

Deep Learning with Neural Networks

Some important CNN variants

Pablo Martínez Olmos, pamartin@ing.uc3m.es

Style transfer



This CVPR paper is the Open Access version, provided by the Computer Vision Foundation.
Except for this watermark, it is identical to the version available on IEEE Xplore.

Image Style Transfer Using Convolutional Neural Networks

Leon A. Gatys

Centre for Integrative Neuroscience, University of Tübingen, Germany

Bernstein Center for Computational Neuroscience, Tübingen, Germany

Graduate School of Neural Information Processing, University of Tübingen, Germany

leon.gatys@bethgelab.org

Alexander S. Ecker

Centre for Integrative Neuroscience, University of Tübingen, Germany

Bernstein Center for Computational Neuroscience, Tübingen, Germany

Max Planck Institute for Biological Cybernetics, Tübingen, Germany

Baylor College of Medicine, Houston, TX, USA

Matthias Bethge

Centre for Integrative Neuroscience, University of Tübingen, Germany

Bernstein Center for Computational Neuroscience, Tübingen, Germany

Max Planck Institute for Biological Cybernetics, Tübingen, Germany

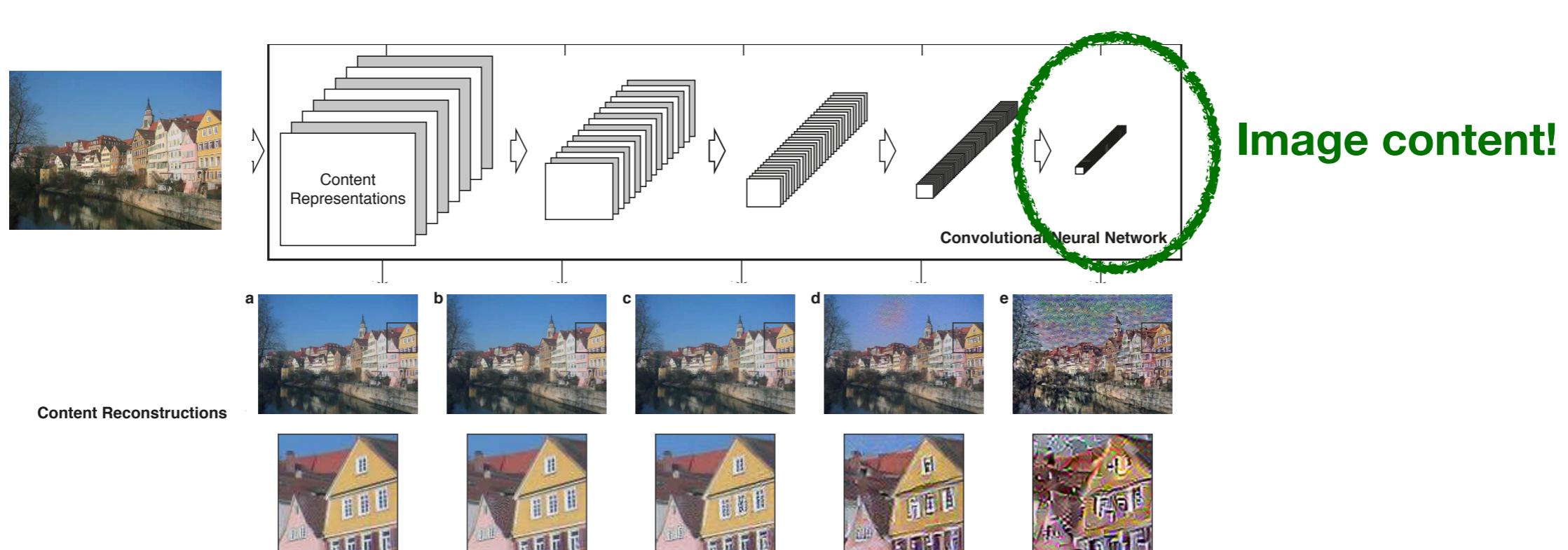
Style Transfer

- “We use image representations derived from CNNs optimised for object recognition, which make **high level image information explicit**.”
- “Our algorithm can **separate and recombine the image content and style** of images”
- “The algorithm allows us to produce new images of high perceptual quality that **combine the content** of an arbitrary photograph **with the appearance** of numerous well-known artworks”
- “Our results provide new insights into the deep image representations learned by CNNs and demonstrate their **potential for high level image synthesis and manipulation**”



Extracting content from an image

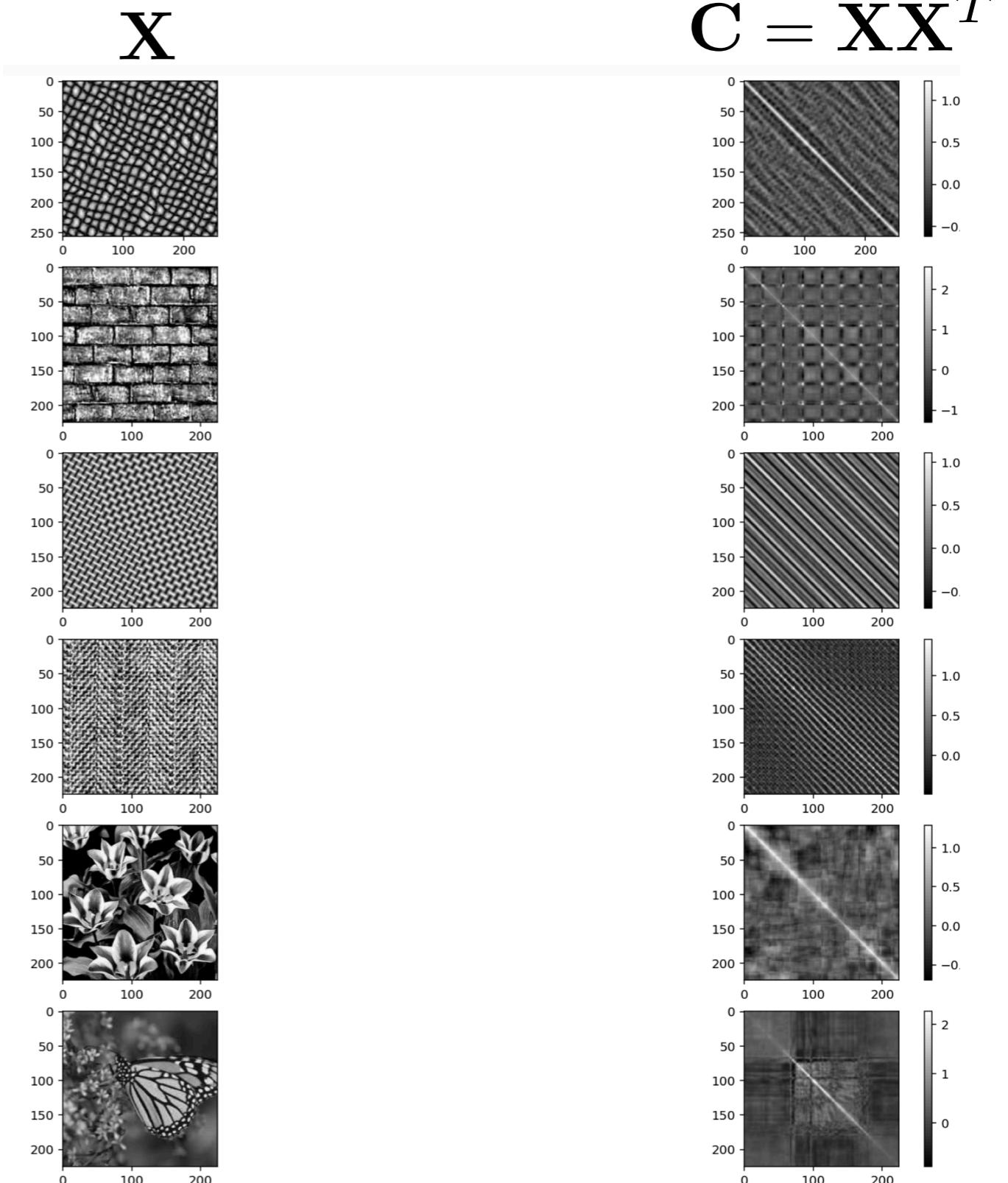
- The input image is transformed into **representations that are increasingly sensitive** to the actual content of the image, but become relatively invariant to its precise appearance.
- **Higher layers** in the network capture the high-level content in terms of objects, but do not constrain the exact pixel values of the reconstruction very much.
- Reconstructions from the **lower layers simply reproduce the exact pixel values of the original image**.



