

Basics of AI

And Some X-Ray or CT Examples

Marc Kachelrieß

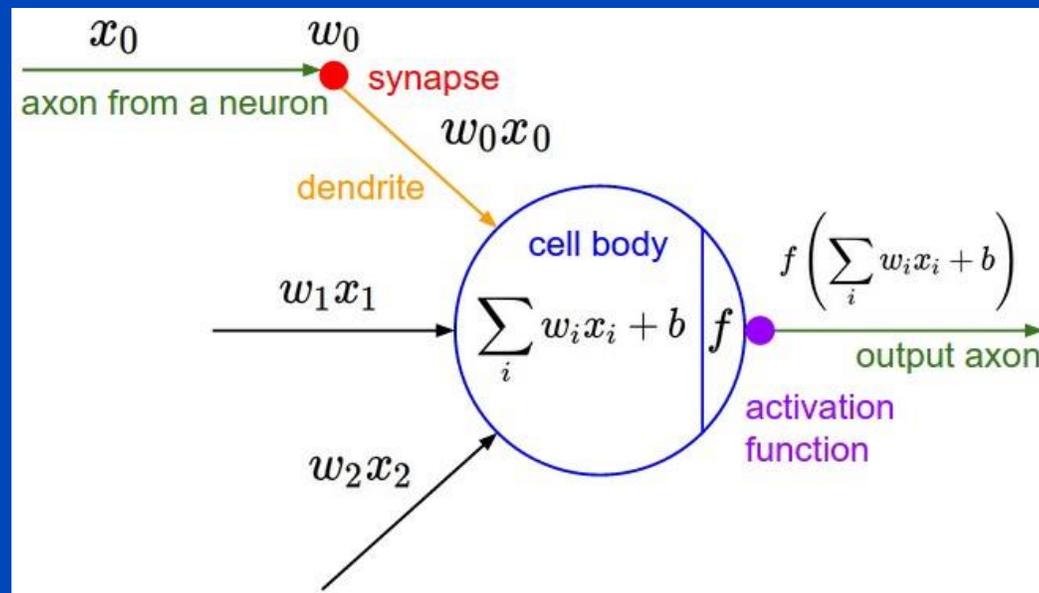
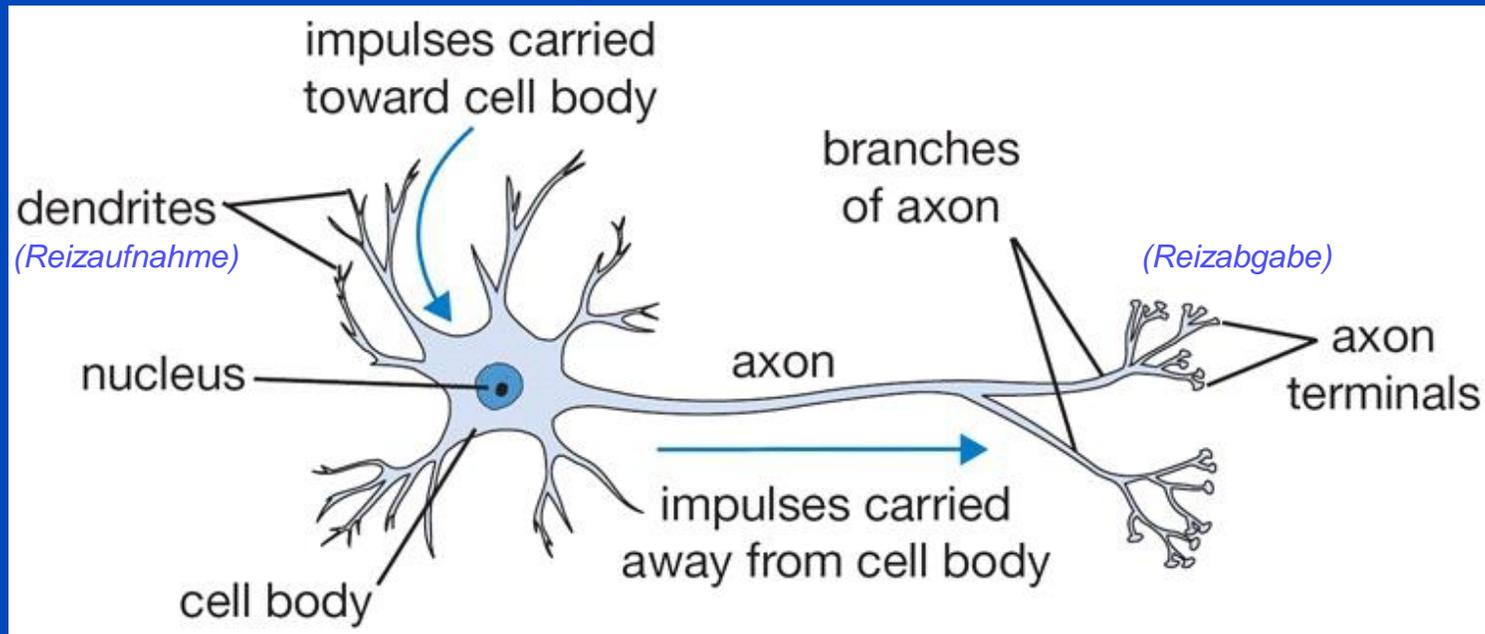
German Cancer Research Center (DKFZ)

