# Convolutional Neural Networks

Training course: Training course: Machine learning and Destination Earth

Presented by Maria Luisa Taccari
with contributions and materials from Vera Gahlen

ECMWF

marialuisa.taccari@ecmwf.int

**ECMWF**

# Outline

- **Motivation for convolutions**

- **What is a convolution?**

- **Convolution's arithmetic**

- **Building a Convolutional Neural Network**

- **Connecting image structure and CNN architecture**

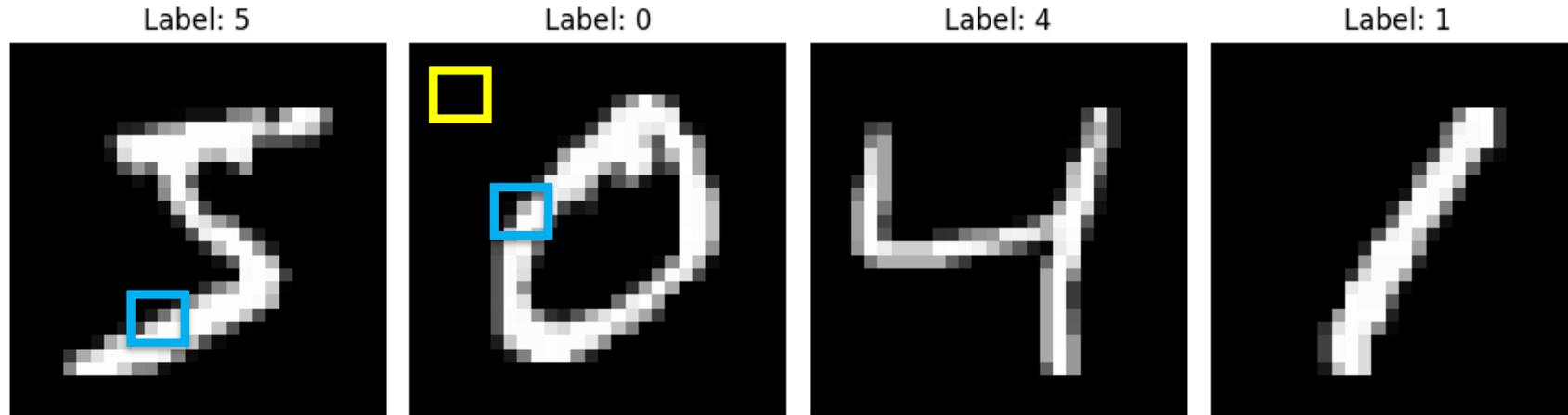- **Popular CNN-based architectures - ResNets, U-nets**

# Motivation for convolutions

# Why not use an MLP?

- Expressiveness?
  - MLPs are universal approximators — special architectures aren't needed to *represent* good solutions.

- Compute?
  - Well-designed architectures encode inductive biases that constrain the hypothesis space, enabling similar performance with less data and less computation.
  - Historically, before CNNs, training large networks for vision tasks was **infeasible — more compute alone wasn't enough**.

- Optimization!
  - Architectures shape the loss landscape, making good solutions ***discoverable* by gradient descent**.
  - Modern DL relies more on architectures than on optimizers to make training feasible.

- Generalization
  - Constraining the hypothesis space also helps prevent overfitting.
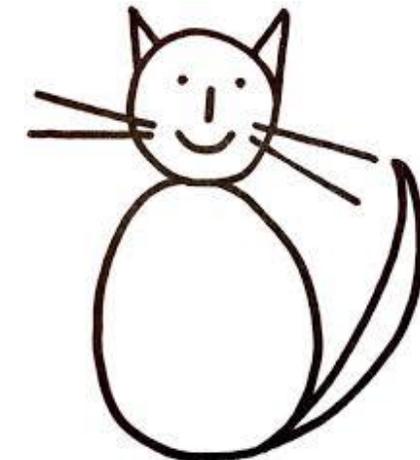
>> Architectures don't just tell us *what can be represented* — they determine *what can be **found***.

# What structure do images have?



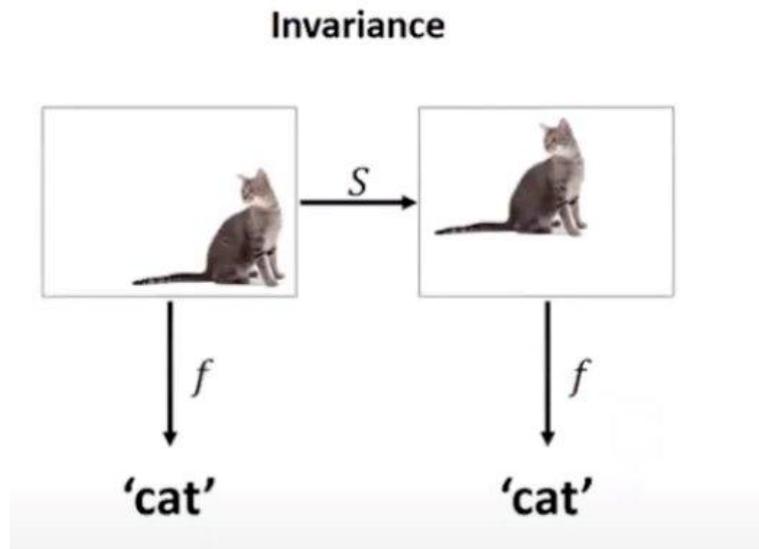Label: 5     Label: 0     Label: 4     Label: 1

MNIST Dataset of handwritten digits

- **Spatial locality**
  To make sense of a pixel, we need (only!) the surrounding pixels