Textures as image correlations

- The **correlation matrix is sensitive to textures** (spatial patterns that happen all over the image).

- **Right:** normalised images and the correlation matrix.



Extracting style from an image

- Style feature space can be built on top of the filter responses in any layer of the network. It consists of the **correlations between the different filter responses**.
- Feature correlations are given by the Gram matrix G , where G_{ij}^l is the inner product between the vectorised featured maps i and j in layer l :

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l.$$

Extracting style from an image

- Style feature space can be built on top of the filter responses in any layer of the network. It consists of the **correlations between the different filter responses**.
- Feature correlations are given by the Gram matrix G , where G_{ij}^l is the inner product between the vectorised featured maps i and j in layer l :

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l.$$


Sum across all pixels in feature maps

Extracting style from an image

- Style feature space can be built on top of the filter responses in any layer of the network. It consists of the **correlations between the different filter responses**.
- Feature correlations are given by the Gram matrix G , where G_{ij}^l is the inner product between the vectorised featured maps i and j in layer l :

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l.$$


Sum across all pixels in feature maps

- We obtain a stationary, multi-scale representation of the input image, which **captures its texture information but not the global arrangement**.

Extracting style from an image

- Style feature space can be built on top of the filter responses in any layer of the network. It consists of the **correlations between the different filter responses**.

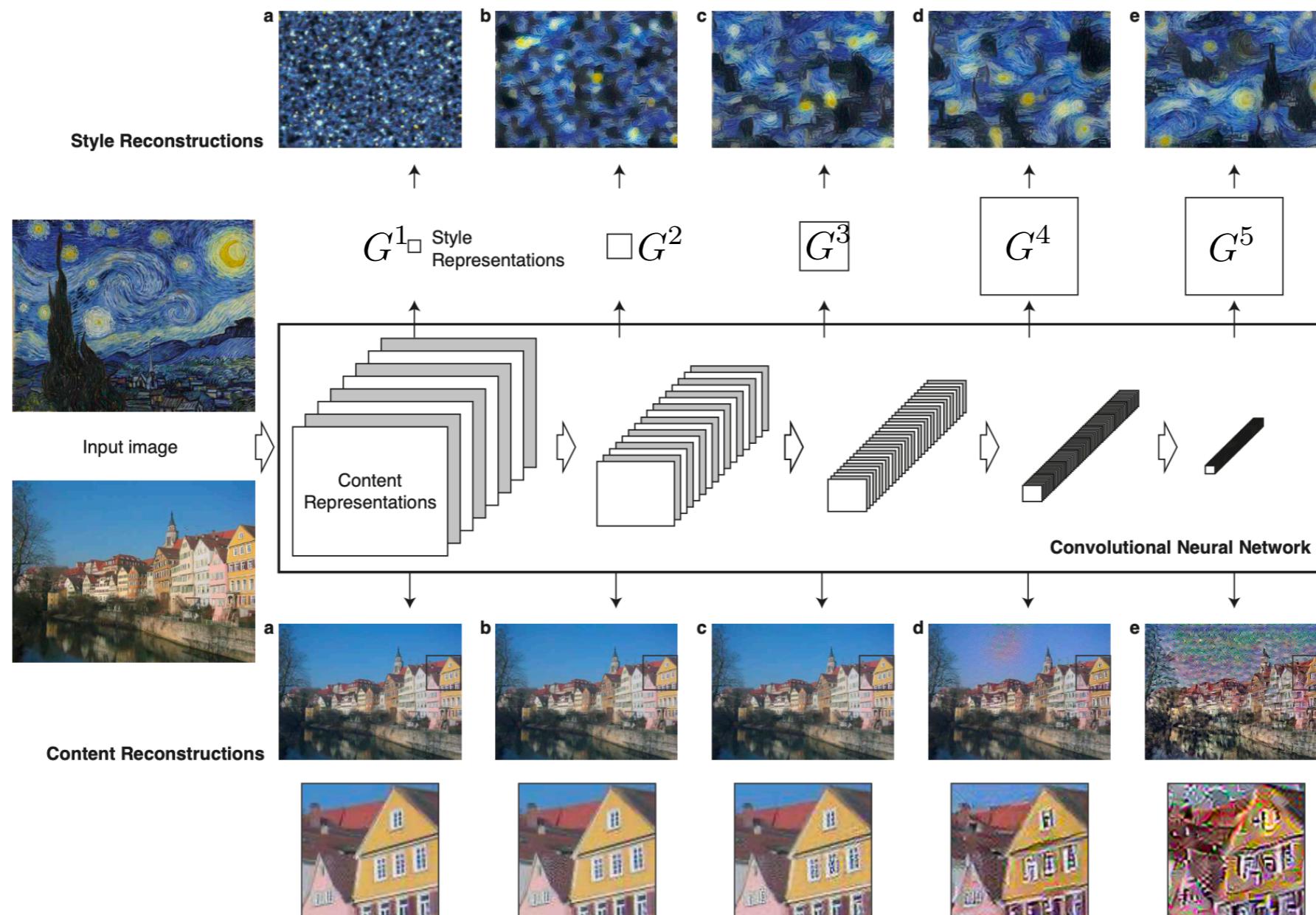
- Feature correlations are given by the Gram matrix G , where G_{ij}^l is the inner product between the vectorised featured maps i and j in layer l :

