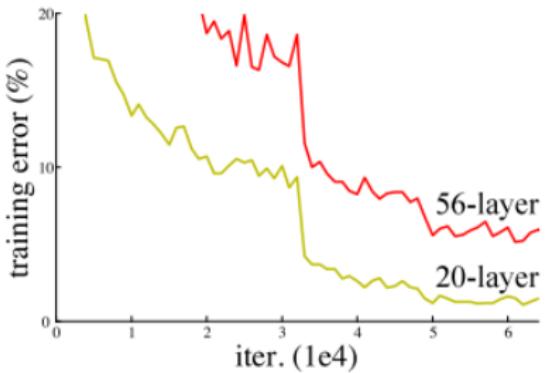
Entrada $28 \times 28 \times 192 \rightarrow 16$ filtros $1 \times 1 \rightarrow$ camada intermediária $28 \times 28 \times 16 \rightarrow$ 32 convoluções $5 \times 5 \rightarrow$ saída $28 \times 28 \times 32$. Total de multiplicações: cada ponto na camada intermediária requer $1 \times 1 \times 192$ multiplicações, e portanto $1 \times 1 \times 192 \times 28 \times 28 \times 16 = 2.4$ milhões para cálculo da camada intermediária; cada ponto na camada de saída requer $5 \times 5 \times 16$ multiplicações e portanto $5 \times 5 \times 16 \times 28 \times 28 \times 32 = 10$ milhões para o cálculo da camada de saída. Total = 12.4 milhões.

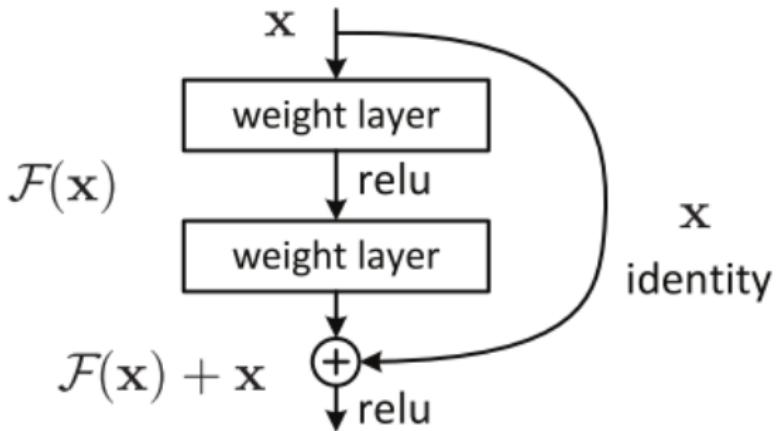
GoogLeNet (aka “Inception”) Architecture



- concept of “network in network”
- inception module (“choices” made during training)
- use of 1x1 convolution
- auxiliary output to reinforce activation of intermediary layers
- more aggressive data augmentation than AlexNet
- 4 million parameters (?)
- variants (e.g., change 5x5 with two 3x3, or change 3x3 with a 1x3 and a 3x1)

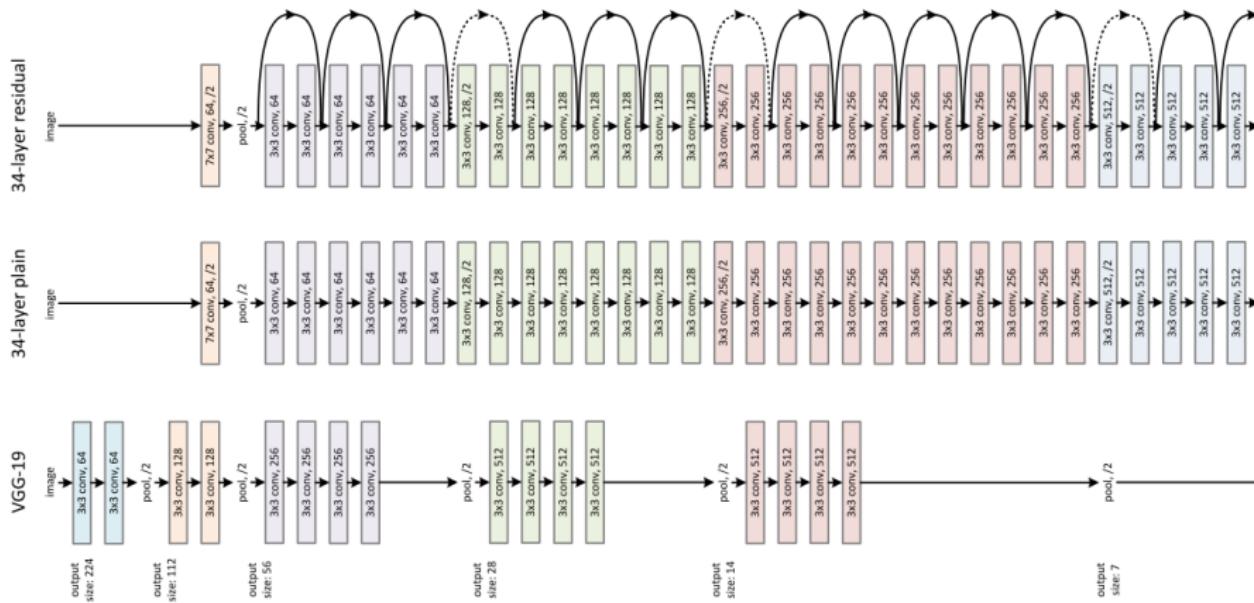
Apesar de ser aceito que profundidade da rede é importante, observou-se que camadas adicionais degradavam a acurácia no conjunto de treinamento





Entrada x . Pode ser conveniente aprender um mapeamento ($H(x)$) que seja a identidade. Mas representar identidade por meio de transformações não lineares não é simples. Assim, em vez de H , aprende-se o resíduo $\mathcal{F}(x)$, de modo que $H(x) = \mathcal{F}(x) + x$ e assim a identidade corresponde a resíduo zero

ResNet (winner ILSVRC 2015)



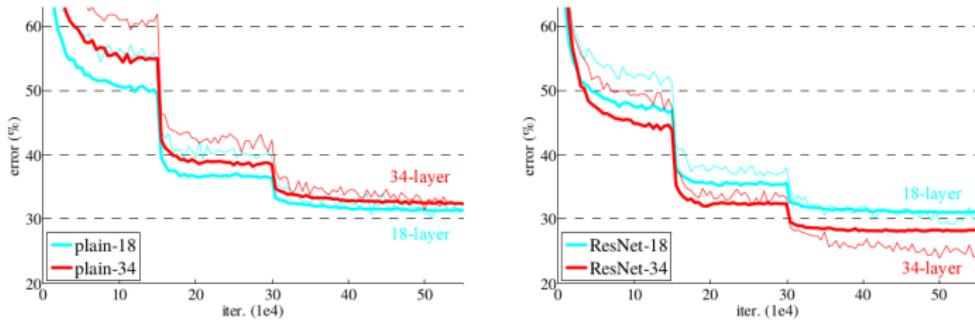


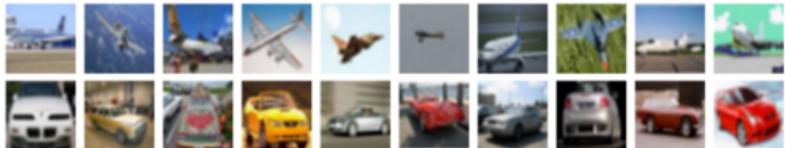
Figure 4. Training on **ImageNet**. Thin curves denote training error, and bold curves denote validation error of the center crops. Left: plain networks of 18 and 34 layers. Right: ResNets of 18 and 34 layers. In this plot, the residual networks have no extra parameter compared to their plain counterparts.

