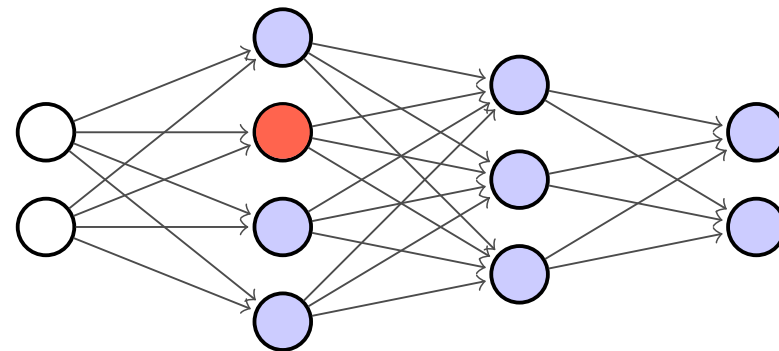


EPFL

Image Processing I

Chapter 6 Convolutional neural networks

Prof. Michael Unser, LIB



December 2025

OUTLINE

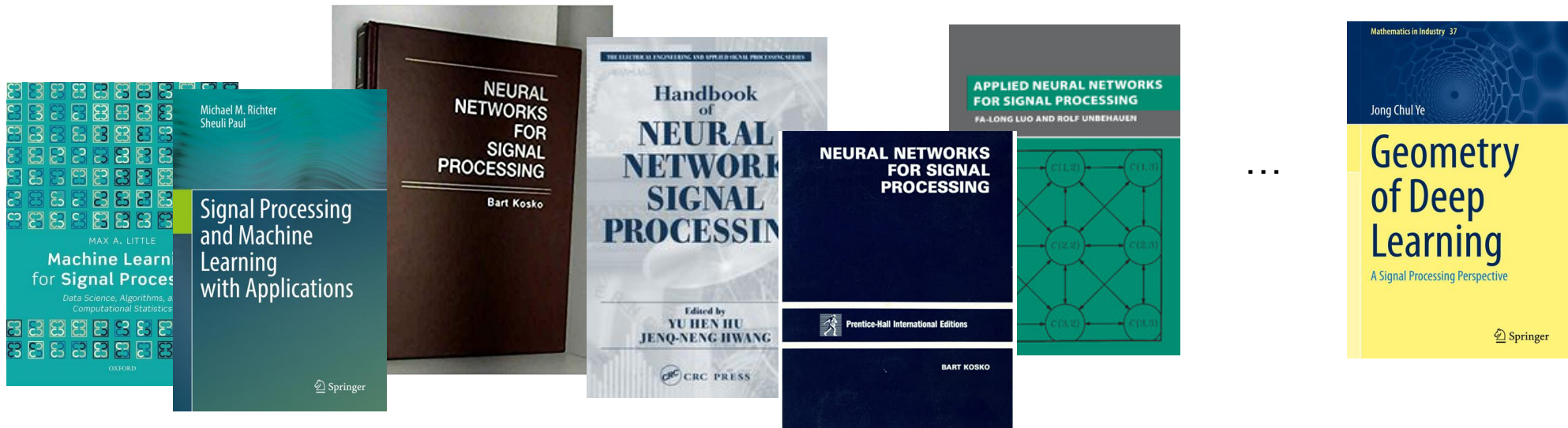
- Introduction
 - The (deep) learning (r)evolution in signal processing
 - Artificial neurons
 - Neural networks architectures for image processing
- Basic Components of CNNs
 - Operator-based formalism
 - Composition properties
 - Pooling
 - Continuity properties
- CNNs in Practice
 - Deep learning pipeline
 - Denoising
 - Segmentation

The (deep) learning (r)evolution in image processing

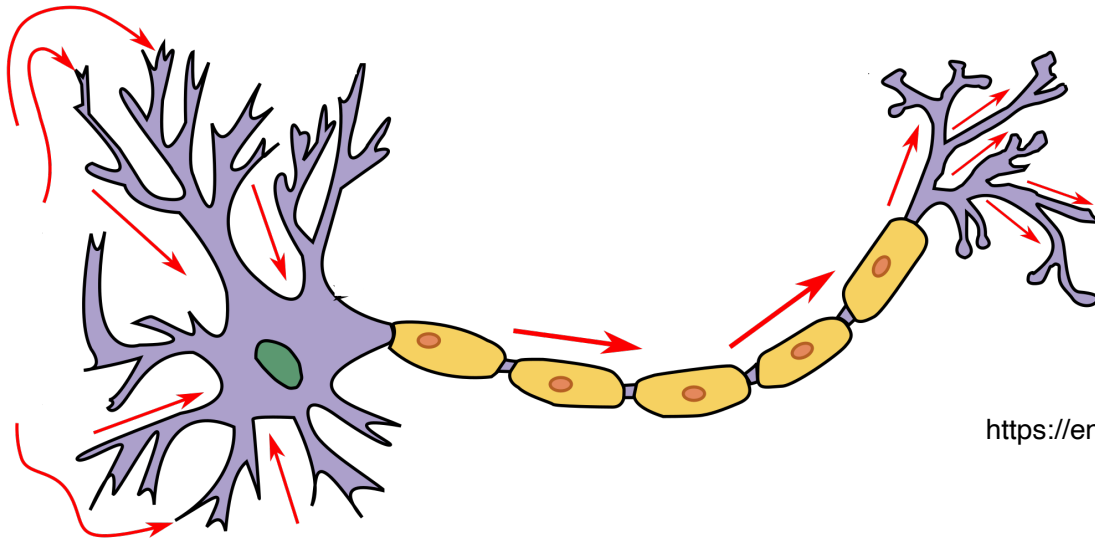
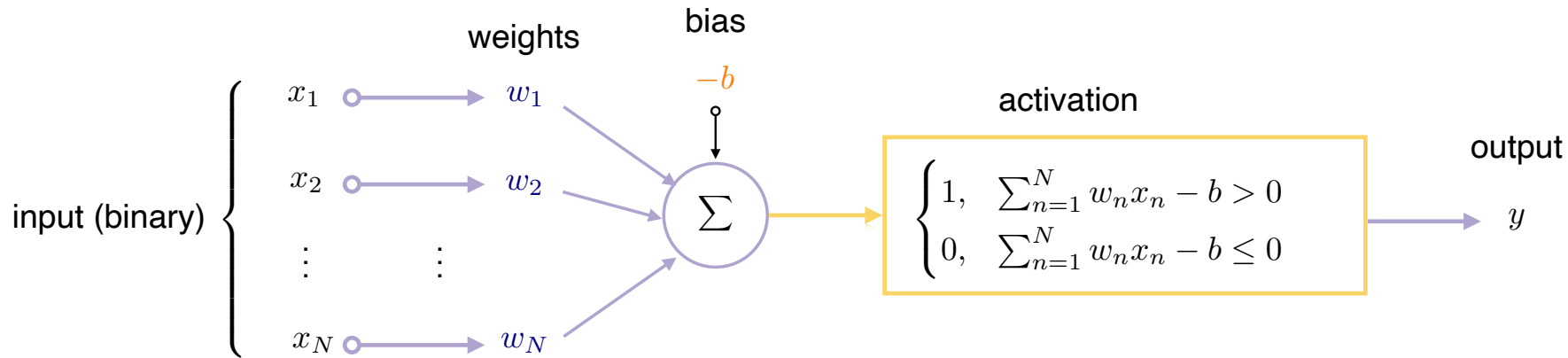
Special issues



Flurry of new textbooks on neural networks

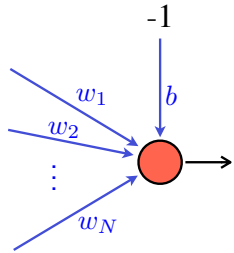


Formal model of neuron (McCulloch & Pitt)



https://en.wikipedia.org/wiki/Artificial_neuron

Artificial neurons



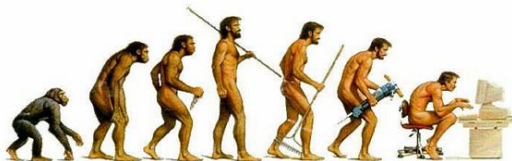
Definition: An artificial neuron with **weights** $\mathbf{w} = (w_1, \dots, w_N) \in \mathbb{R}^N$, **bias** $b \in \mathbb{R}$ and **activation function** $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is defined as the function $f : \mathbb{R}^N \rightarrow \mathbb{R}$

$$f(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} - b) = \sigma\left(\sum_{n=1}^N w_n x_n - b\right).$$

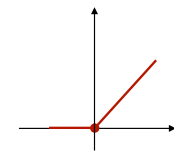
■ Examples of activation functions

■ Threshold Logic Unit (Heaviside): $\text{TLU}(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases}$ (McCulloch & Pitt 1943; Rosenblatt 1957)

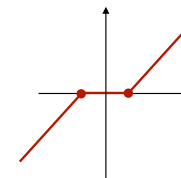
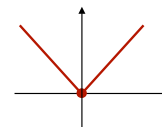
■ Sigmoid function: $\sigma(x) = \frac{1}{1 + e^{-x}}$ (Rumelhart 1986, ...)