- **Compositional structure**
  Simple features (edges) combine into complex ones (textures, shapes)
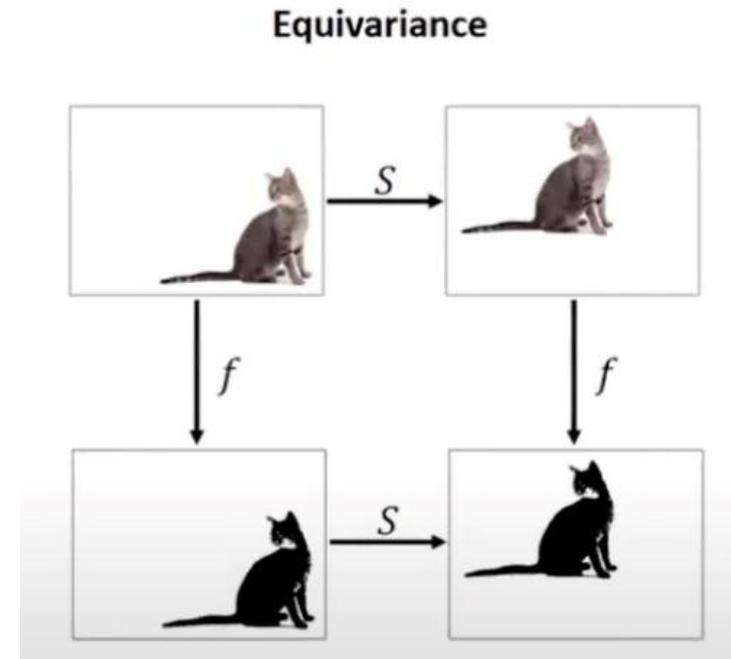
# What structure do images have?



Translation Invariance

Invariance

'cat' → 'cat'

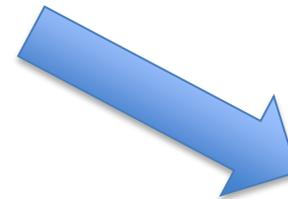Translation Equivariance

Equivariance

# Motivation for convolutions

## What structure do images have?

- **Locality** — Neighboring pixels tend to form meaningful local structures (edges, corners, textures)
- **Compositional structure** — Simple features combine into complex ones
- **Translation equivariance** — Shifting the input shifts the output the same way

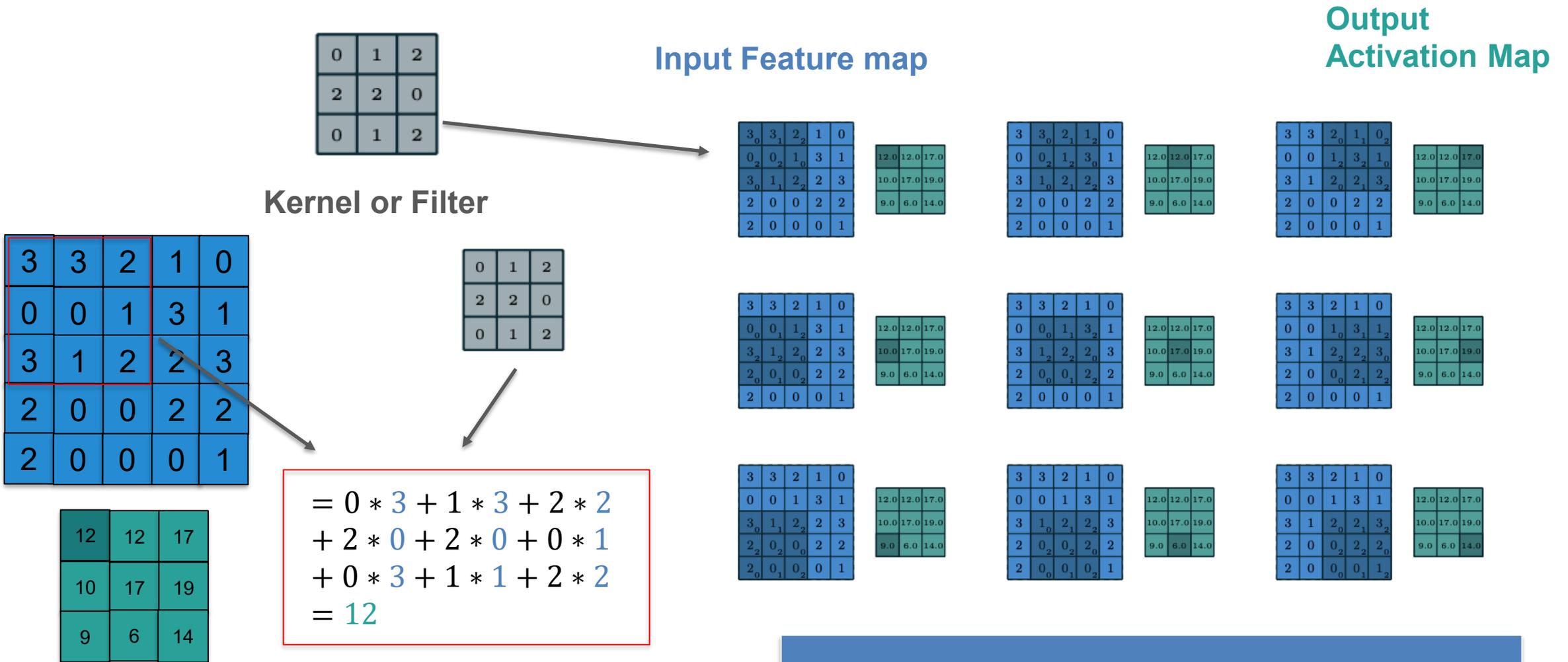We want the model to perform **localized, translation-equivariant** *aggregation/extraction* of visual information