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l.$$

Sum across all pixels in feature maps

- We obtain a stationary, multi-scale representation of the input image, which **captures its texture information but not the global arrangement**.
- As with content, we can visualise the information captured by these style feature spaces built on different layers of the network by constructing an image that matches the style representation of a given input image

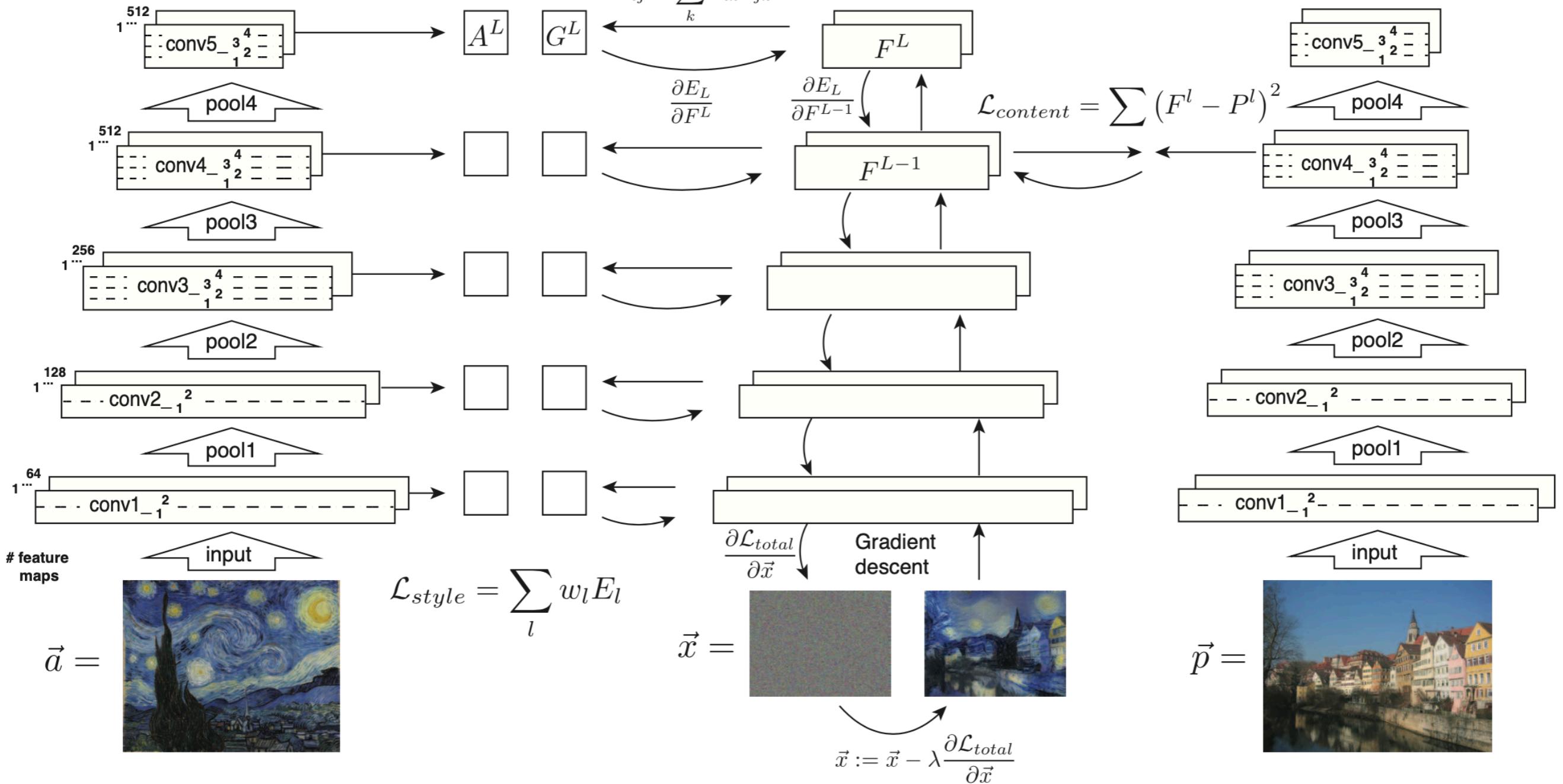
Extracting style and content from an image



Style transfer algorithm

$$E_L = \sum (G^L - A^L)^2$$

$$\mathcal{L}_{total} = \alpha \mathcal{L}_{content} + \beta \mathcal{L}_{style}$$



Style transfer

A



B



C



D



[Github repository \(one of many\)](#)

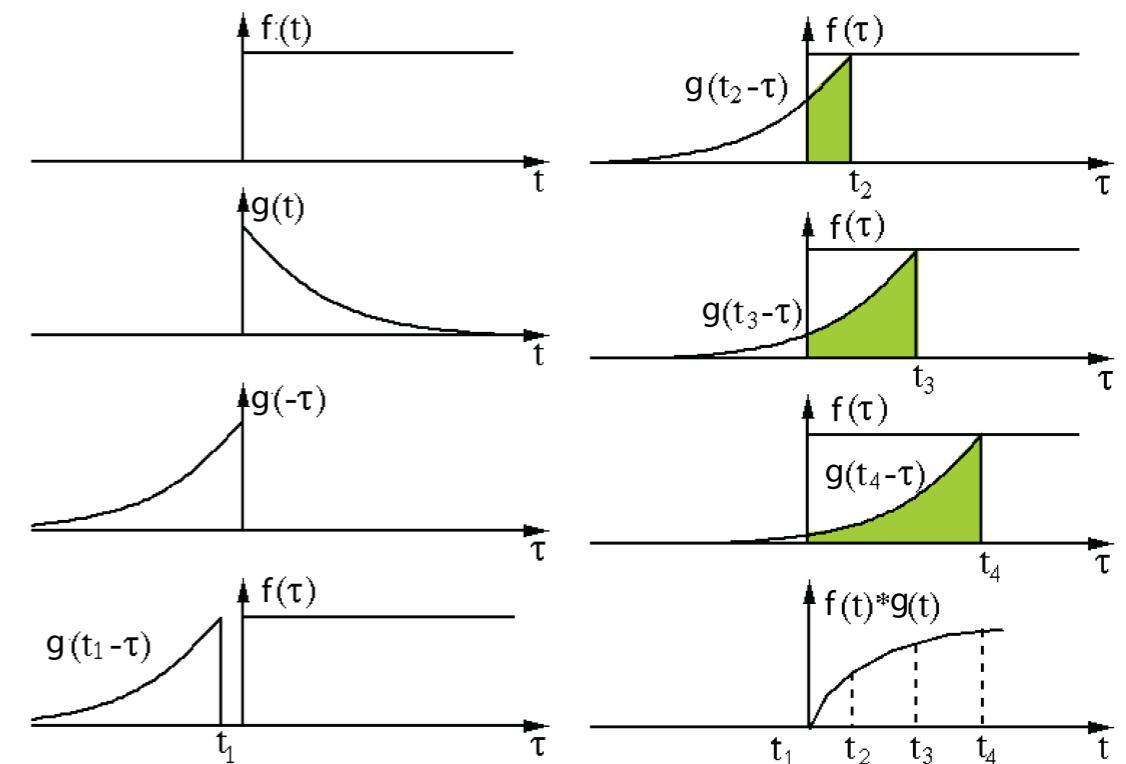
1D CNNs

1D CNNs

- CNNs perform well on computer vision problems, due to their ability to operate convolutionally **extract features from local input patches**.
- The same properties that make CNNs excel at computer vision also make them **highly relevant to sequence processing**
- 1D CNNs, typically used with **dilated kernels**, have been used with great success for audio generation and machine translation
- They can offer a **fast alternative to RNNs** for simple tasks such as text classification and timeseries forecasting