■ Rectified Linear Unit: $\text{ReLU}(x) = x_+ = \max(0, x)$



■ And variants



Neural network architectures for image processing

Neural networks are constructed from the **composition of basic modules** that can be chained at will.

Convolutional neural networks (CNN), in particular, are inspired by the structure of the primary visual cortex.

They have an architecture that is well suited for image processing.

■ Basic modules

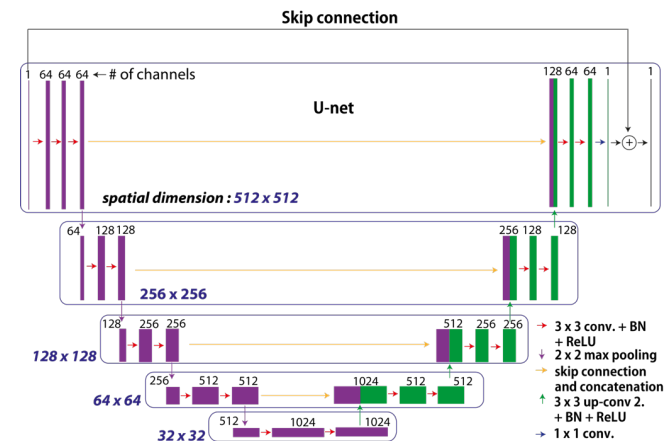
- Multi-channel convolution operators (filters)
- Pointwise nonlinearities
- Pooling: linear combination, flattening, sub-sampling, ...

Some modules—in particular, the filters—are adjustable.

The parameters of the CNN (weights) are set during the training procedure.

■ Training (not covered in this chapter)

- Requires a comprehensive collection of reference input-output pairs.
The larger the training set, the better!
- Formulated as a large-scale optimization problem
- Solved using some form of stochastic gradient algorithm (ADAM)
- Requires a lot of computational resources (GPU)



Basic Components of CNNs

- Operator-based formalism
- Composition
- Pooling
- Continuity and stability estimates

Operator-based formalism

Generic operator $T : \mathcal{X} \rightarrow \mathcal{Y}$ where \mathcal{X} and \mathcal{Y} are complete normed vector spaces.

$$\forall x \in \mathcal{X} : \quad y = T\{x\} \in \mathcal{Y}$$

- General multivariate nonlinear operator (for image patches)

$$T : \mathbb{R}^M \rightarrow \mathbb{R}^N$$

- Pointwise nonlinearity = activation function of neuron

$$T_{\text{acti}} : \mathbb{R} \rightarrow \mathbb{R} \quad \text{with} \quad T_{\text{acti}}\{x\} = \sigma(x + b), \quad b \in \mathbb{R}$$

- General multivariate linear operator (for fully connected layer)

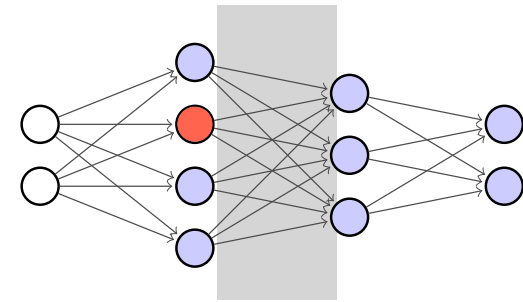
$$T_{\text{lin}} : \mathbb{R}^M \rightarrow \mathbb{R}^N \quad \text{with} \quad T_{\text{lin}}\{\mathbf{x}\} = \mathbf{W}\mathbf{x} \quad \text{where} \quad \mathbf{W} \in \mathbb{R}^{N \times M}$$

- Image-to-image operator

$$T : \ell_2(\mathbb{Z}^d) \rightarrow \ell_2(\mathbb{Z}^d)$$

Most operators of interest are **shift-invariant**

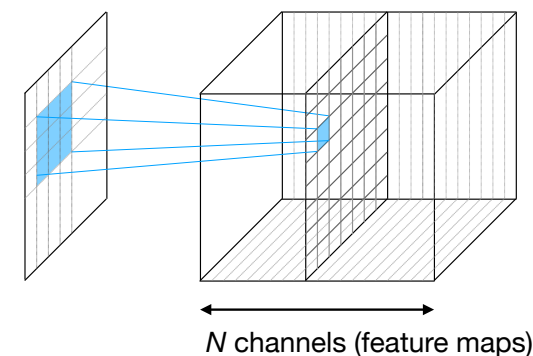
Example: $T_{\text{LSI}}\{f\} = h * f$ with $h = T_{\text{LSI}}\{\delta[\cdot]\}$ (discrete convolution)



Convolutional layer

- Patch extraction operator: $\ell_2(\mathbb{Z}^d) \times \mathbb{Z}^d \rightarrow \mathbb{R}^M$ with $M = \#W$

$$\text{Vect}(f[\cdot], \mathbf{k}) = \mathbf{f}_W[\mathbf{k}] = (f[\mathbf{k} - \mathbf{k}_0])_{\mathbf{k}_0 \in W}$$



- Convolution layer : $\ell_2(\mathbb{Z}^d) \rightarrow \ell_2^N(\mathbb{Z}^d)$ with N channels = feature maps

Shared operator $T_{\text{patch}} : \mathbb{R}^M \rightarrow \mathbb{R}^N$

$$T_{\text{patch}}(\mathbf{f}_W[\mathbf{k}]) = \boldsymbol{\sigma}(\mathbf{W}\mathbf{f}_W[\mathbf{k}]) \quad \text{where} \quad \mathbf{W} = \begin{pmatrix} \mathbf{w}_1^\top \\ \vdots \\ \mathbf{w}_N^\top \end{pmatrix} \quad \text{and} \quad \boldsymbol{\sigma} = \begin{pmatrix} \sigma_1 : \mathbb{R} \rightarrow \mathbb{R} \\ \vdots \\ \sigma_N : \mathbb{R} \rightarrow \mathbb{R} \end{pmatrix}$$

convolution masks

pointwise nonlinearities

Typically: $\sigma_n(x) = \text{ReLU}(x - b_n)$

Implementation of N channel convolution layer:

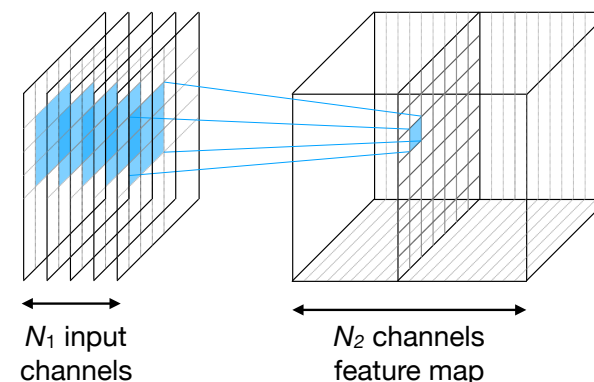
$$\begin{pmatrix} \sigma_1((h_1 * f)[\mathbf{k}]) \\ \vdots \\ \sigma_N((h_N * f)[\mathbf{k}]) \end{pmatrix} = T_{\text{patch}}(\mathbf{f}_W[\mathbf{k}]) = \boldsymbol{\sigma}(\mathbf{W}\mathbf{f}_W[\mathbf{k}])$$

Vector-valued convolutional layer (tensor)

Input feature map: $\mathbf{f}[\cdot] = (f_1[\cdot], \dots, f_{N_1}[\cdot])$

- Tensor patch extraction: $\ell_2^{N_1}(\mathbb{Z}^d) \times \mathbb{Z}^d \rightarrow \mathbb{R}^{M \times N_1}$ with $M = \#W$

$$\text{Tensor}(\mathbf{f}[\cdot], \mathbf{k}) = \mathbf{F}[\mathbf{k}] = (f_i[\mathbf{k} - \mathbf{k}_0])_{\mathbf{k}_0 \in W, i \in \{1, \dots, N_1\}}$$