The architecture needs to suit the problem

*Sliding filters* apply the same local pattern extraction everywhere in the image

# What is a convolution?

# What is a convolution?



Kernel or Filter

Input Feature map

Output Activation Map

$$= 0 * 3 + 1 * 3 + 2 * 2$$
$$+ 2 * 0 + 2 * 0 + 0 * 1$$
$$+ 0 * 3 + 1 * 1 + 2 * 2$$
$$= 12$$

Convolve a filter with the image = spatially sliding it over the image and computing the dot product

*A guide to convolution arithmetic for deep Learning*
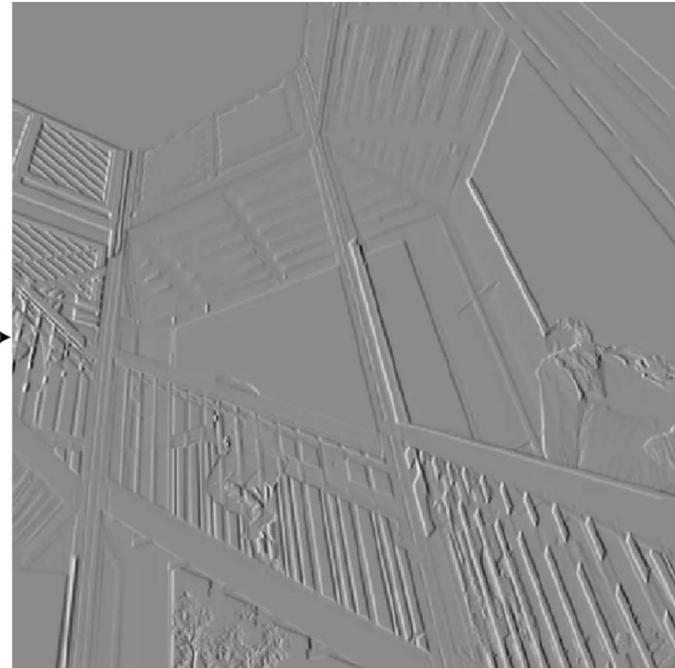*Dumoulin V., Visin. F, 2018, arXiv:1603.07285*

# What is a convolution?

$$\begin{bmatrix} +1 & 0 & -1 \\ +2 & 0 & -2 \\ +1 & 0 & -1 \end{bmatrix}$$

Horizontal Sobel kernel

Applying a vertical edge detector kernel

https://setosa.io/ev/image-kernels/

The kernel/filter is the **trainable** part of the convolutional layer

**ECMWF**

**EUROPEAN CENTRE FOR MEDIUM-RANGE WEATHER FORECASTS**

# Convolution versus MLP

**MLP(FC) Architecture**

Flatten   Matmul   Output class

Learned weights

32 x 32   1024   1024 x 10   10

- Poor scaling with image size
- Inefficient weight use – no "weight sharing"
- FC do not provide translation invariance nor equivariance

**CNN Architecture**

Flatten   Dot product   Output

5 x 5 patch

**Weight Matrix**

**Output Activation Map**

32 x 32   25   25 x 1   28x28x1

*Lecture 2A: Convolutional Neural Networks (Full Stack Deep Learning - Spring 2021)*

# Convolution's Arithmetic

# Convolution's arithmetic



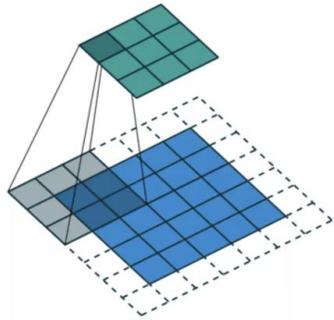The shape of the output feature map (W,H) is defined based on:
- Shape of the input feature map (W,H)
- The Kernel size (w,h)
- The stride (s)
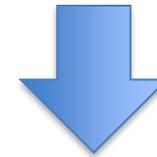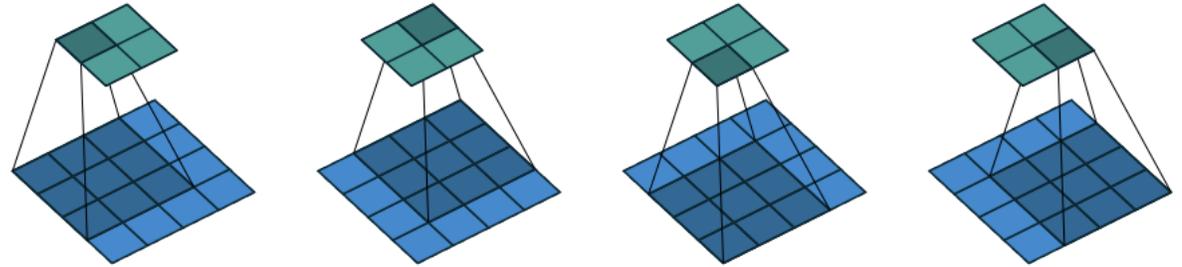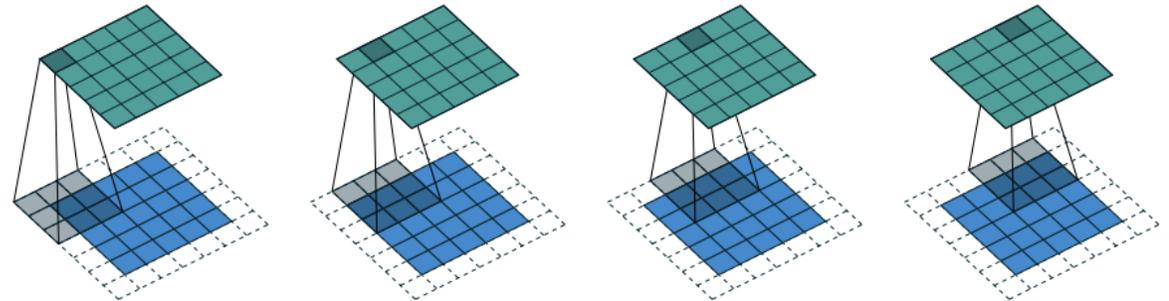- The padding (p)