- identificam uma degradação associada ao aumento de camadas, que aparentemente não é *overfitting* nem *vanishing gradient*
- propõe módulo para aprender o resíduo da transformação desejada
- aplica *batch normalization* após cada convolução e antes da ativação
- não usa *dropout*
- conseguem treinar redes com mais de 100 camadas
- testado em outros datasets (CIFAR, PASCAL, MS-COCO)
- ResNet 56 layers: 0.85M parameters (?)
- ResNet 110 layers: 1.7M parameters (?)

DL is being used to address multiple tasks in CV

Image classification

airplane



automobile



bird



cat



deer



dog



frog



horse



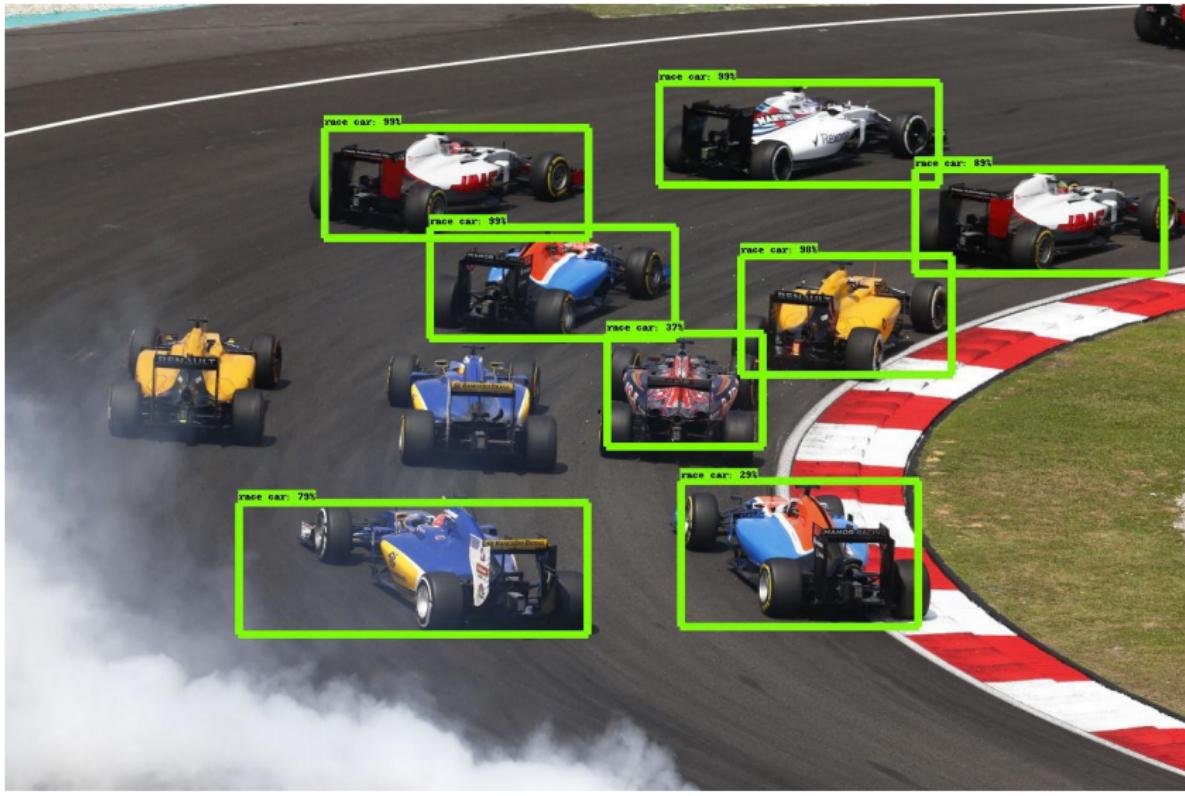
ship



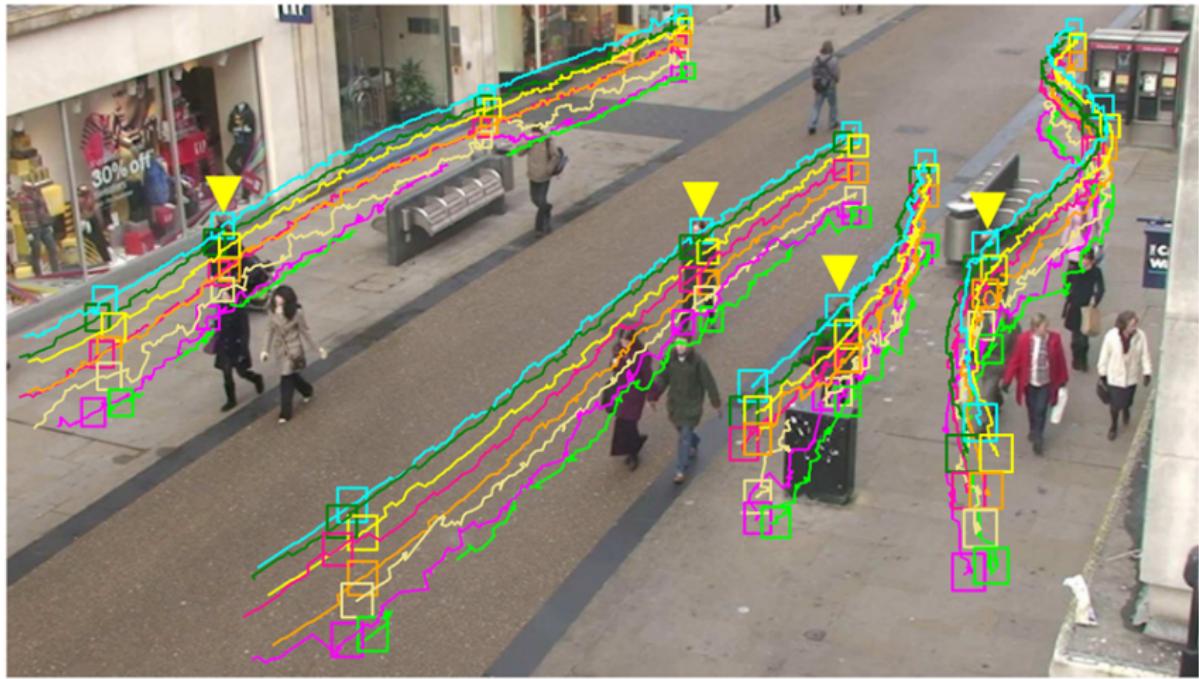
truck



Object detection



Object tracking in video



2

4

7

Semantic segmentation



30 FPS

Scene description



A female tennis player in action on the court.