- Vector-valued convolution layer : $\ell_2^{N_1}(\mathbb{Z}^d) \rightarrow \ell_2^{N_2}(\mathbb{Z}^d)$

with $\mathbf{F}[\mathbf{k}], \mathbf{w}_n \in \mathbb{R}^{(M \times N_1)}$

Shared operator $T_{\text{tensor}} : \mathbb{R}^{(M \times N_1)} \rightarrow \mathbb{R}^{N_2}$ convolution tensors pointwise nonlinearities

$$T_{\text{tensor}}(\mathbf{F}[\mathbf{k}]) = \boldsymbol{\sigma}(\mathbf{W}\mathbf{F}[\mathbf{k}]) \quad \text{where} \quad \mathbf{W} = \begin{pmatrix} \mathbf{w}_1^\top \\ \vdots \\ \mathbf{w}_{N_2}^\top \end{pmatrix} \quad \text{and} \quad \boldsymbol{\sigma} = \begin{pmatrix} \sigma_1 : \mathbb{R} \rightarrow \mathbb{R} \\ \vdots \\ \sigma_{N_2} : \mathbb{R} \rightarrow \mathbb{R} \end{pmatrix}$$

$$= \boldsymbol{\sigma}\left((\mathbf{H} * \mathbf{f})[\mathbf{k}]\right) \quad \text{with} \quad \mathbf{H}[\cdot] = N_2 \times N_1 \text{ array of filters}$$

Number of parameters: $(M \times N_1) \times N_2$

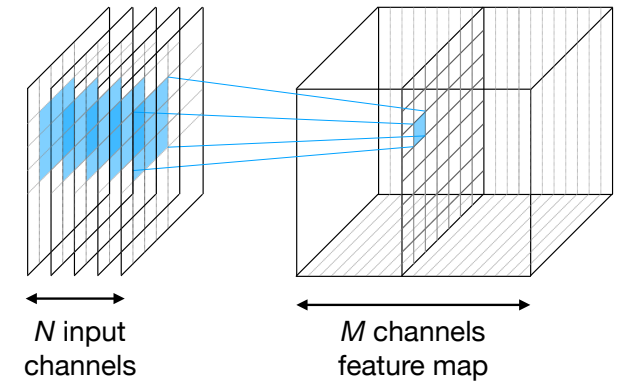
Convolution of vector-valued images

$$\text{Input feature map: } \mathbf{f}[\cdot] = \begin{pmatrix} f_1[\cdot] \\ \vdots \\ f_N[\cdot] \end{pmatrix}$$

$$\text{Matrix-valued filter: } \mathbf{H}[\cdot] = \begin{pmatrix} h_{1,1}[\cdot] & h_{1,2}[\cdot] & \dots & h_{1,N}[\cdot] \\ \vdots & \vdots & & \vdots \\ h_{M,1}[\cdot] & h_{M,2}[\cdot] & \dots & h_{M,N}[\cdot] \end{pmatrix}$$

- Vector-valued filterbank: $\ell_2^N(\mathbb{Z}^d) \rightarrow \ell_2^M(\mathbb{Z}^d)$

$$(\mathbf{H} * \mathbf{f})[\cdot] = \begin{pmatrix} (h_{1,1} * f_1)[\cdot] + (h_{1,2} * f_2)[\cdot] + \dots + (h_{1,N} * f_N)[\cdot] \\ \vdots \\ (h_{M,1} * f_1)[\cdot] + (h_{M,2} * f_2)[\cdot] + \dots + (h_{M,N} * f_N)[\cdot] \end{pmatrix}$$



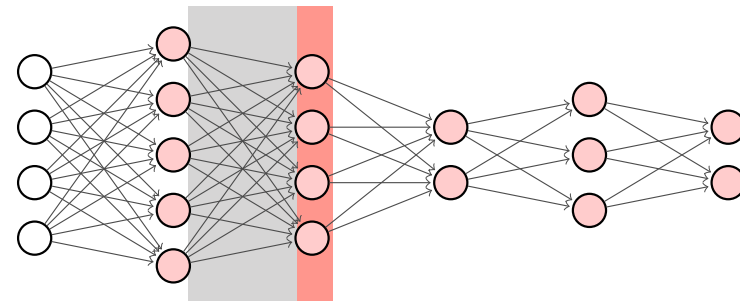
Composition

Series of (nonlinear) operators $T_\ell : \mathcal{X}_\ell \rightarrow \mathcal{Y}_\ell$ where \mathcal{X}_ℓ and \mathcal{Y}_ℓ are complete normed vector spaces.

Hypothesis: T_ℓ and $T_{\ell-1}$ have compatible domain and range; i.e., $\mathcal{X}_\ell = \mathcal{Y}_{\ell-1}$

- Composition of $T_1 : \mathcal{X}_1 \rightarrow \mathcal{X}_2$ and $T_2 : \mathcal{X}_2 \rightarrow \mathcal{X}_3$

$$T_2 \circ T_1 : \mathcal{X}_1 \rightarrow \mathcal{X}_3 \quad \text{with} \quad T_2 \circ T_1\{x\} = T_2\{T_1\{x\}\}$$



- Deep neural network

$$\mathbf{f}_{\text{deep}}(\mathbf{x}) = (\sigma_L \circ A_L \circ \sigma_{L-1} \circ \cdots \circ \sigma_2 \circ A_2 \circ \sigma_1 \circ A_1)(\mathbf{x})$$

- Affine layers $A_\ell : \mathbb{R}^{N_{\ell-1}} \rightarrow \mathbb{R}^{N_\ell}$

$$A_\ell(\mathbf{x}) = \mathbf{W}_\ell \mathbf{x} + \mathbf{b}_\ell \quad \text{with trainable weight matrix } \mathbf{W}_\ell \in \mathbb{R}^{N_\ell \times N_{\ell-1}} \text{ and bias } \mathbf{b}_\ell \in \mathbb{R}^{N_\ell}$$

- Pointwise nonlinearities $\sigma_\ell : \mathbb{R}^{N_\ell} \rightarrow \mathbb{R}^{N_\ell}$

$$\sigma_\ell(\mathbf{x}) = (\sigma(x_1), \dots, \sigma(x_{N_\ell})) \quad \text{with common activation function } \sigma : \mathbb{R} \rightarrow \mathbb{R}$$

Composition: Properties

- Preservation of linearity (affiness)

$$A_1 : \mathbf{x} \mapsto \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1 \text{ and } A_2 : \mathbf{y} \mapsto \mathbf{W}_2 \mathbf{y} + \mathbf{b}_2 \quad \Rightarrow \quad A_2 \circ A_1 : \mathbf{x} \mapsto (\mathbf{W}_2 \mathbf{W}_1) \mathbf{x} + (\mathbf{b}_2 + \mathbf{W}_2 \mathbf{b}_1)$$

- Preservation of convolutional structure

T_1 and T_2 are LSI with impulse responses $h_1, h_2 \in \ell_1(\mathbb{Z}^d)$

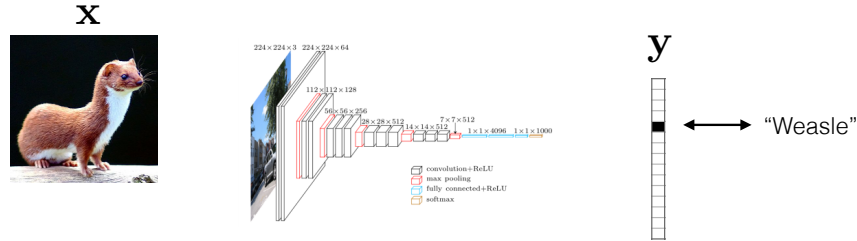
$$\Rightarrow \quad T_2 \circ T_1 : f \mapsto h * f \quad \text{with} \quad h = h_2 * h_1 \in \ell_1(\mathbb{Z}^d)$$

Application: Construction of larger receptive fields

- Preservation of continuity

The composition of two continuous functions is continuous

Pooling



■ Down-sampling: $\downarrow_m \{f\}[\mathbf{k}] = f[m\mathbf{k}]$

■ Max pooling: $\mathbb{R}^N \rightarrow \mathbb{R}$

$$\mathbf{u} \mapsto \max(\mathbf{u}) = \max(u_1, \dots, u_N)$$

■ Softmax: $\mathbb{R}^N \rightarrow \mathbb{R}^N$ (transforms output of CNN in “probabilities”)