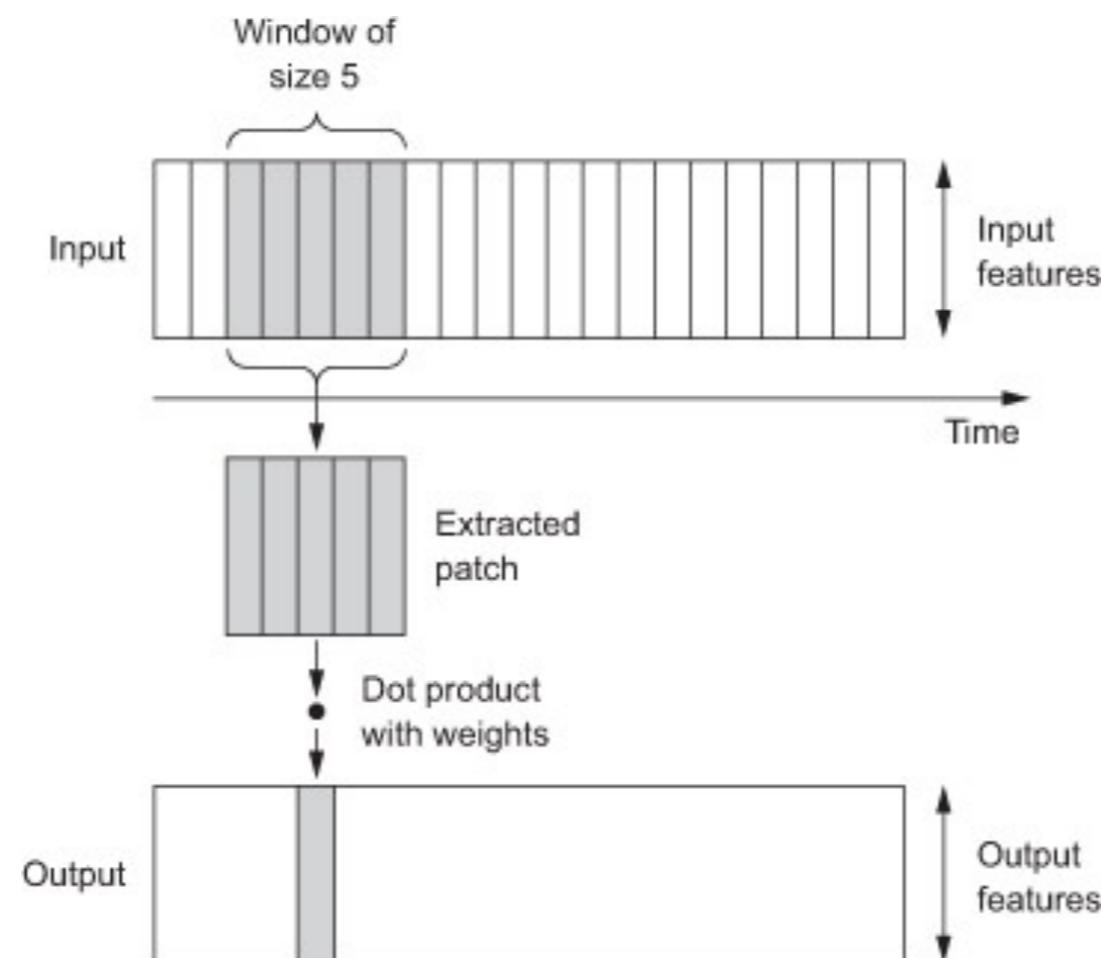
$$y(t) = f(t) * g(t) = \int f(\tau)g(t - \tau)d\tau$$