A group of young men playing a game of soccer



A man riding a wave on top of a surfboard.



A baseball game in progress with the batter up to plate.



A brown bear standing on top of a lush green field.



A person holding a cell phone in their hand.

CNN for image segmentation

More: <http://blog.cure.ai/notes/semantic-segmentation-deep-learning-review>

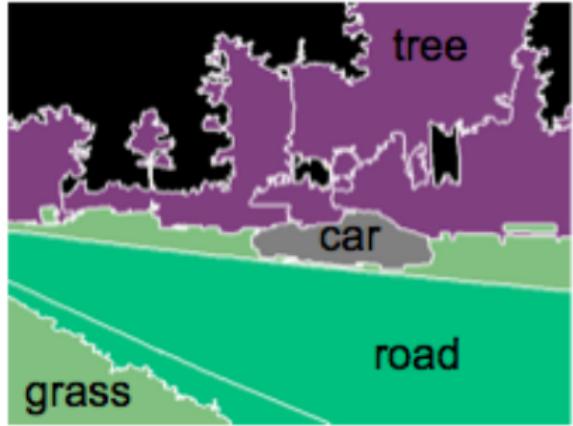


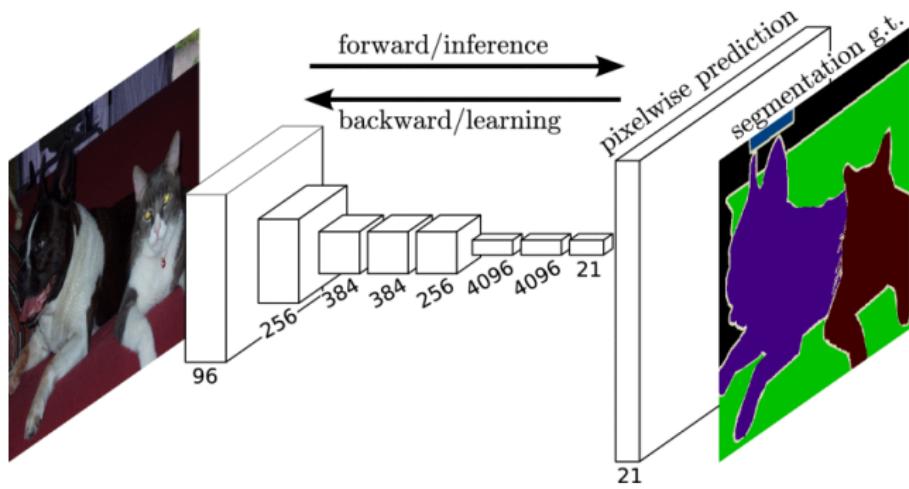
Image segmentation

Objetivo: Delinear o contorno dos objetos/componentes de interesse presentes na imagem (ou atribuir um rótulo único para todos os *pixels* de um objeto)

A arquitetura das redes segue a forma *encoder-decoder*

A parte *encoder* em geral é a parte convolucional de uma CNN.

A parte *decoder* busca restaurar a resolução inicial.



U-Net

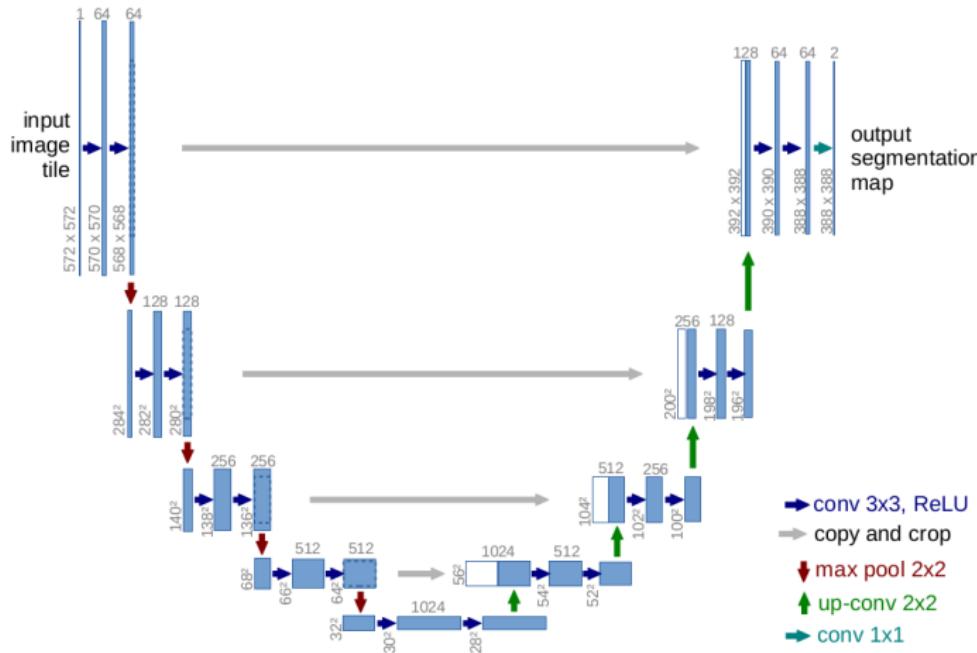
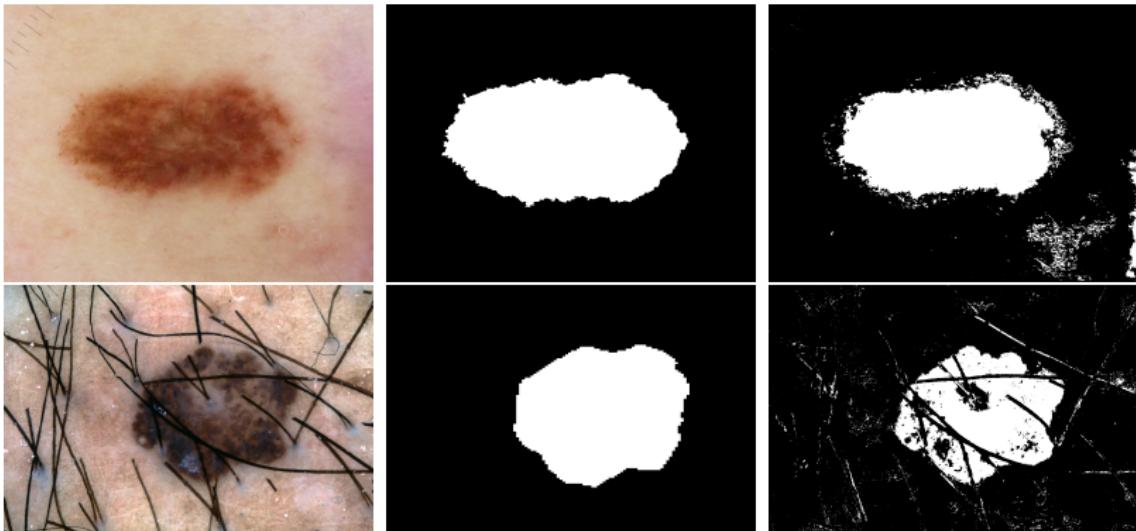


Fig. 1. U-net architecture (example for 32x32 pixels in the lowest resolution). Each blue box corresponds to a multi-channel feature map. The number of channels is denoted on top of the box. The x-y-size is provided at the lower left edge of the box. White boxes represent copied feature maps. The arrows denote the different operations.