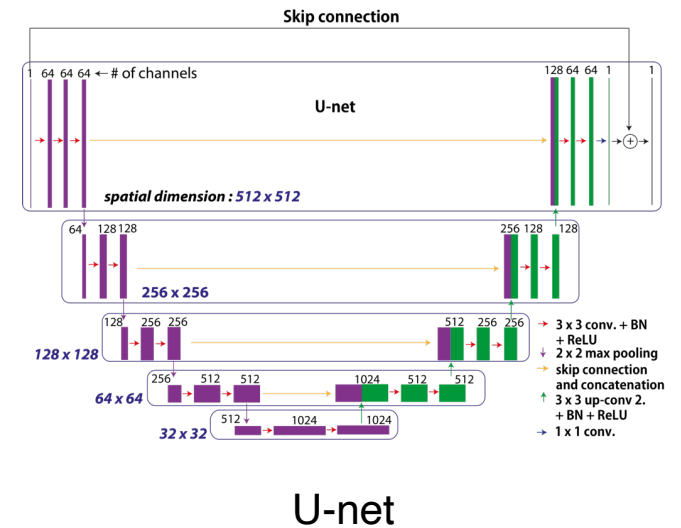
$$\mathbf{u} \mapsto (p_k = \frac{\exp(u_k)}{\sum_{n=1}^N \exp(u_n)})$$

■ Up-sampling

$$\uparrow_m \{f\}[\mathbf{k}] = \begin{cases} f[\mathbf{n}], & \mathbf{k} = m\mathbf{n} \\ 0, & \text{otherwise} \end{cases}$$

■ Up-sampling with repetition

$$\text{UP}_m \{f\}[\mathbf{k}] = f[\mathbf{k}/m] = (\uparrow_m \{f\} * \mathbb{1}_{[0, m-1]^d})[\mathbf{k}]$$



U-net

Continuity requirements

Generic operator $T : \mathcal{X} \rightarrow \mathcal{Y}$ where $(\mathcal{X}, \|\cdot\|_{\mathcal{X}})$ and $(\mathcal{Y}, \|\cdot\|_{\mathcal{Y}})$ are complete normed vector spaces.

Definition

The operator $T : \mathcal{X} \rightarrow \mathcal{Y}$ is **continuous** if $\lim_{i \rightarrow \infty} T\{x_i\} = T\{\lim_{i \rightarrow \infty} x_i\} = T\{x\}$ for any sequence $(x_i)_{i \in \mathbb{N}}$ that is converging in \mathcal{X} with $\lim_{i \rightarrow \infty} x_i = x$

Definition

The operator $T : \mathcal{X} \rightarrow \mathcal{Y}$ is **Lipschitz continuous** if there exists a constant $L > 0$ such that $\|T\{x_1\} - T\{x_2\}\|_{\mathcal{Y}} \leq L\|x_1 - x_2\|_{\mathcal{X}}$ for any $x_1, x_2 \in \mathcal{X}$.

■ Lipschitz constant

$\text{Lip}(T) = L$ where L is the smallest constant such that the Lipschitz inequality holds.

To avoid instabilities, **all modules** of a CNN should be **Lipschitz continuous**.

This implies that they are a.e. **differentiable**, which is desirable for training with backpropagation.

Counterexample: TLU networks are discontinuous, and therefore very hard to train unless the architecture is shallow (perceptron).

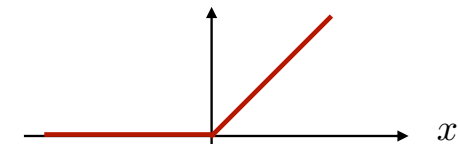
Lipschitz constant of primary modules

■ Pointwise nonlinearity

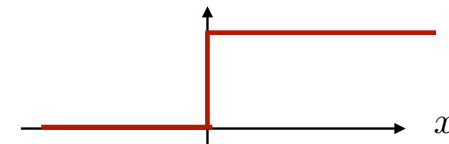
$\sigma : \mathbb{R} \rightarrow \mathbb{R}$ where σ is differentiable

$$\text{Lip}(\sigma) = \sup_{x \in \mathbb{R}} \left| \frac{d\sigma(x)}{dx} \right| = \|\sigma'\|_{L_\infty} \quad (\text{cf. Mean Value Theorem})$$

Example: $\text{Lip}(\text{ReLU}) = \sup_{x \in \mathbb{R}} |u(x)| = 1$



Heaviside



■ LSI operator (convolution channel)

$T_{\text{LSI}}\{f\} = h * f$ where $h \in \ell_1(\mathbb{Z}^d)$

$$\text{Lip}(T_{\text{LSI}}) = H_{\max} = \sup_{\omega \in [0, \pi]^d} |H(e^{j\omega})| \leq \|h\|_{\ell_1}$$

Justification (Parseval)

$$\begin{aligned} \|h * f - h * g\|_{\ell_2}^2 &= \|h * (f - g)\|_{\ell_2}^2 = \frac{1}{(2\pi)^d} \int_{[-\pi, \pi]^d} |H(e^{j\omega})|^2 |F(e^{j\omega}) - G(e^{j\omega})|^2 d\omega \\ &\leq \frac{H_{\max}^2}{(2\pi)^d} \int_{[-\pi, \pi]^d} |F(e^{j\omega}) - G(e^{j\omega})|^2 d\omega = H_{\max}^2 \|f - g\|_{\ell_2}^2 \end{aligned}$$

Lipschitz constant of primary modules (cont'd)

■ Linear (resp. affine) transform

$$T_{\text{lin}} : \mathbb{R}^M \rightarrow \mathbb{R}^N \quad \text{with} \quad \mathbf{x} \mapsto \mathbf{A}\mathbf{x} \text{ (linear)}$$

$$\text{or} \quad \mathbf{x} \mapsto \mathbf{A}\mathbf{x} + \mathbf{b} \text{ (affine)} \quad \text{where} \quad \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{b} \in \mathbb{R}^M$$

$$\text{Lip}(T_{\text{lin}}) = \sup_{\|\mathbf{x}\|_2 \leq 1} \|\mathbf{A}\mathbf{x}\|_2 = \rho(\mathbf{A}) \quad (\text{spectral norm} = \text{largest singular value of } \mathbf{A})$$

■ Imposing Lip-1 layer by spectral normalization

$$T_{\text{normal}}\{\mathbf{x}\} = \frac{1}{\rho(\mathbf{A})} \mathbf{A}\mathbf{x}$$

Estimation by power method: For $k = 0, \dots, K$

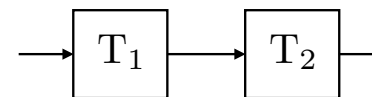
$$\mathbf{u}_k = \frac{1}{\|\mathbf{A}^\top(\mathbf{A}\mathbf{u}_k)\|_2} \mathbf{A}^\top(\mathbf{A}\mathbf{u}_k)$$

Upon convergence, $\mathbf{u} = \lim_k \mathbf{u}_k$ is the dominant eigenvector of $\mathbf{A}^\top \mathbf{A}$.

$$\text{Finally, } \rho(\mathbf{A}) = \sqrt{(\mathbf{A}\mathbf{u})^\top(\mathbf{A}\mathbf{u})}$$

Stability and combination of modules

■ Composition



$$\text{Lip}(T_1) = L_1 \ \& \ \text{Lip}(T_2) = L_2 \quad \Rightarrow \quad \text{Lip}(T_2 \circ T_1) \leq L_2 L_1$$

$$\forall f, g \in \mathcal{X}_1 : \|\text{T}_2\{\text{T}_1\{f\}\} - \text{T}_2\{\text{T}_1\{g\}\}\|_{\mathcal{Y}_2} \leq L_2 \|\text{T}_1\{f\} - \text{T}_1\{g\}\|_{\mathcal{Y}_1 = \mathcal{X}_2} \leq L_2 L_1 \|f - g\|_{\mathcal{X}_1}$$

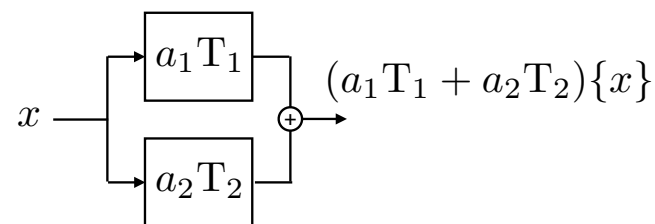
Deep neural network with Lip-1 activations (e.g., ReLU) :

$$\mathbf{f}_{\text{deep}}(\mathbf{x}) = (\sigma_L \circ A_L \circ \sigma_{L-1} \circ \dots \circ \sigma_2 \circ A_2 \circ \sigma_1 \circ A_1)(\mathbf{x}) \quad \Rightarrow \quad \text{Lip}(\mathbf{f}_{\text{deep}}) = \prod_{\ell=1}^L \left(\overbrace{\text{Lip}(\sigma_\ell)}^1 \text{Lip}(A_\ell) \right)$$