Heidelberg, Germany

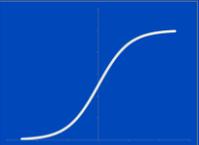
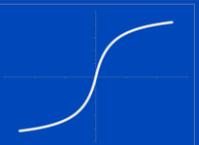
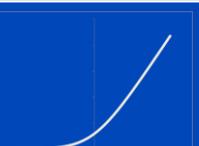
www.dkfz.de/ct

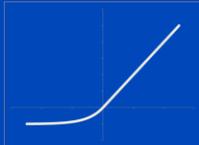
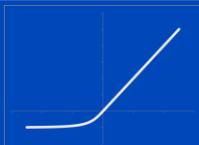
Nomenclature

- Iteration = Epoch
- Batch = Subset (randomly changing for each epoch)
- Loss function = Cost function
- Learning rate = η



Activation Functions

Function	Equation	Plot
Identity	$f(x) = x$	
Sigmoid	$f(x) = \frac{1}{1 + e^{-x}}$	
Hard sigmoid	$f(x) = \begin{cases} 0 & \text{for } x < -\alpha \\ \frac{\alpha+x}{2\alpha} & \text{for } -\alpha \leq x < \alpha \\ 1 & \text{for } x \geq \alpha \end{cases}$	
Tanh	$f(x) = \frac{2}{1 + e^{-2x}} - 1$	
Softsign	$f(x) = \frac{x}{1 + x }$	
Softplus	$f(x) = \log(1 + \exp x)$	

Function	Equation	Plot
ReLU	$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	
Leaky ReLU	$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	
ELU	$f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	
Inverse square root LU	$f(x) = \begin{cases} \frac{x}{\sqrt{1+\alpha x^2}} & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	

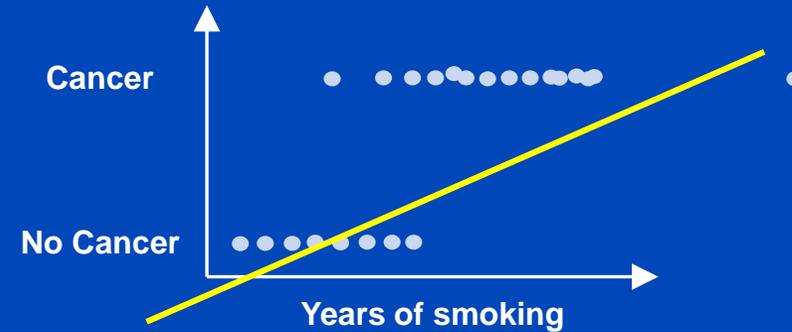
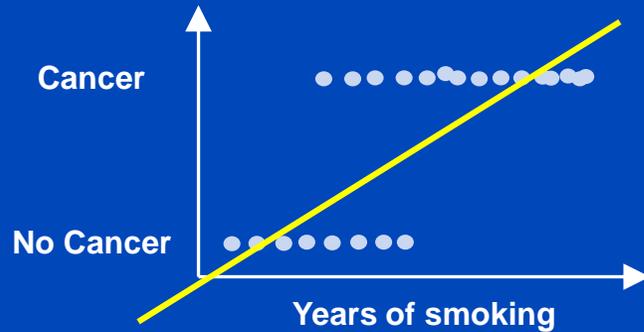
...