# Convolution's arithmetic

**Input Feature map 4x4**

Stride

s=1

- Convolutions can subsample the image by jumping across some locations — this is called 'stride'

**Input Feature map 5x5**

s=2

# Convolution's arithmetic

**Padding**

- Padding solves the problem of filters running out of image
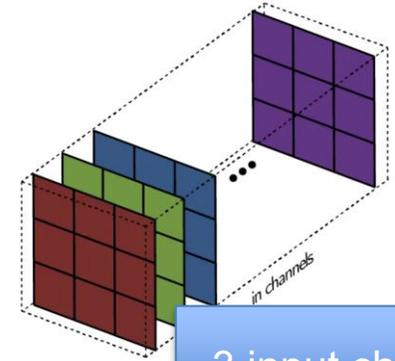- Done by adding extra rows/cols to the input (usually set to 0)
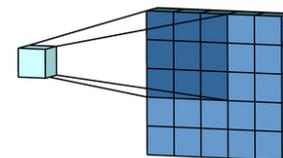
p=0

p=1

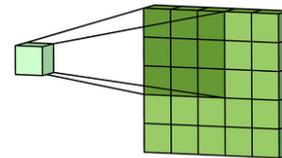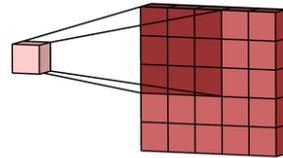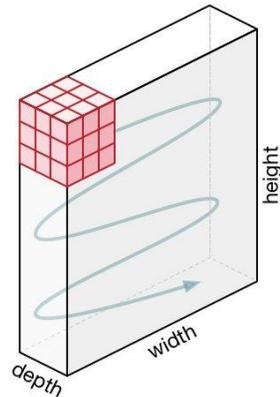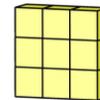# Convolution's arithmetic



**RGB Image**
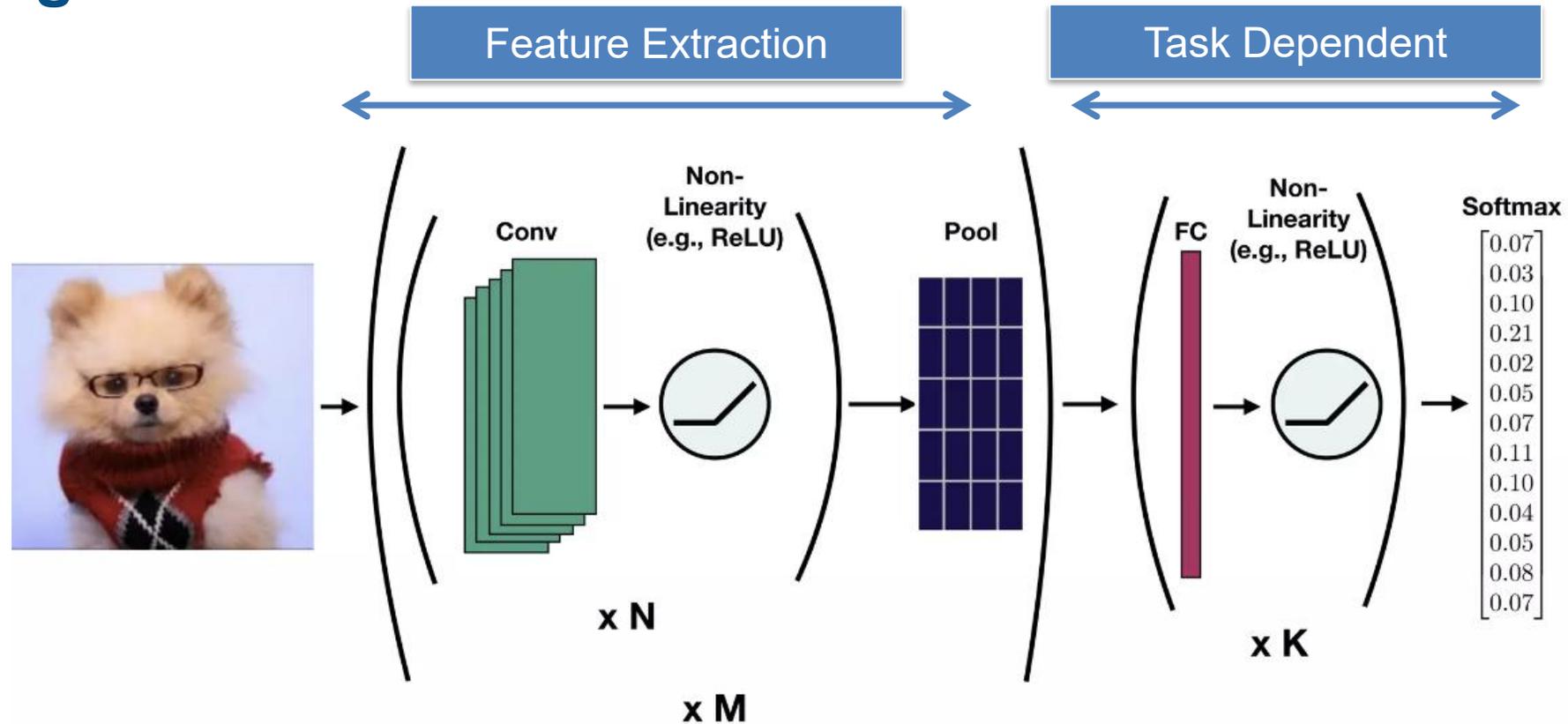
3 input channels

1 filter with 3 kernels

**Can't forget the bias term!**

Multiple filters in one layer >> the same number of output channels (~ different views of the image)