➔ Deeper networks tend to be less stable

■ Linear combination

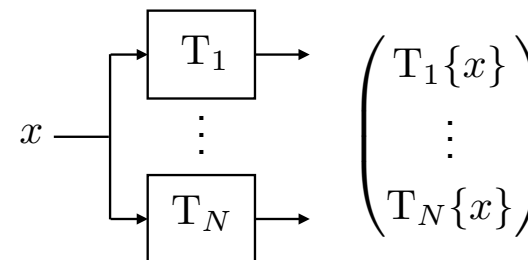
$$\text{Lip}\left(\sum_{i=1}^I a_i T_i\right) \leq \sum_{i=1}^I |a_i| \text{Lip}(T_i) \quad (\text{by triangle inequality})$$



■ Parallel feature maps

$$\mathbf{T} = \begin{pmatrix} T_1 \\ \vdots \\ T_N \end{pmatrix} : \ell_2(\mathbb{Z}^d) \rightarrow \ell_2^N(\mathbb{Z}^d)$$

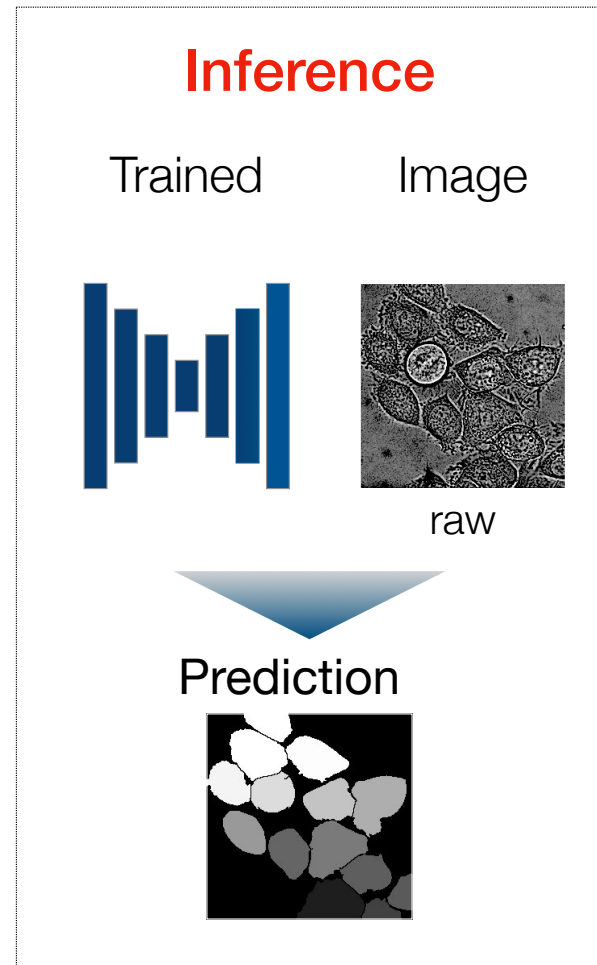
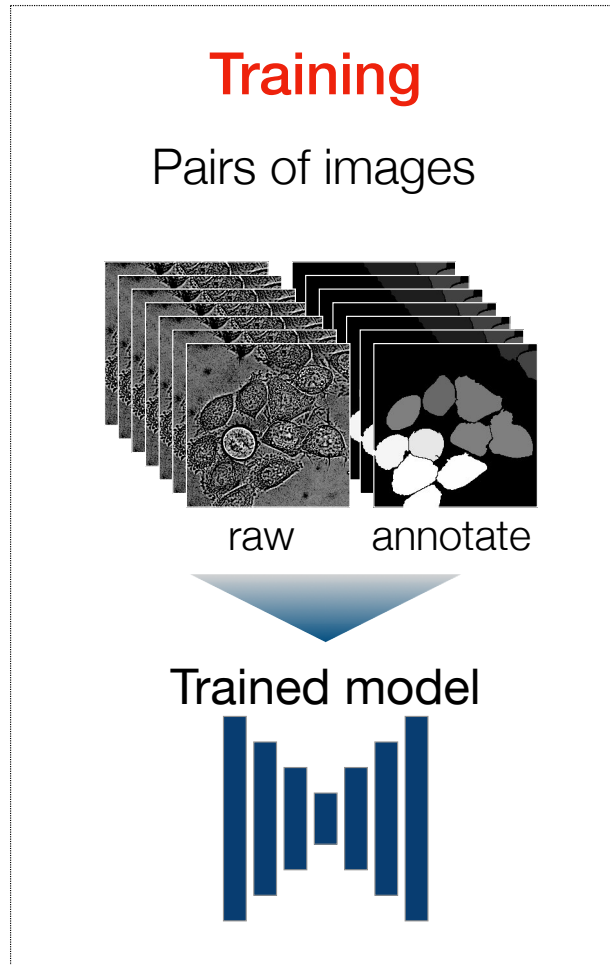
$$\text{Lip}(\mathbf{T}) \leq \sqrt{L_1^2 + \dots + L_N^2} \leq L_1 + \dots + L_N$$



CNNs in Practice

- Deep learning pipeline
- Denoising
- Segmentation

Deep Learning Pipeline



Denoising: “Handcrafted” ancestor of Resnet



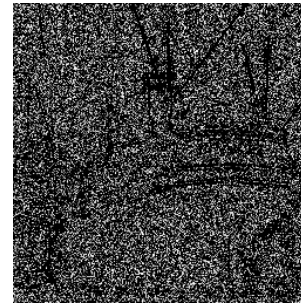
Ground-truth



LP



HP



HP-clip

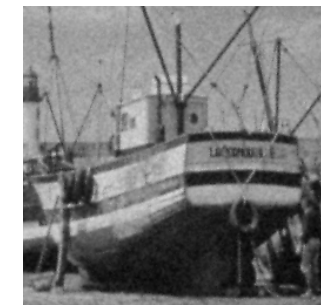
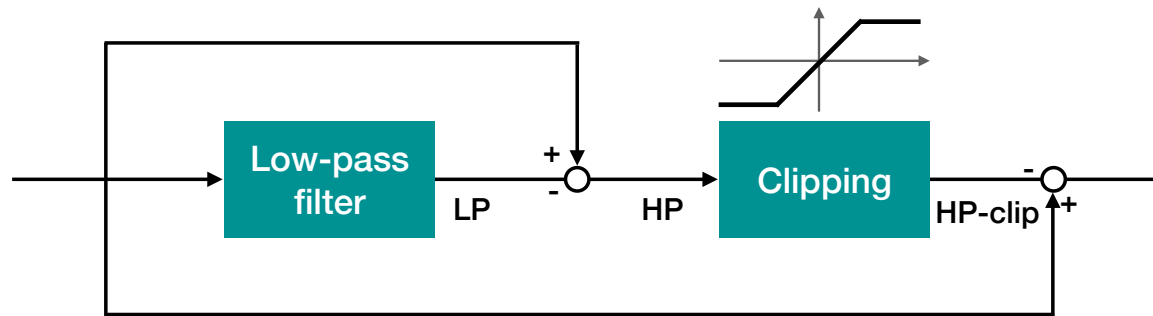
noise estimator