$$y[n] = f[n] * g[n] = \sum_{k=-\infty}^{\infty} f[k]g[n - k]$$



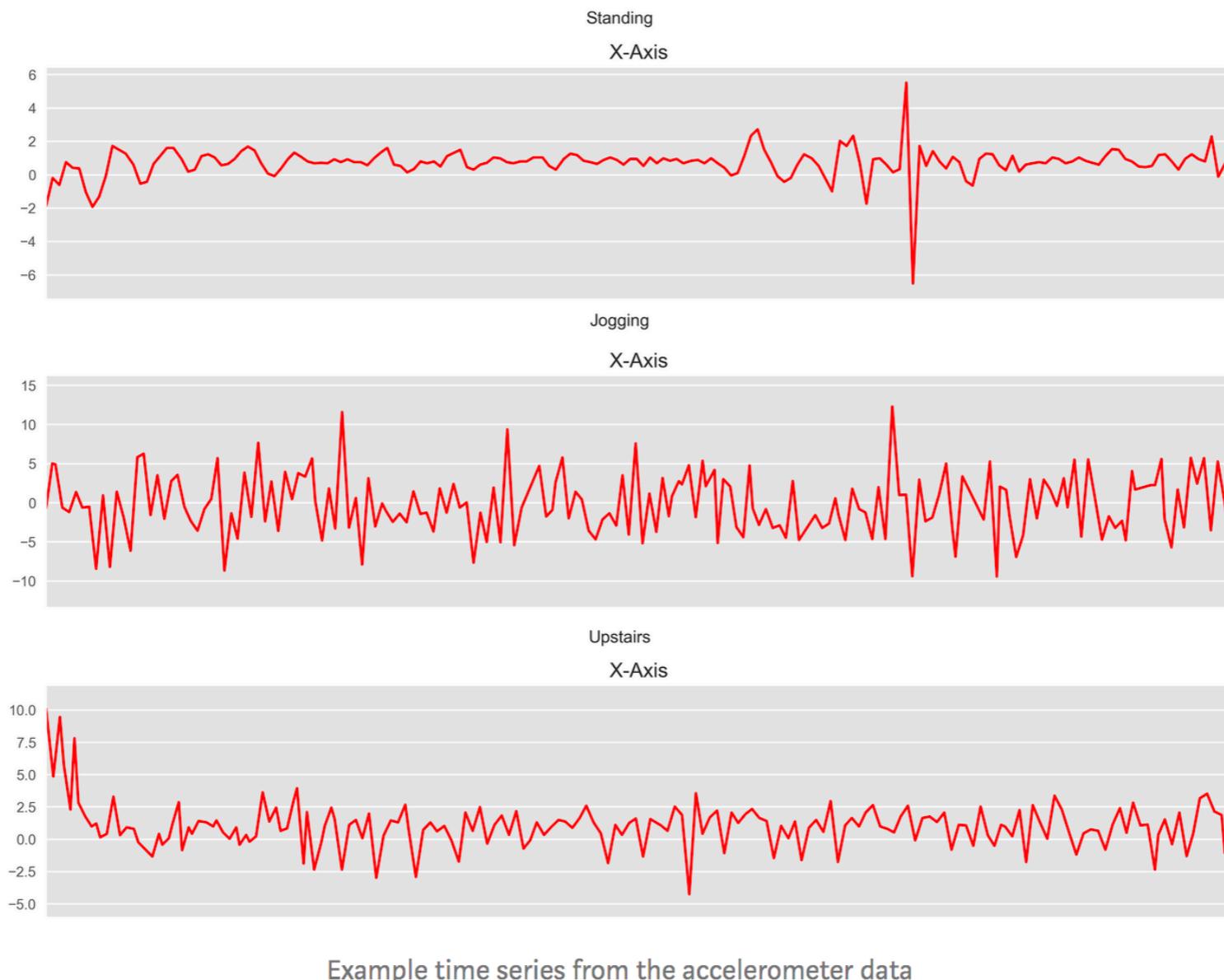
1D CNNs

- Combine finite-length filters with 1D patches (subsequences) from sequences
- Each output timestep is obtained from a temporal patch in the input sequence
- Convolution layers can recognize local patterns in a sequence
- 1D pooling layers, batch normalization,

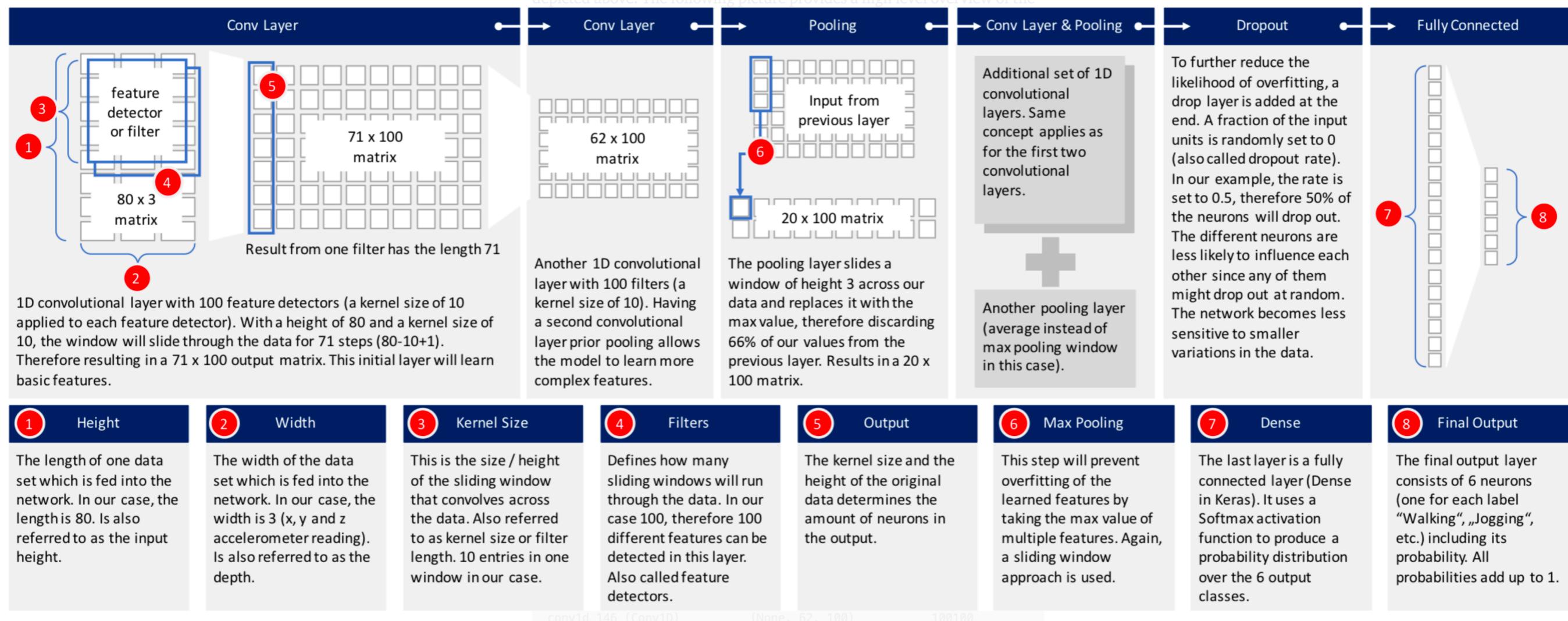


1D CNNs. An Example

- Example borrowed from [**this blog**](#) (also posted in Aula Global)



Predict the type of activity a user is performing (“Walking”, “Jogging” or “Standing”, ...)
There are 6 activities in total.



Layer (type)	Output Shape	Param #
reshape_45 (Reshape)	(None, 80, 3)	0
conv1d_145 (Conv1D)	(None, 71, 100)	3100
conv1d_146 (Conv1D)	(None, 62, 100)	100100
max_pooling1d_39 (MaxPooling)	(None, 20, 100)	0
conv1d_147 (Conv1D)	(None, 11, 160)	160160
conv1d_148 (Conv1D)	(None, 2, 160)	256160
global_average_pooling1d_29	(None, 160)	0
dropout_29 (Dropout)	(None, 160)	0
dense_29 (Dense)	(None, 6)	966

Total params: 520,486
Trainable params: 520,486
Non-trainable params: 0

None

Accuracy on test data: 0.92

Loss on test data: 0.39

Graph CNNs

Published as a conference paper at ICLR 2017

SEMI-SUPERVISED CLASSIFICATION WITH GRAPH CONVOLUTIONAL NETWORKS

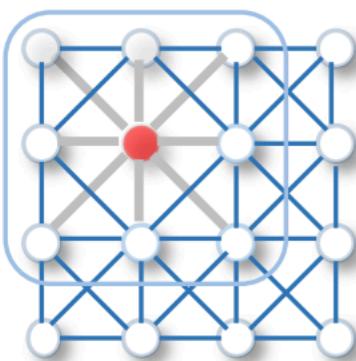
Thomas N. Kipf
University of Amsterdam
T.N.Kipf@uva.nl

Max Welling
University of Amsterdam
Canadian Institute for Advanced Research (CIFAR)
M.Welling@uva.nl

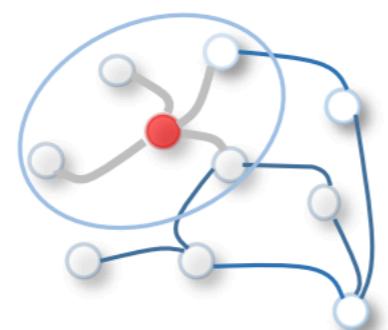
uc3m

Graph CNNs (GCN)

- A graph convolutional network (GCN) is a neural network that operates on graphs.
- Given a graph $G = (V, E)$, a GCN takes as input
 1. An input feature matrix \mathbf{X} of dimension $N \times F$, where N is the number of nodes and F is the number of features per node
 2. An $N \times N$ matrix representation of the graph structure (adjacency matrix \mathbf{A}).
- The first hidden layer in the GCN process each node in G according to \mathbf{X} and \mathbf{A} , $\mathbf{H} = f(\mathbf{X}, \mathbf{A})$
- \mathbf{H} is a node feature matrix $N \times F'$ where each row is a new feature representation of the node



(a) 2D Convolution. Analogous to a graph, each pixel in an image is taken as a node where neighbors are determined by the filter size. The 2D convolution takes the weighted average of pixel values of the red node along with its neighbors. The neighbors of a node are ordered and have a fixed size.