Application example



(ISIC 2018 dataset)

CNN for object detection/recognition

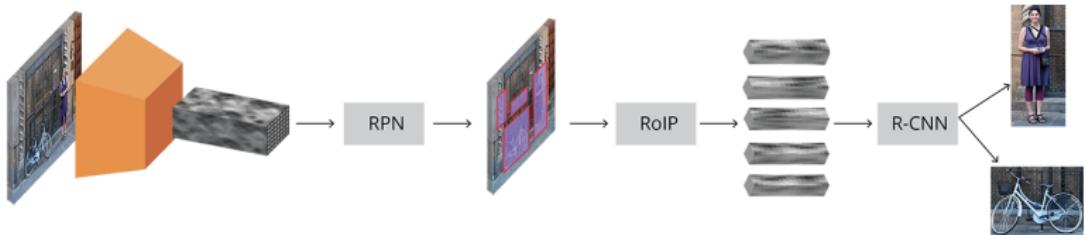
<https://towardsdatascience.com/deep-learning-for-object-detection-a-comprehensive-review-73930816d8d9>

O conjunto de dados de treinamento consiste de imagens mais uma lista de *bounding box* dos objetos com os respectivos rótulos de classe.

Função objetivo é uma combinação de classificação (acertar rótulos) e regressão (acertar tamanho e localização dos *boxes*).

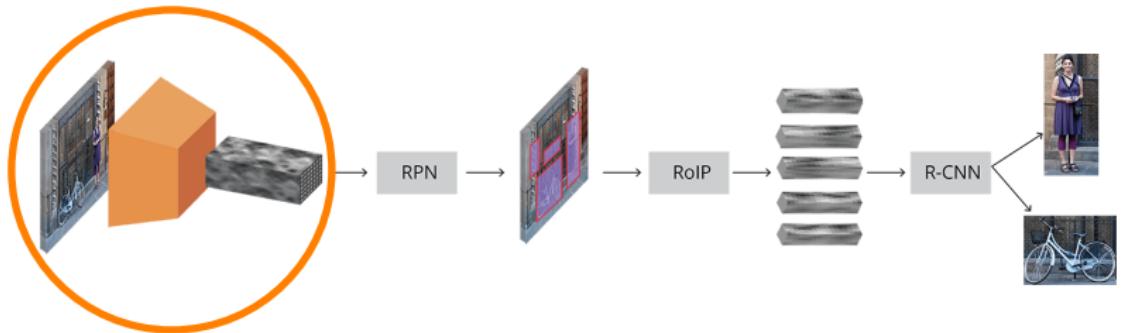
Existem os modelos two-step (ex., Faster R-CNN) e *single shot* (SSD ou YOLO)