	Input	LP	Denoising
RMSE	17.30	18.19	12.08
PSNR	21.89 dB	21.45 dB	25.01dB

with respect to ground truth

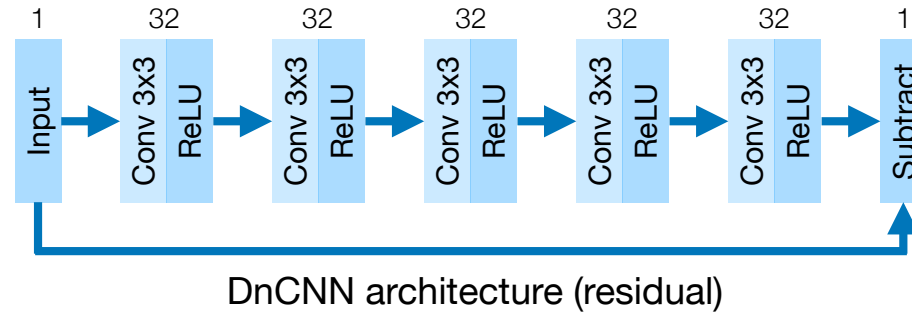


Input
noisy image



Output
denoised image

Deep CNN for residual image denoising



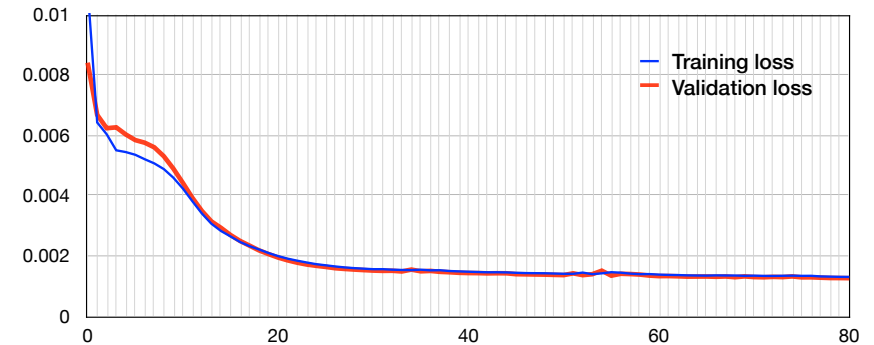
Datasets

- 80 natural graylevel images 384×384, artificially degraded by adding noise



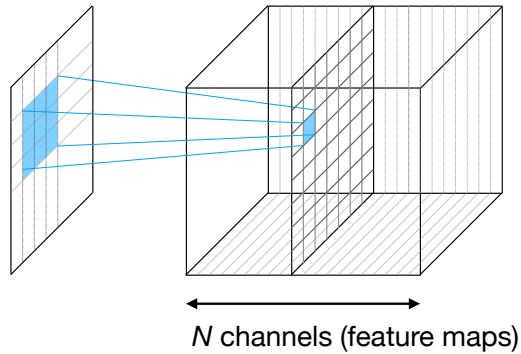
Training CNN

- 20 images (25%) in validation set
- **37'601 parameters** (weights and biases)
- Loss function is the MSE
- 80 epochs, batch_size = 16, lr = 0.001
- Global normalization: $\mu=0$, $\text{Var}=1$



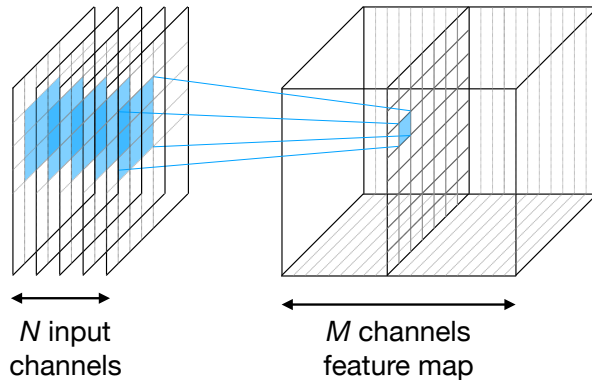
DnCNN architecture

First conv layer:



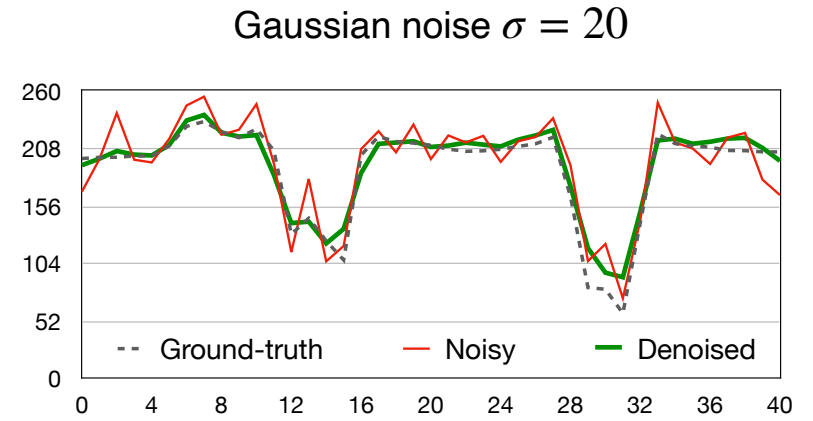
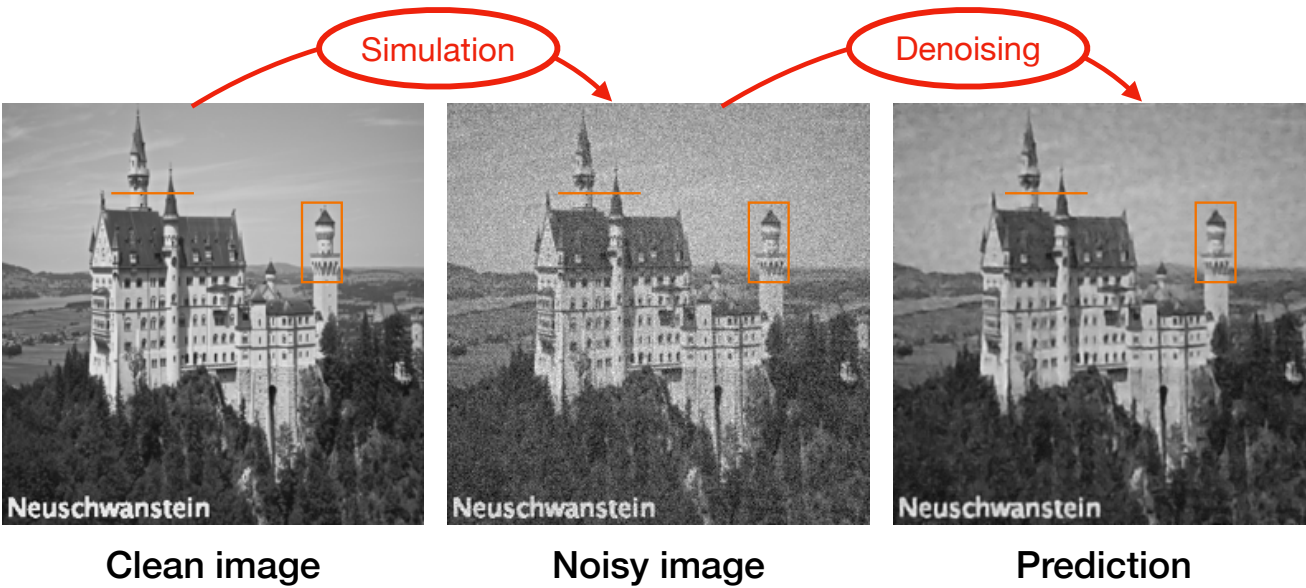
Number of learnable parameters: $(3 \times 3) + 1$ (bias) per output channel

Internal conv layers:



Number of learnable parameters: $(3 \times 3) \times 32 + 1$ (bias) per output channel

Layer	Shape	Parameters
Input	384 x 384 x 1	0
Conv2D 3x3 + b, ReLU	384 x 384 x 32	320
Conv2D 3x3 + b, ReLU	384 x 384 x 32	9248
Conv2D 3x3 + b, ReLU	384 x 384 x 32	9248
Conv2D 3x3 + b, ReLU	384 x 384 x 32	9248
Conv2D 3x3 + b, ReLU	384 x 384 x 32	9248
Conv2D 3x3 + b, ReLU	384 x 384 x 1	289
Subtract	384 x 384 x 1	0
Total		37601