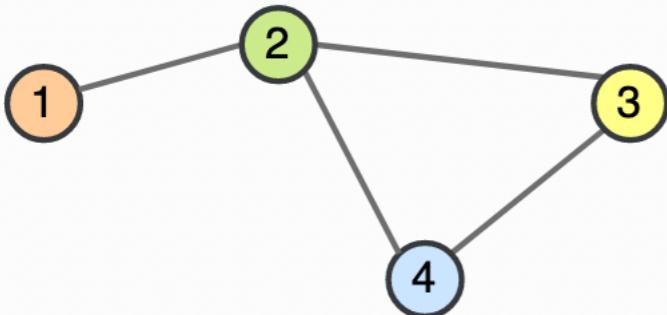


(b) Graph Convolution. To get a hidden representation of the red node, one simple solution of the graph convolutional operation is to take the average value of the node features of the red node along with its neighbors. Different from image data, the neighbors of a node are unordered and variable in size.

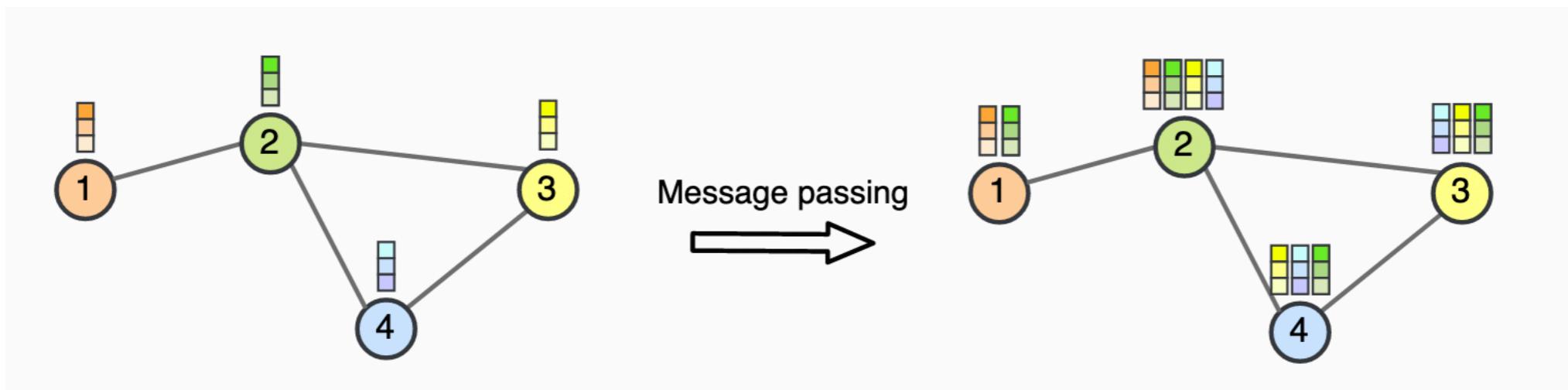
- At each layer, these features are aggregated to form the next layer's features using a propagation rule f .

Fig. 1: 2D Convolution vs. Graph Convolution.

Graph CNNs (GCN)



$$\mathbf{x}_v^{(\ell+1)} = \mathbf{W}^{(\ell+1)} \sum_{w \in \mathcal{N}(v) \cup \{v\}} \frac{1}{c_{w,v}} \cdot \mathbf{x}_w^{(\ell)}$$



Graph Neural Networks

Graph CNNs (GCN)

A Comprehensive Survey on Graph Neural Networks

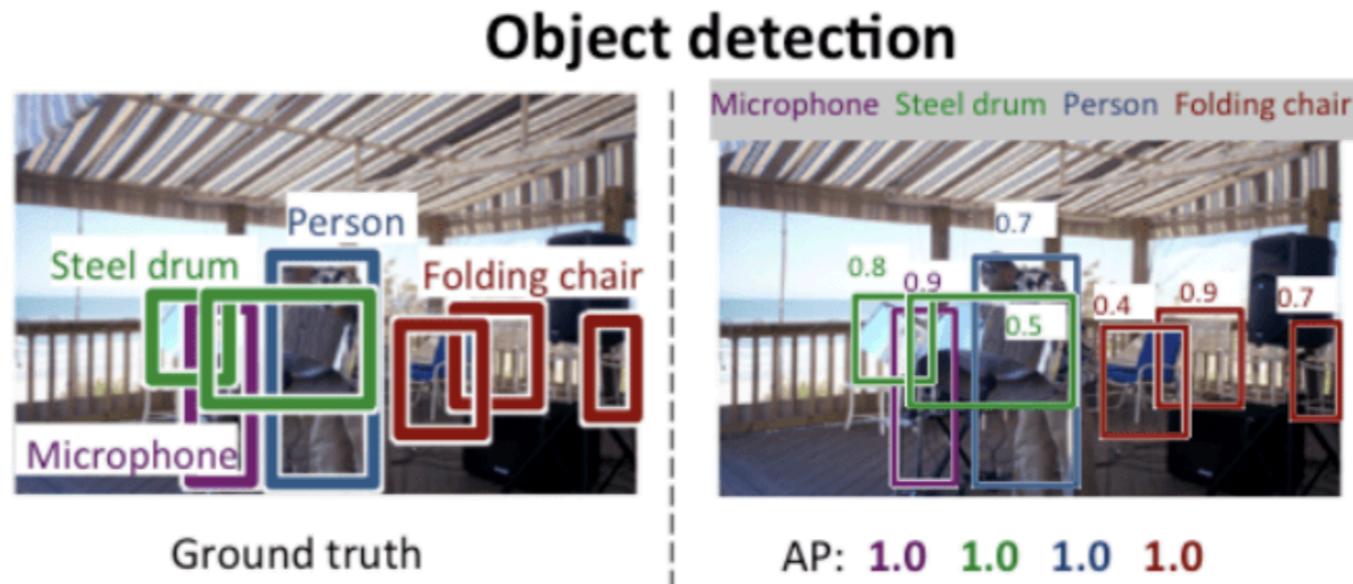
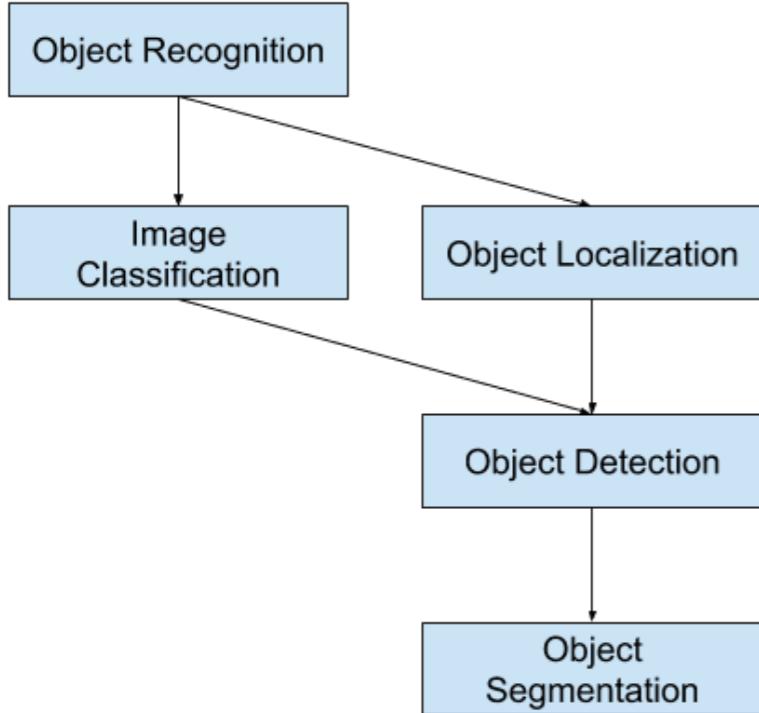
Zonghan Wu^{ID}, Shirui Pan^{ID}, Member, IEEE, Fengwen Chen, Guodong Long^{ID}, Chengqi Zhang^{ID}, Senior Member, IEEE, and Philip S. Yu, Life Fellow, IEEE