1 output channel

*Intuitively Understanding Convolutions for Deep Learning Towards Data Science*

# Building a
# Convolutional Neural Network (CNN)
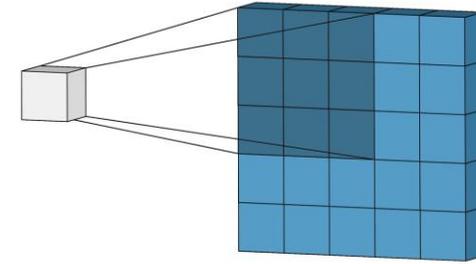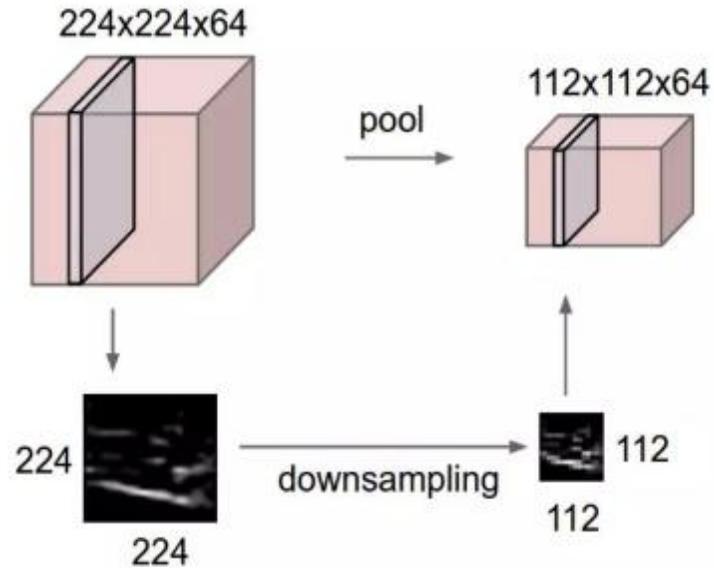
# Building a CNN



Feature Extraction | Task Dependent

Feature Extraction Blocks:
- Convolutional layers extract local features
- Pooling layers aggregate the extracted local features

# Building a CNN
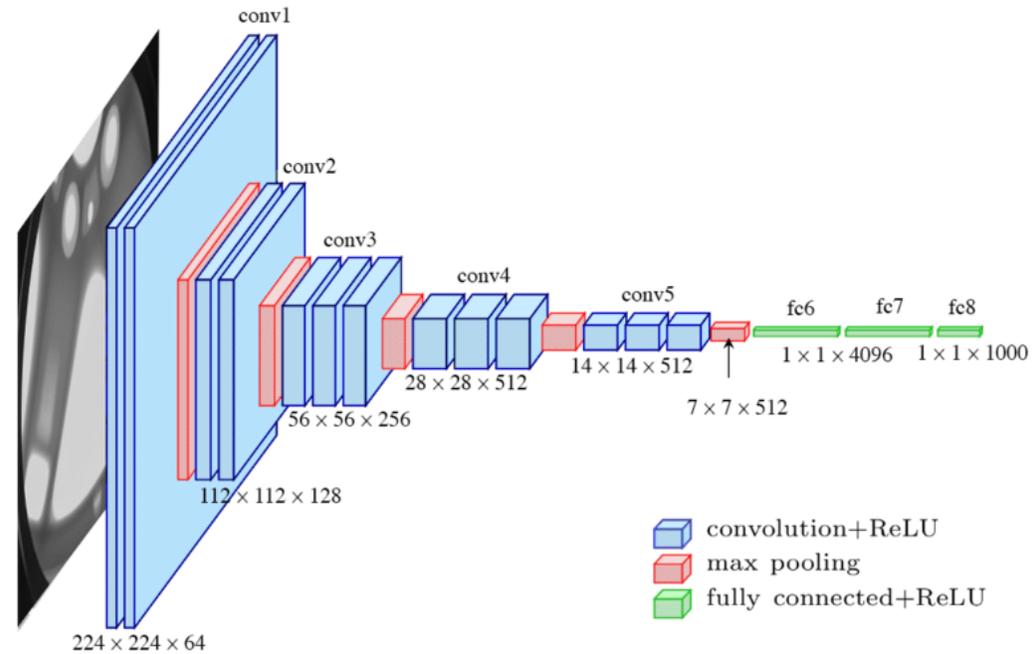


## Pooling



224x224x64

pool →
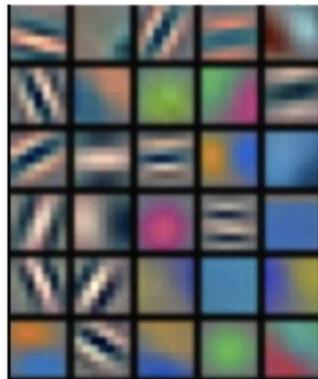
112x112x64

224

224

downsampling →

112

112

- Idea: downsample output activation maps from convolutional layers
- Similar to convolutional layer: Sliding window
- Instead of the dot product with the kernel, apply an aggregation function (e.g. max, avg)