Neuschwanstein

Clean image
Ground-truth



Neuschwanstein

Noisy image 15.36 dB
Additive Gaussian noise



Neuschwanstein

High-pass filter
Soft Clipping
19.97 dB



Neuschwanstein

Resnet (100 epochs)
3 layers, 16 channels
20.91 dB



Neuschwanstein

Gaussian filter
 $\sigma = 1$
18.23 dB



Neuschwanstein

Median filter
radius = 3
18.45 dB



Neuschwanstein

Resnet (200 epochs)
5 layers, 32 channels
21.01 dB



Neuschwanstein

Unet (100 epochs)
3 Pooling steps, 32 channels
20.99 dB

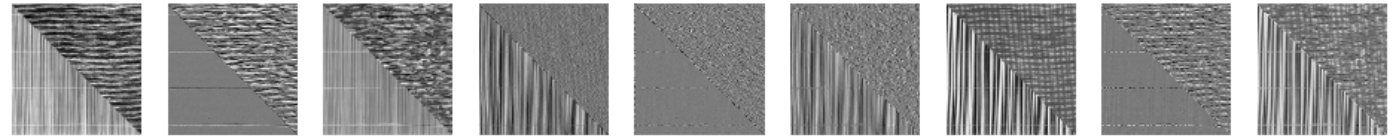
Simulation

Segmentation: “Handcrafted” texture discriminator

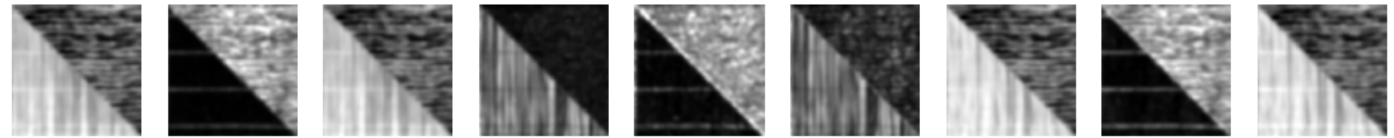
DCT filters (3x3)



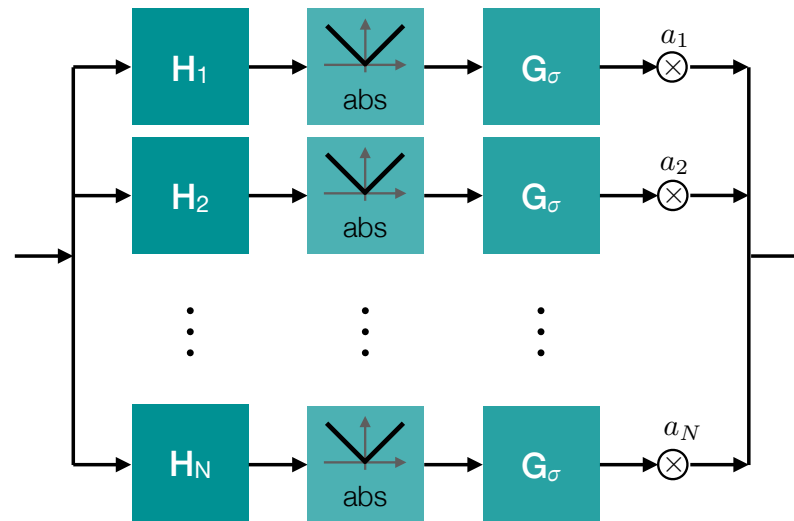
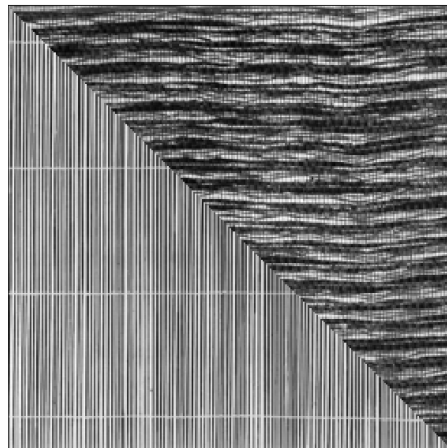
Output of DCT filterbank



Feature maps

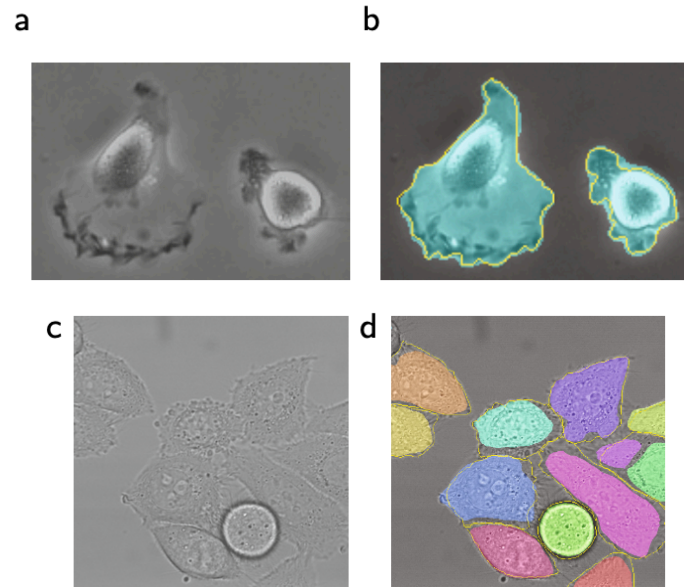
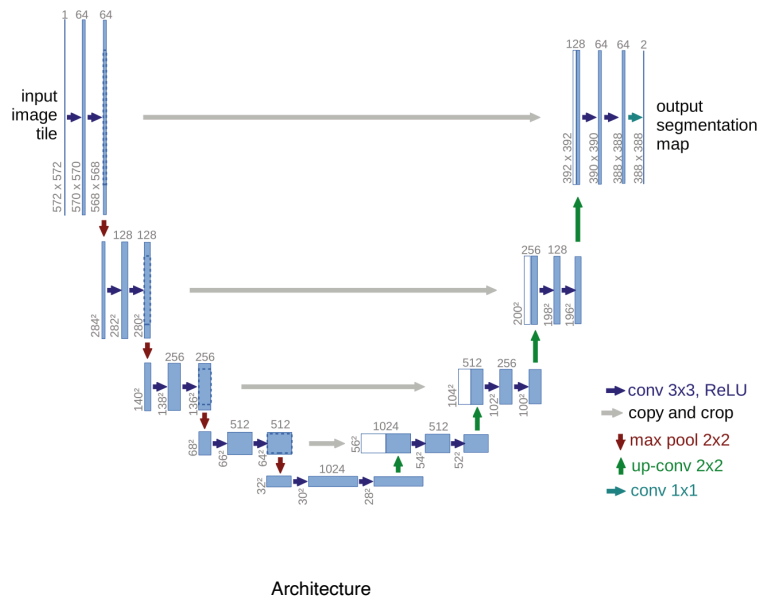