A screenshot of a YouTube video player. The video title is "An Introduction to Graph Neural Networks: Models and Applications" by Microsoft Research. The video has 66,840 views and was uploaded on May 8, 2020. The thumbnail image shows a person's hands holding a blue backpack. The video player interface includes a play button, volume control, and a progress bar showing 0:05 / 58:59.

Learning Convolutional Neural Networks for Graphs

Object detection with Region CNNs (RCNNs)

Object Recognition



- **Region-based CNNs** (R-CNN) and extensions such as Fast R-CNN and YOLO (you only look once) are the most dominant approach for **object localisation and recognition**.

Rich feature hierarchies for accurate object detection and semantic segmentation

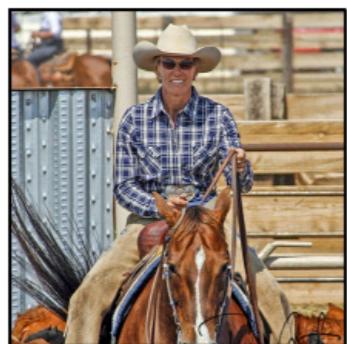
Tech report (v5)

Ross Girshick Jeff Donahue Trevor Darrell Jitendra Malik
UC Berkeley

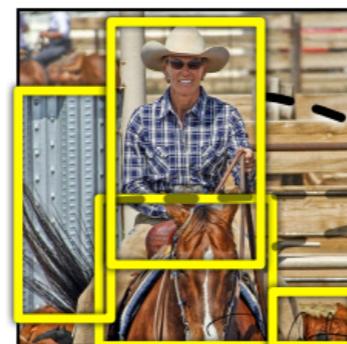
{rgb, jdonahue, trevor, malik}@eecs.berkeley.edu

- Locate regions of interest
- Classify using a pre-trained CNN for object clarification