- Reduces (H,W) by using bigger strides
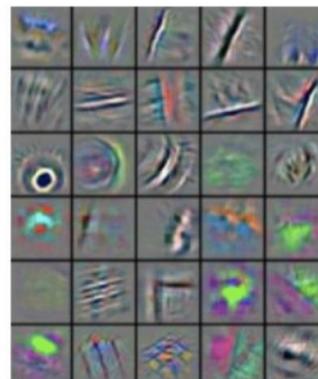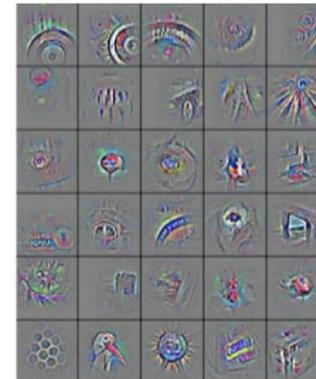- Preserves the number of channels

# Building a CNN

# Connecting image structure and CNN architecture

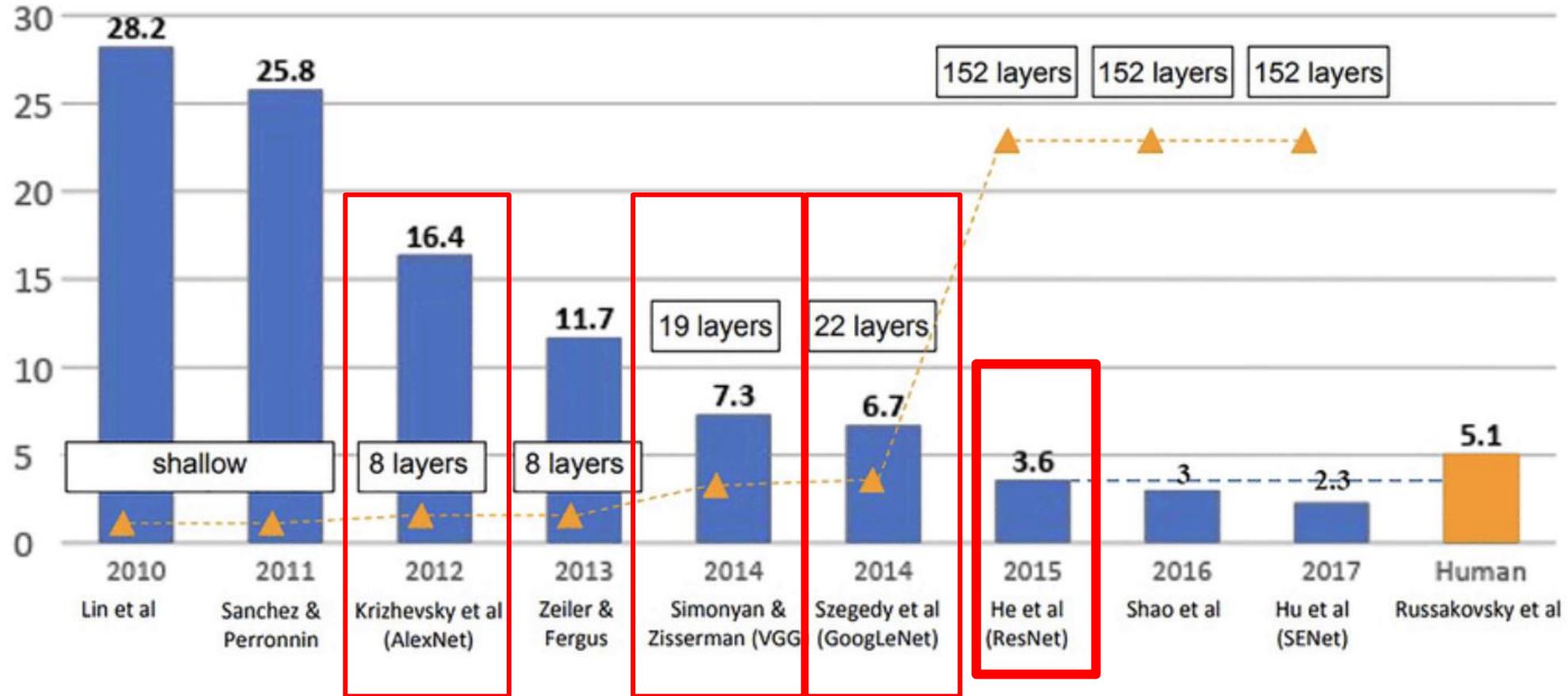| Image property | Architectural reflection | Emergent result |
|---|---|---|
| **Locality** | Sparse local connections with shared filters | Early layers detect local features (edges, textures) |
| **Translation equivariance** | Sliding windows (convolution, pooling) | Feature maps shift with input |
| **Compositional structure** | Layered composition | Hierarchical representations (simple → complex) |

# Popular Convolutional Neural Network Architectures

# Popular CNN-based architectures

*Lecture 2B: Convolutional Neural Networks*
*(Full Stack Deep Learning - Spring 2021)*