N filters nonlinearity aggregation



Popular CNN architecture for image segmentation

U-net introduced by Ronneberger in 2015 for biomedical image segmentation

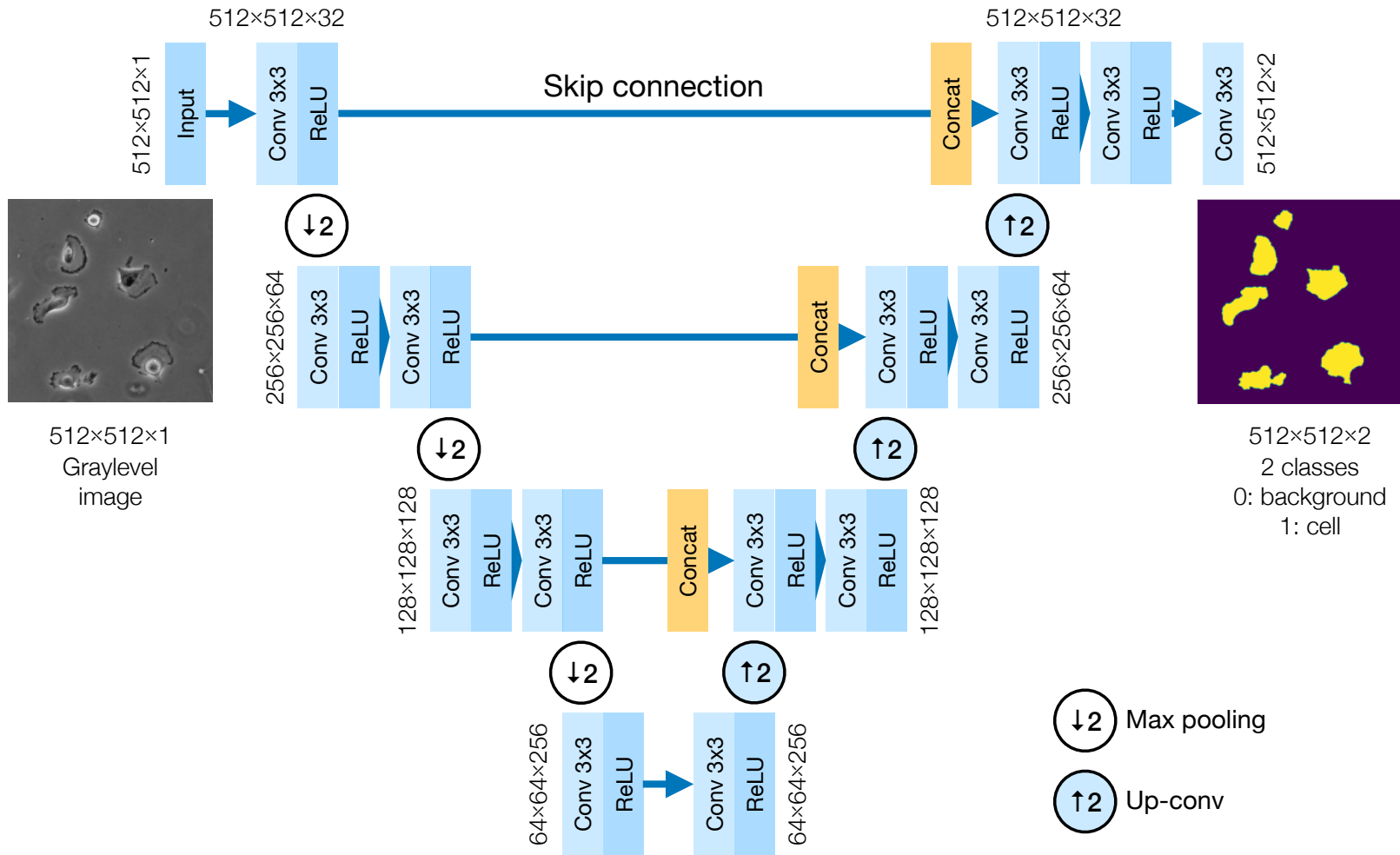


Winner of ISBI 2015 cell tracking challenge

1. <https://arxiv.org/abs/1505.04597>

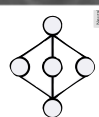
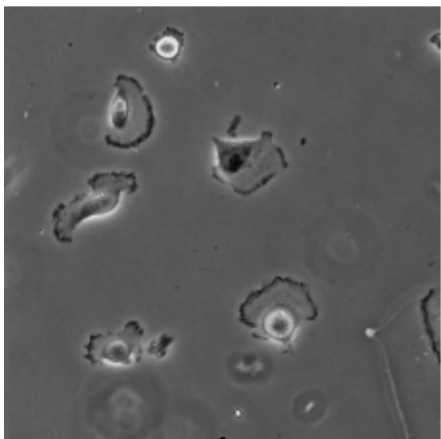
U-Net Architecture

For "semantic" segmentation

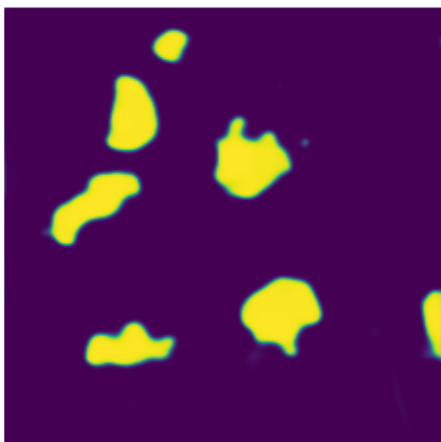


Layer	Shape	Params
Input	512x512x1	0
Conv2D 3x3	512x512x32	320
Conv2D 3x3, Activation	512x512x32	9248
MaxPooling2D	256x256x32	0
Conv2D 3x3, Activation	256x256x64	18496
Conv2D 3x3, Activation	256x256x64	36928
MaxPooling2D	128x128x64	0
Conv2D 3x3, Activation	128x128x128	73856
Conv2D 3x3, Activation	128x128x128	147584
MaxPooling2D	64x64x128	0
Conv2D 3x3, Activation	64x64x256	295168
Conv2D 3x3, Activation	64x64x256	590080
Up-conv	128x128x128	131200
Concatenate	128x128x256	0
Conv2D 3x3, Activation	128x128x128	295040
Conv2D 3x3, Activation	128x128x128	147584
Up-conv	256x256x64	32832
Concatenate	256x256x128	0
Conv2D 3x3, Activation	256x256x64	73792
Activation	256x256x64	0
Conv2D 3x3, Activation	256x256x64	36928
Up-conv	512x512x32	8224
Concatenate	512x512x64	0
Conv2D 3x3, Activation	512x512x32	18464
Conv2D 3x3, Activation	512x512x32	9248
Conv2D 3x3	512x512x2	66
Total		1'925'058

Raw input image



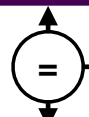
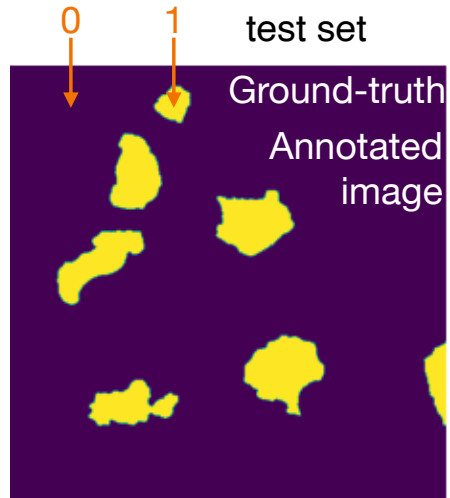
Prediction



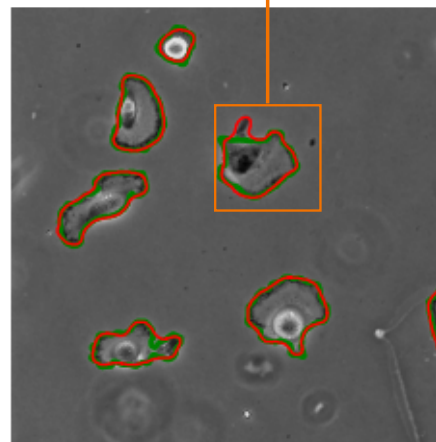
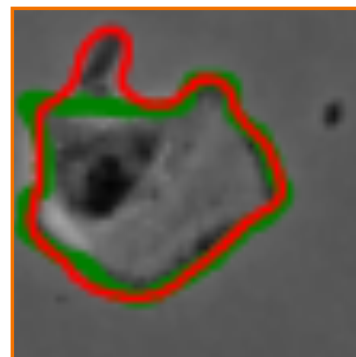
Probability of class 1

$p > 0.5$

Prediction of class 1



IoU = 0.91 (Jaccard)



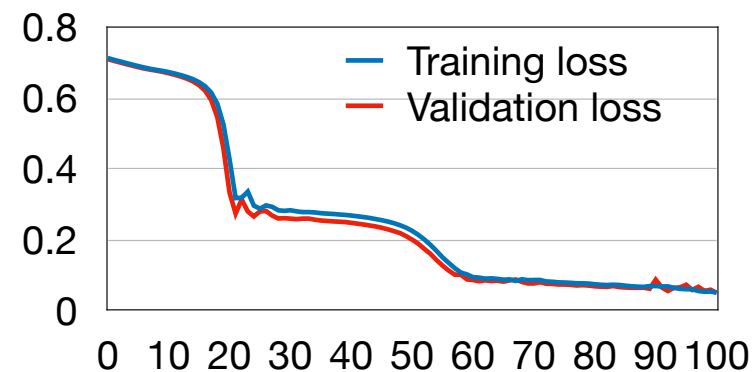
Results of segmentation
GT (green) Pred (red)

Datasets

- 66 graylevel images (1 channel)
- 66 manually annotated images, 2 classes

Training CNN

- 66 annotated images / 25% for validation
- **1'925'058 parameters** (weights and bias)
- Loss function is the binary cross entropy
- 100 epochs, batch_size = 16, lr = 0.001



CONCLUSION: The Pros and Cons of CNNs

- Advantages of (deep) CNNs
 - Are sufficiently flexible to implement most image-processing tasks
 - Can benefit from hardware acceleration (GPU)
 - Performance of **deep CNNs** is often **spectacular**
 - Once trained, they are very fast to deploy (inference)
- Downsides of CNNs
 - Require **huge amounts** of **data** and **computation** for training
 - While the individual components are simple, the global behaviour is poorly understood ⇒ **lack of guarantees**
 - Training and fine-tuning is fastidious “*Graduate-student descent*”
 - No free lunch: Trained CNNs are very task/data specific
 - Can **behave erratically** — lack of robustness, subject to adversarial attacks
- How the new trend benefits from the techniques of “traditional” IP
 - **Deeper understanding** of modules
 - Suggestion of **leaner** and **more robust** architectures