Faster R-CNN architecture



<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>

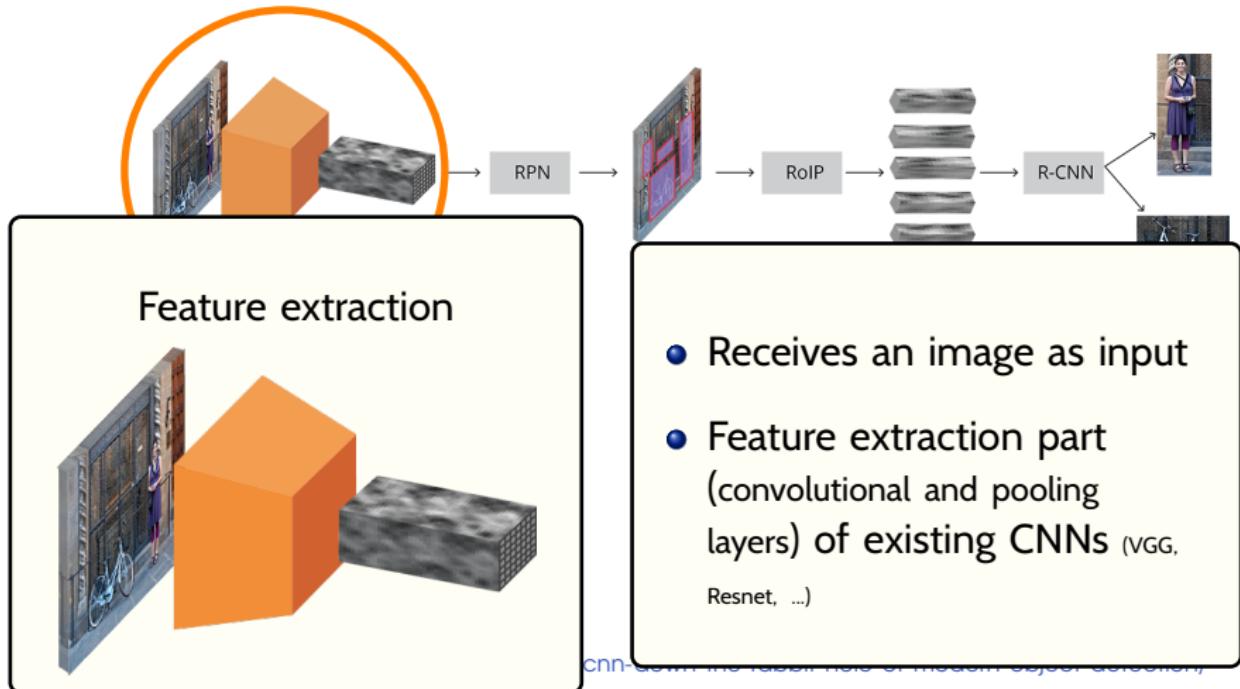
Faster R-CNN architecture



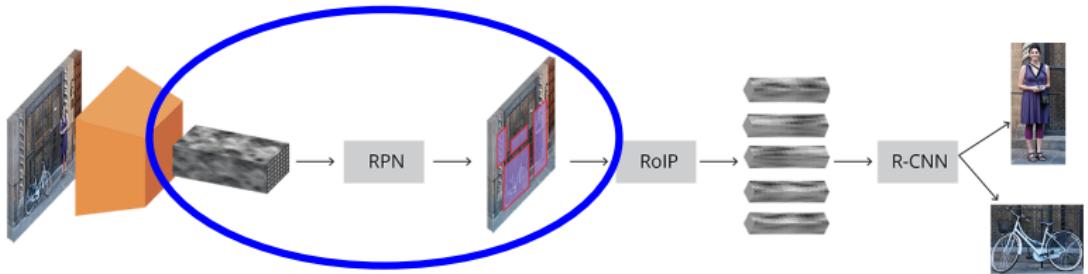
Feature extraction

<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>

Faster R-CNN architecture



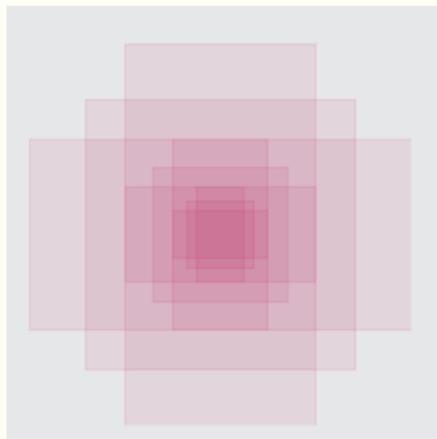
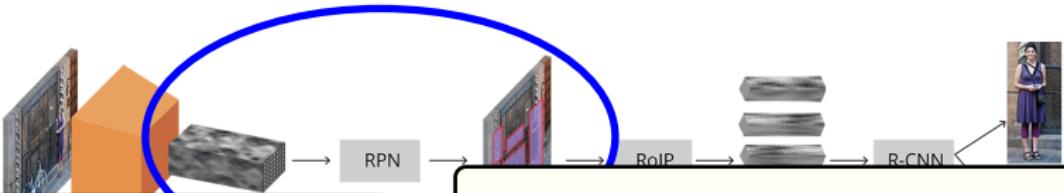
Faster R-CNN architecture



Region proposal network

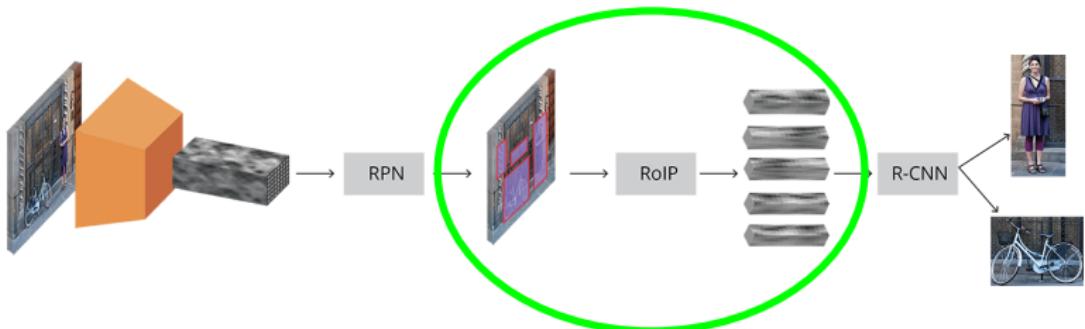
<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>

Faster R-CNN architecture



- Learns if there is an object in the image, in the region enclosed by each of the anchors
- That is done simultaneously over the whole image
- Takes care of the overlaps and returns up to N that are most likely to contain an object

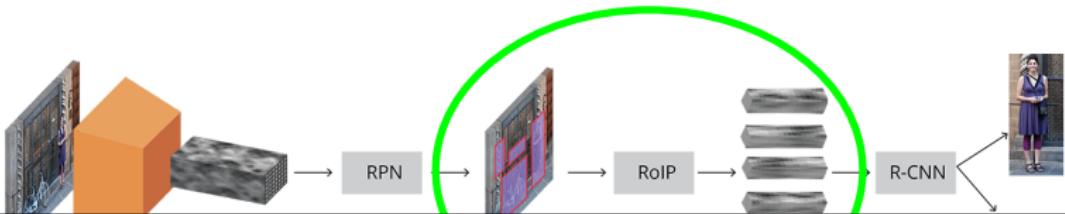
Faster R-CNN architecture



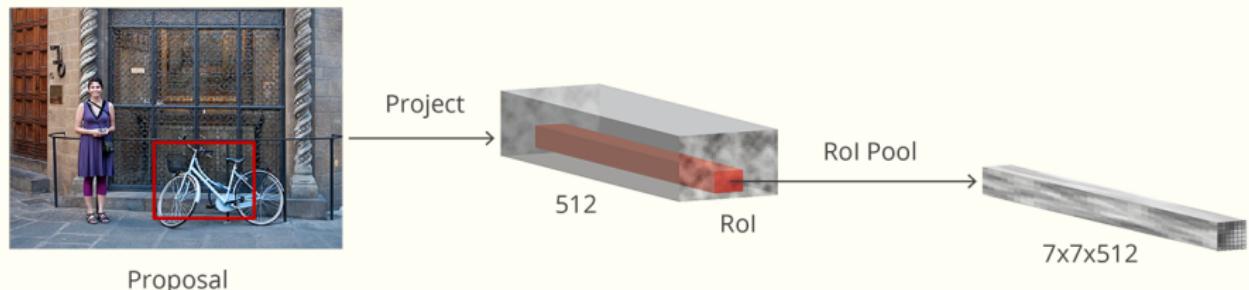
Region of interest pooling

<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>

Faster R-CNN architecture

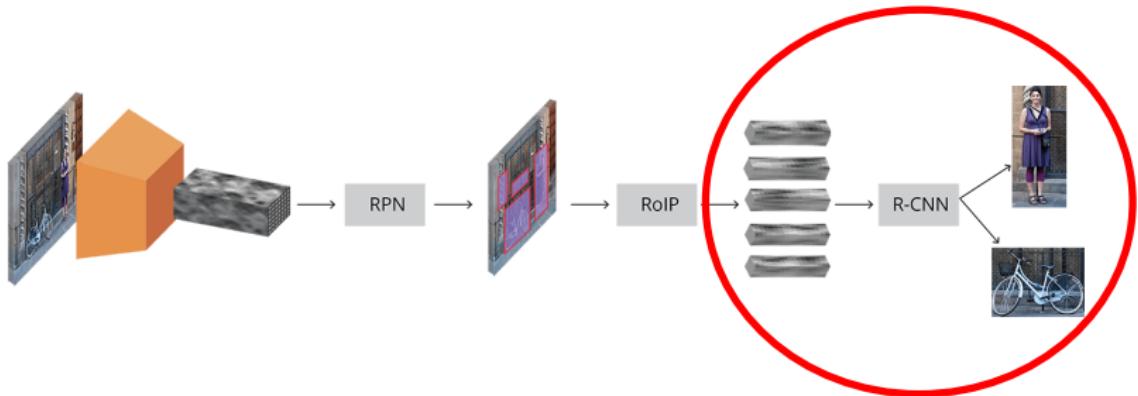


Region of interest pooling



- For each region proposal, crops the corresponding cells from the feature map
 - pools it to a fixed size