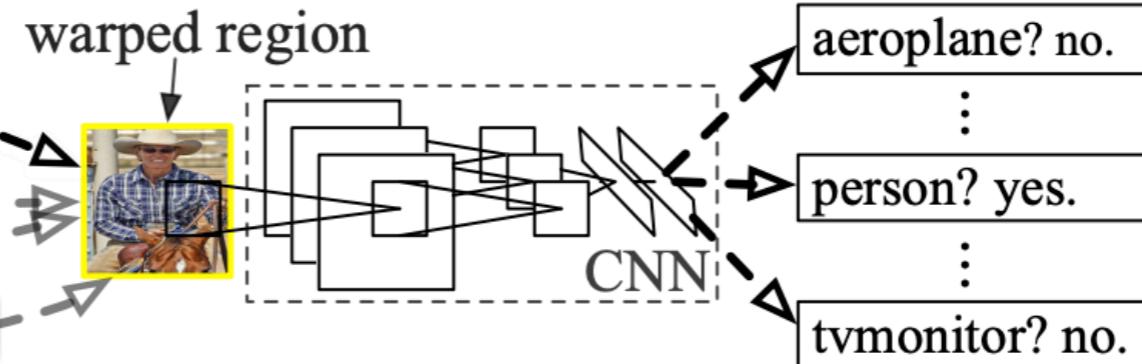
R-CNN: *Regions with CNN features*



1. Input image



2. Extract region proposals (~2k)



3. Compute CNN features

4. Classify regions

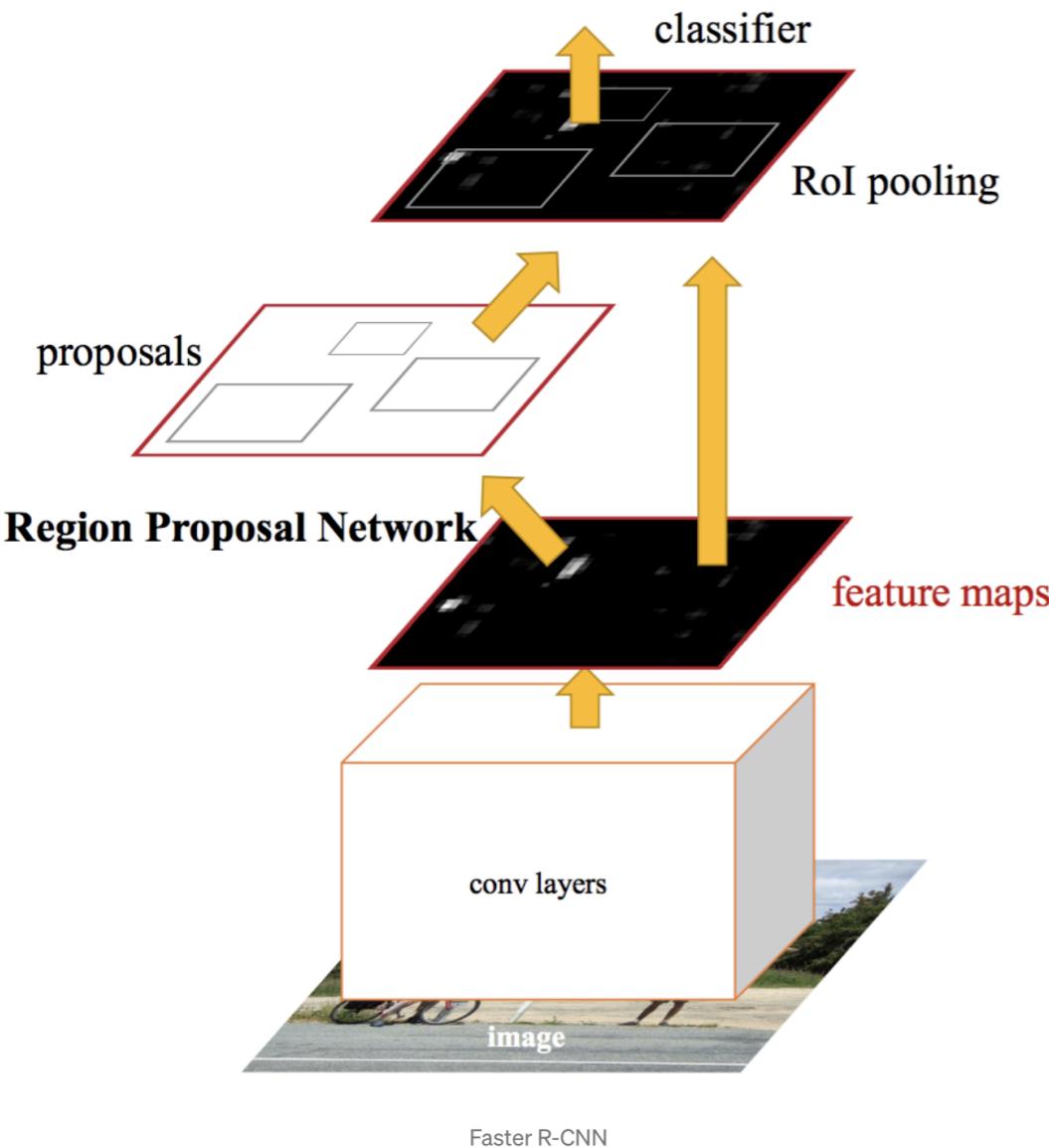
Selective Search for Object Recognition

J. R. R. Uijlings · K. E. A. van de Sande ·
T. Gevers · A. W. M. Smeulders



- One of the simplest
- First, **segment** different areas in the image
- Tons of algorithms for multi-scale **image segmentation** (thresholding, K-means, histogram based, edge detection)
- Second, define “squared” regions based on segmented areas.

Fast/Faster R-CNN: obtain region proposals



- Same idea, but you define regions proposals over the feature maps
- You only have to filter the input image once
- In Faster R-CNN they replace selective search by **regions learned by the network itself**

Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks

Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun

Image segmentation using CNNs

Image Segmentation Using Deep Learning: A Survey

Shervin Minaee, Yuri Boykov, Fatih Porikli, Antonio Plaza, Nasser Kehtarnavaz, and Demetri Terzopoulos

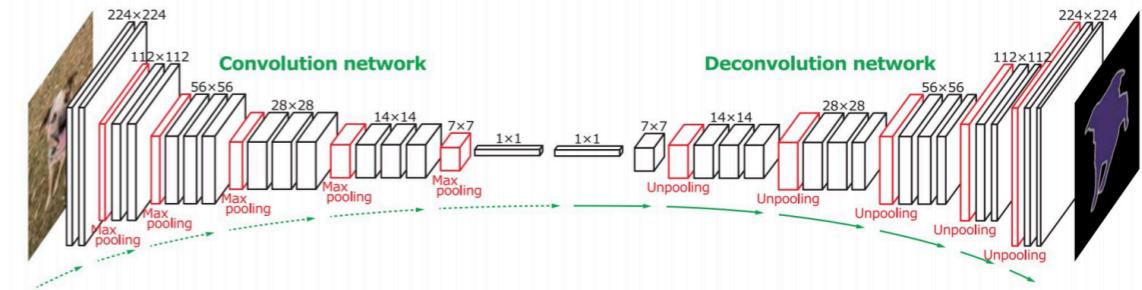


Fig. 11. Deconvolutional semantic segmentation. Following a convolution network based on the VGG 16-layer net, is a multi-layer deconvolution network to generate the accurate segmentation map. From [42].

- Most powerful methods are based on **encoder-decoder networks**
- We will explore this in the block on **unsupervised deep learning**

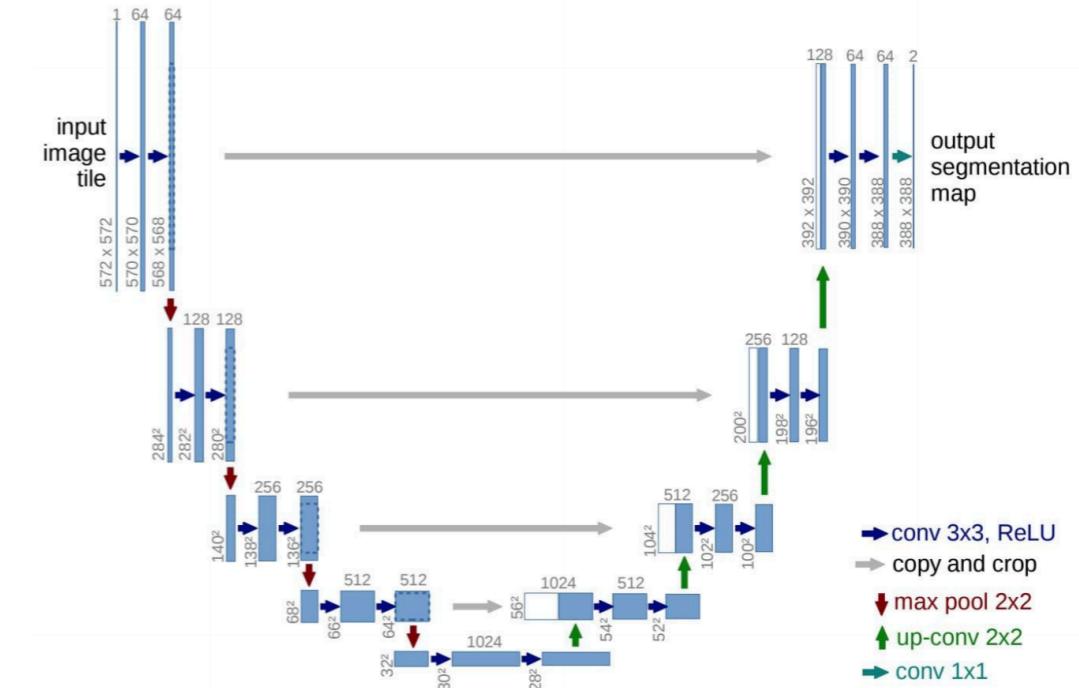


Fig. 14. The U-net model. The blue boxes denote feature map blocks with their indicated shapes. From [49].