**ECMWF**  EUROPEAN CENTRE FOR MEDIUM-RANGE WEATHER FORECASTS
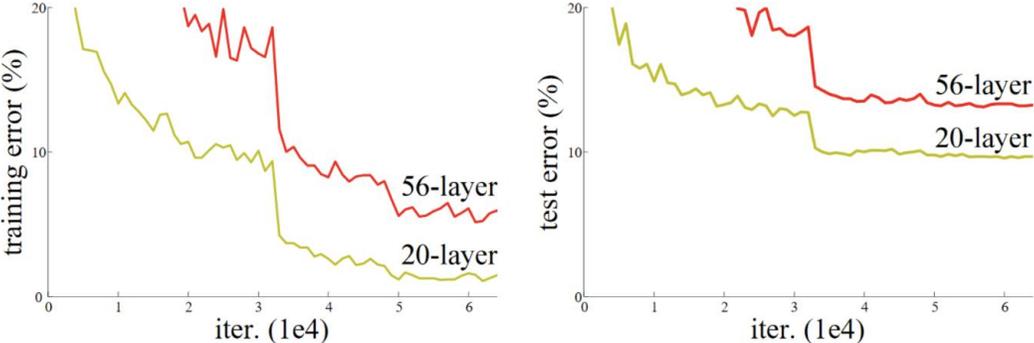
27
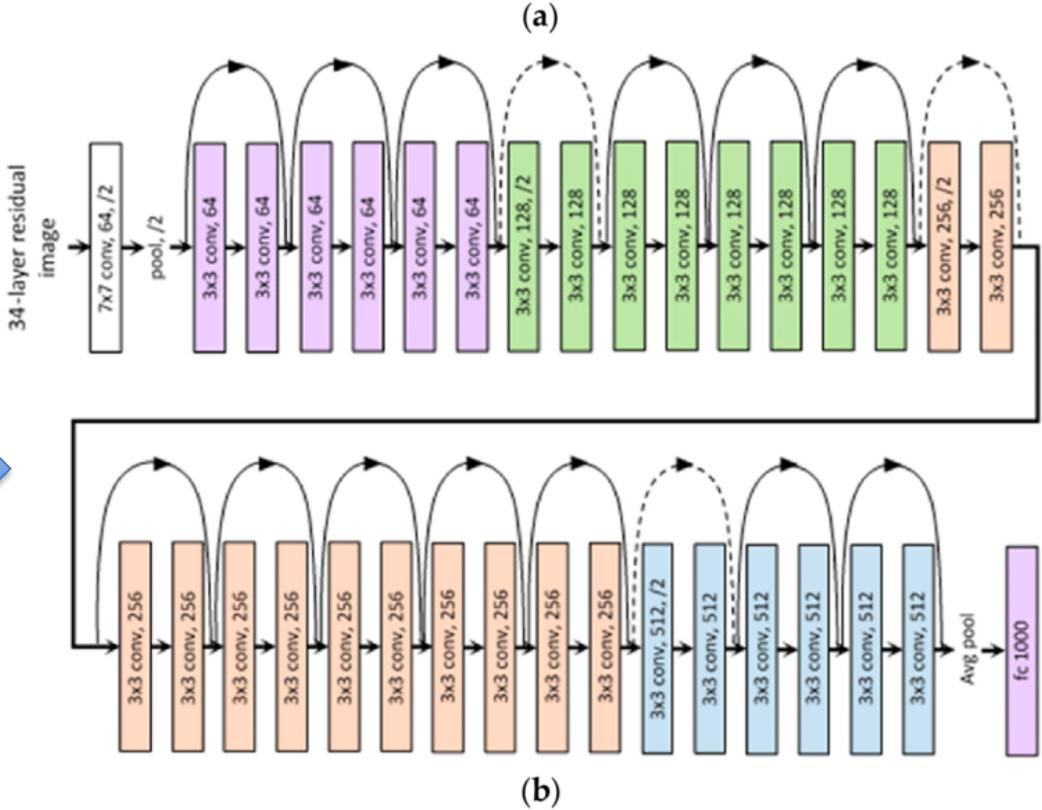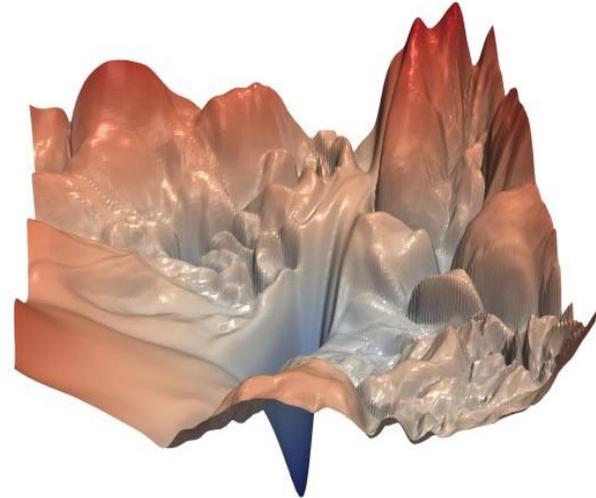
# Popular CNN-based architectures - ResNets

## ResNet



Figure 14.1: Training of networks of different depth (courtesy of Kaiming He et al.)

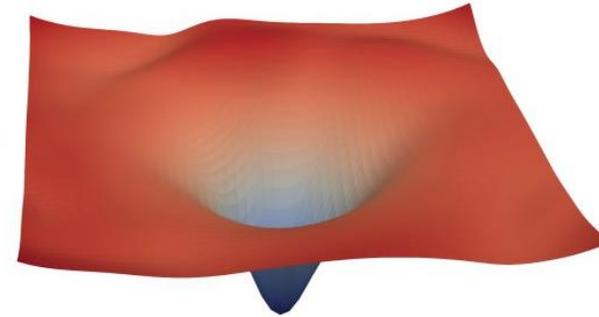"Deep Residual Learning for Image Recognition"