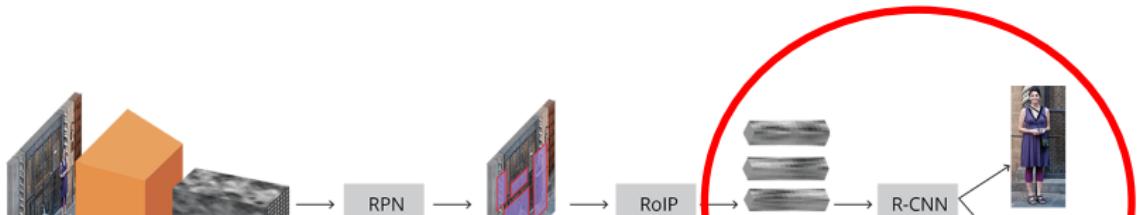
Faster R-CNN architecture



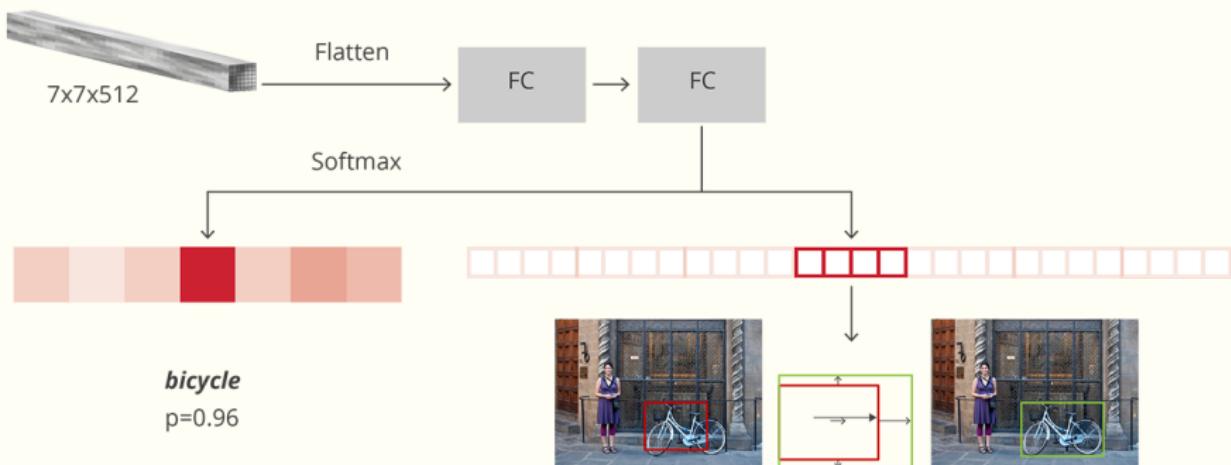
Region classification

<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>

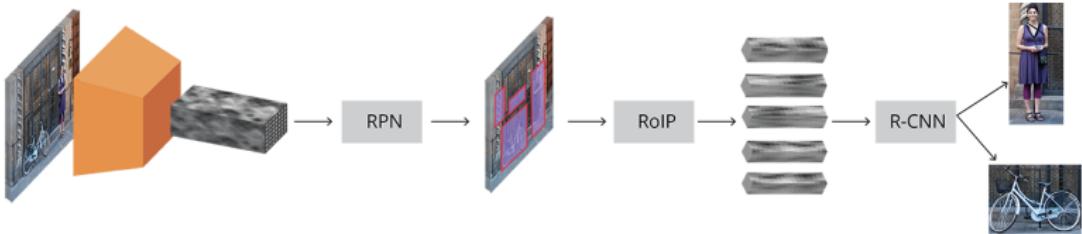
Faster R-CNN architecture



Region classification



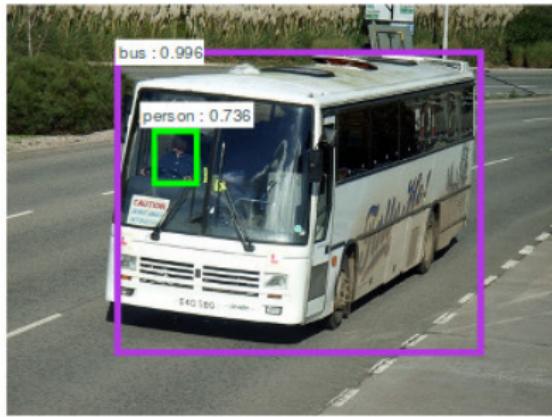
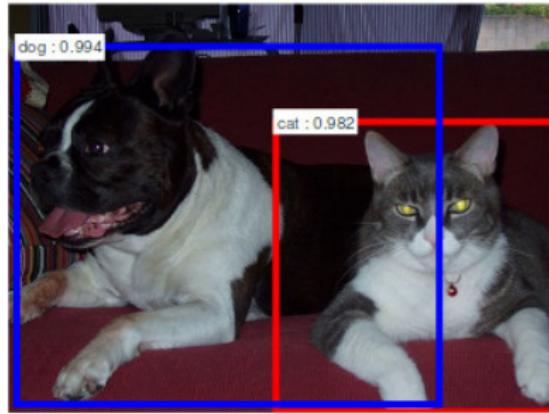
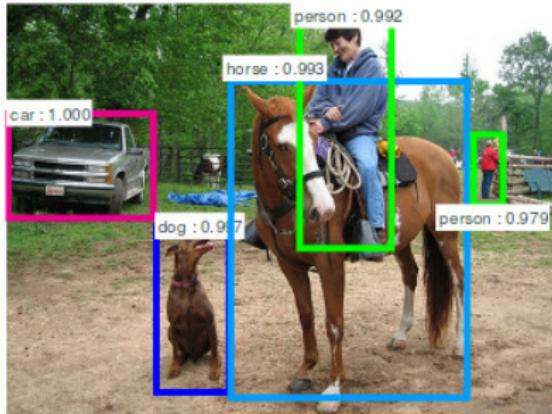
Faster R-CNN architecture



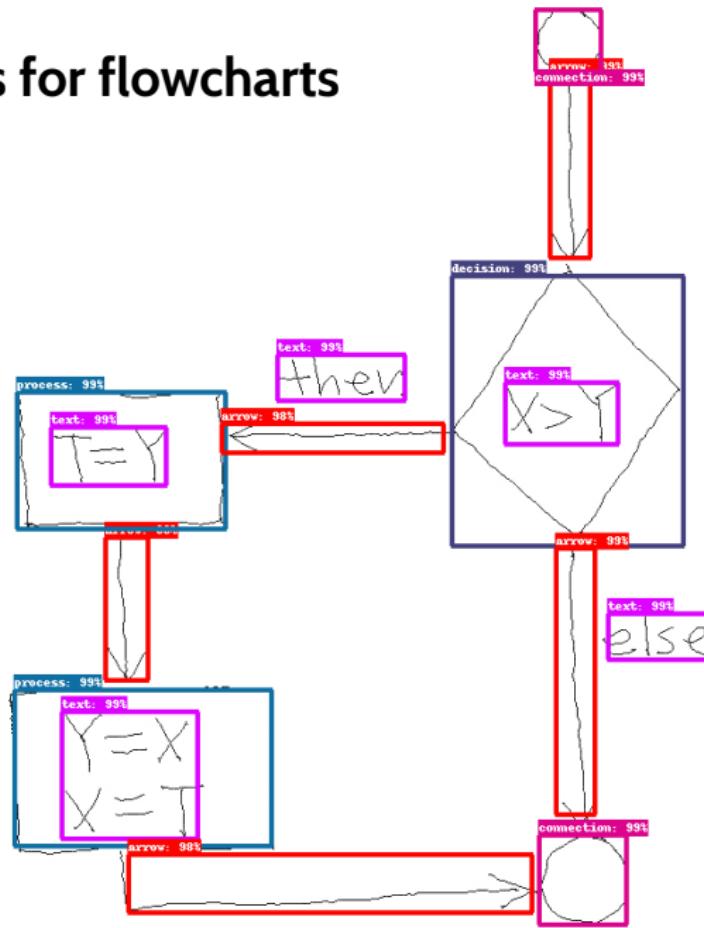
Input is an image

Output is a list of detected objects (bounding box and class probabilities)

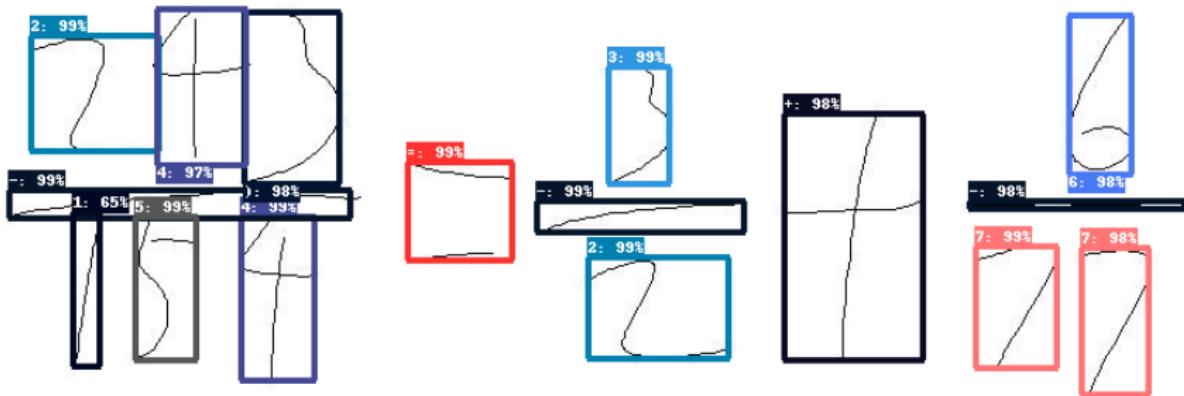
<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>



Results for flowcharts



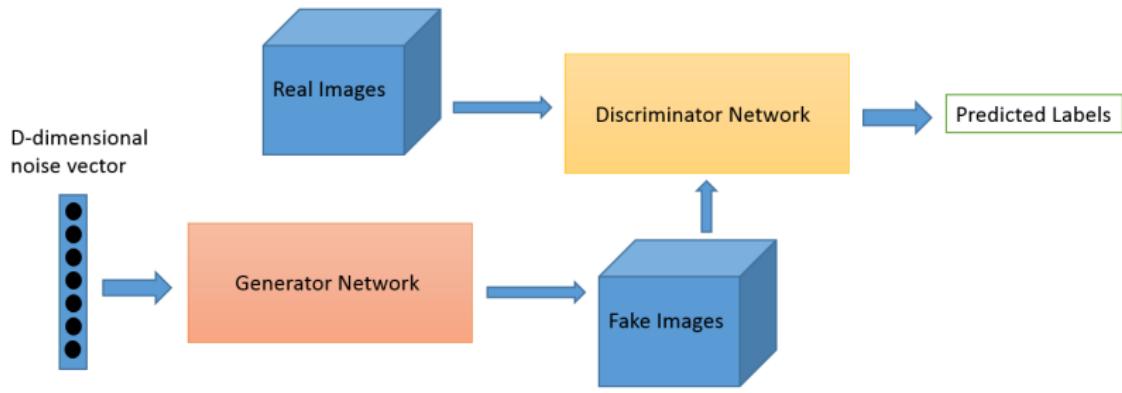
Results for math expressions



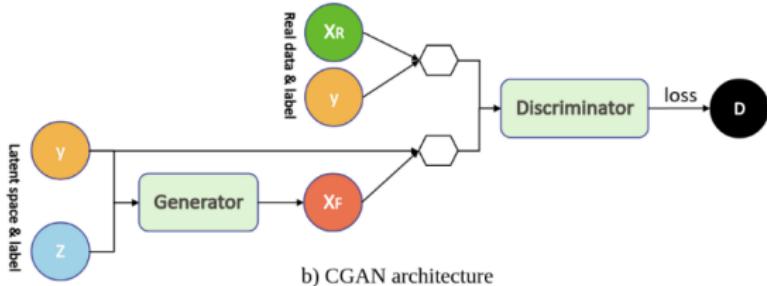
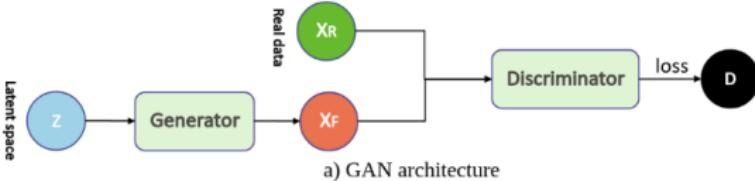
Other DL network models

- **GAN** – image generation / adversarial training
- **Siamese network, triplet network** — image retrieval
- **RNN, LSTM, GRU** – temporal/sequential data processing
- **Geometric Deep Learning** deals with the extension of Deep Learning techniques to graph/manifold structured data.
<http://geometricdeeplearning.com/>

GAN schema



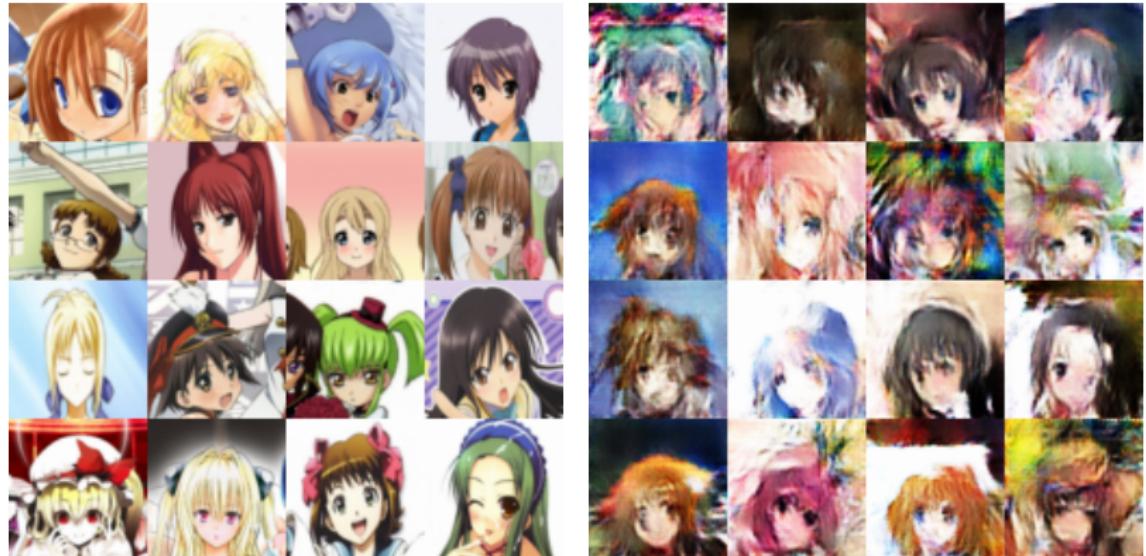
(fonte: O'Reilly)