ResNet-34 Layered architecture

# Popular CNN-based architectures



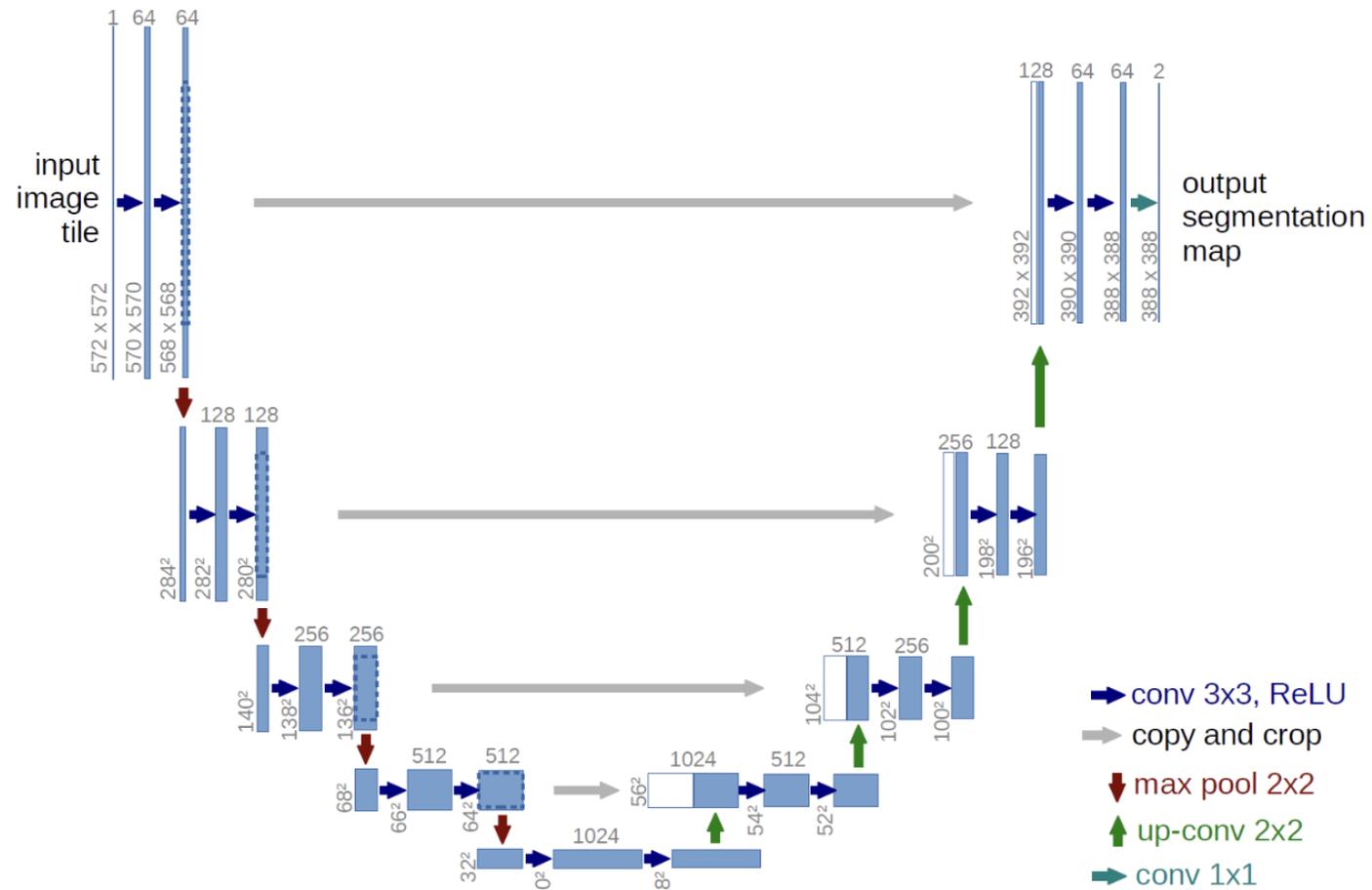(a) without skip connections      (b) with skip connections

The loss surfaces of ResNet-56 with and without skip connections

Using skip connections helps smooth the loss function, which makes training easier as it avoids falling into a very sharp area.

**Visualizing the Loss Landscape of Neural Nets**
https://arxiv.org/abs/1712.09913
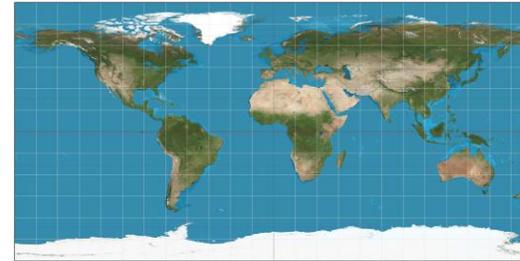
# Popular CNN-based architectures – U-Nets



U-Net: Convolutional Networks for Biomedical Image Segmentation
arxiv.org/abs/1505.04597

# Popular CNN-based architectures



**Periodic Convolutions**

https://github.com/pangeo-data/WeatherBench/blob/master/src/train_nn.py#L102

**Spatio-Temporal Data - ConvLSTMS**



*Convolutional LSTMs for Cloud-Robust Segmentation of Remote Sensing Imagery*
*Marc Rußwurm*

ECMWF    EUROPEAN CENTRE FOR MEDIUM-RANGE WEATHER FORECASTS    31

# References

- Deep Learning book CNNs - https://www.deeplearningbook.org/contents/convnets.html
- Understanding Deep Learning – CNN chapter - https://udlbook.github.io/udlbook/
- CNN feature visualization - https://distill.pub/2017/feature-visualization/
- CNN feature visualization - https://arxiv.org/pdf/1311.2901.pdf
- Intuitively Understanding Convolutions for Deep Learning
- A guide to convolution arithmetic for deep Learning - Dumoulin V., Visin. F, 2018, arXiv:1603.07285
- Lecture 2A and Lecture 2B Convolutional Neural Networks (Full Stack Deep Learning - Spring 2021)
- Invariance and equivariance - https://www.doc.ic.ac.uk/~bkainz/teaching/DL/notes/equivariance.pdf
- Imperial's Deep learning course: Equivariance and Invariance, Bernhard Kainz
- The CNN notebook from Lisa Zhang