<https://mc.ai/a-tutorial-on-conditional-generative-adversarial-nets-keras-implementation/>

$$\min_G \max_D V(D, G) = E_{\mathbf{x} \sim p_{data}(\mathbf{x})}[\log D(\mathbf{x})] + E_{\mathbf{z} \sim p_z(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))]$$

$$\min_G \max_D V(D, G) = E_{\mathbf{x} \sim p_{data}(\mathbf{x})}[\log D(\mathbf{x}|\mathbf{y})] + E_{\mathbf{z} \sim p_z(\mathbf{z})}[\log(1 - D(G(\mathbf{z}|\mathbf{y})))]$$



(<https://github.com/nashory/gans-awesome-applications>)



es rostos ?





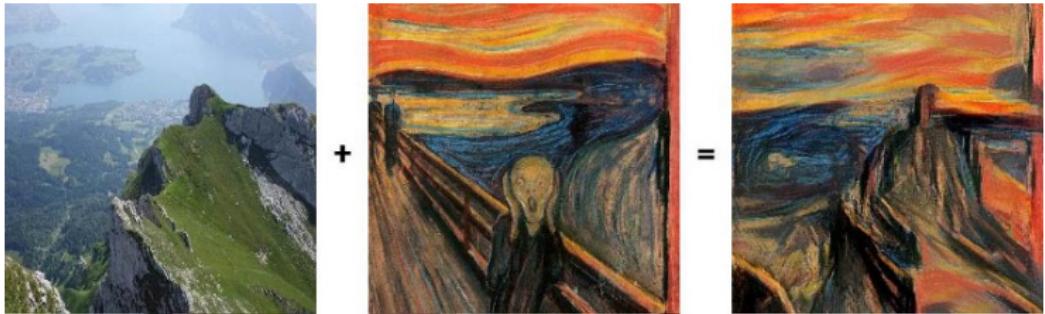




Esses rostos não existem. Foram gerados por computador!

<https://thispersondoesnotexist.com/>

Artistic Style transfer

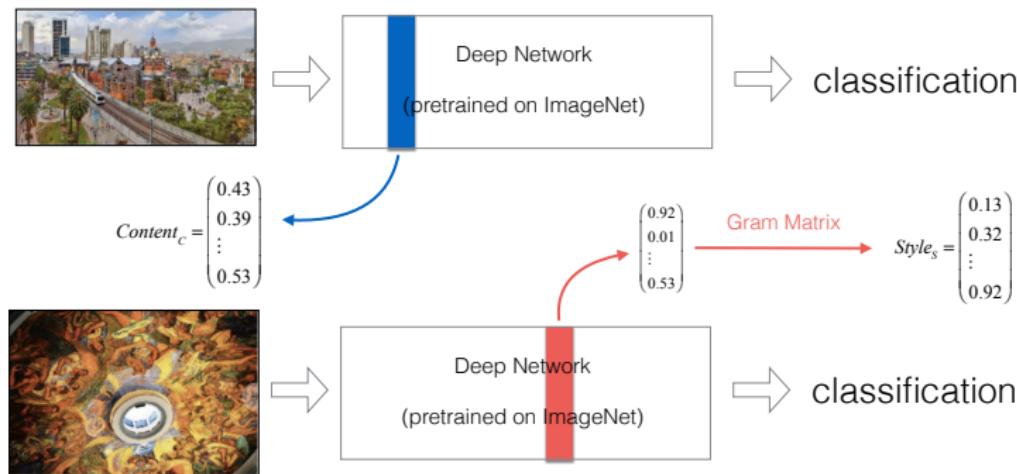


Artistic Style Transfer with Convolutional Neural Network

Case study 3: Art Generation

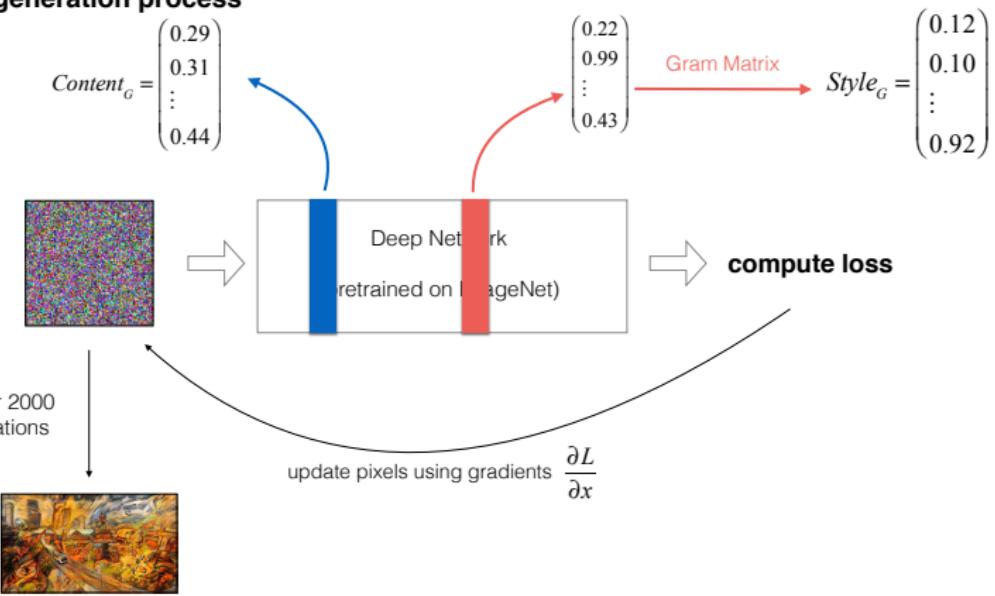
4. Architecture?

We use a **pre-trained model** because it **extracts important information** from images.



Case study 3: Art Generation

Image generation process



$$Content_c = \begin{pmatrix} 0.43 \\ 0.39 \\ \vdots \\ 0.53 \end{pmatrix}$$
$$Style_s = \begin{pmatrix} 0.13 \\ 0.32 \\ \vdots \\ 0.92 \end{pmatrix}$$

Kian Katanforoosh

Some references

Neural nets: <http://neuralnetworksanddeeplearning.com/>

NN and CNN: Stanford CS class <http://cs231n.github.io/convolutional-networks/>

Github (F. Chollet): <https://github.com/fchollet/deep-learning-with-python-notebooks>

Books

- Deep Learning, Ian Goodfellow and Yoshua Bengio and Aaron Courville (<https://www.deeplearningbook.org/>)
- Dive into Deep Learning <https://d2l.ai/>