...

...

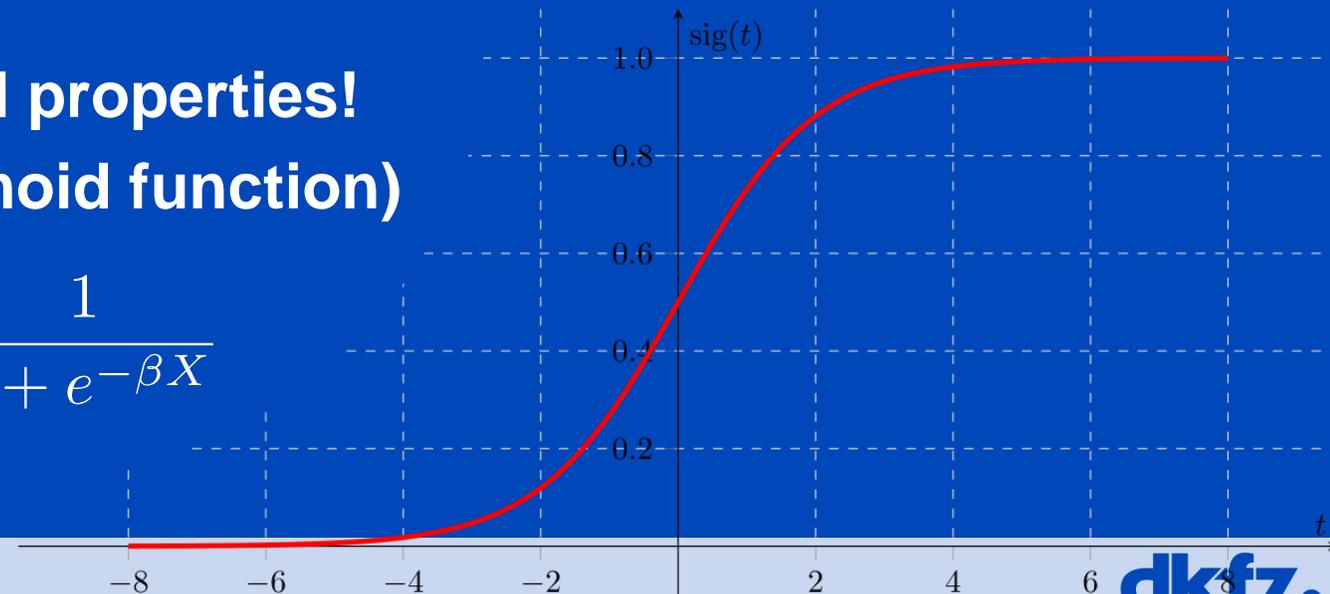
Binary Logistic Regression

- What if y is categoric, e.g. $y \in \{0, 1\}$?



- Linear regression has undesired properties!
- Use logistic model instead (sigmoid function)

$$p(x) = p(t(x)) = \frac{1}{1 + e^{-t(x)}} = \frac{1}{1 + e^{-\beta X}}$$



Loss Function

- The neural network coefficients (weights and biases) \mathbf{c} are chosen by minimizing a loss function (cost function)

$$\mathbf{c} = \arg \min_{\mathbf{c}} \sum_{n=1}^N L(\mathbf{c}, \mathbf{x}_n, \mathbf{y}_n)$$

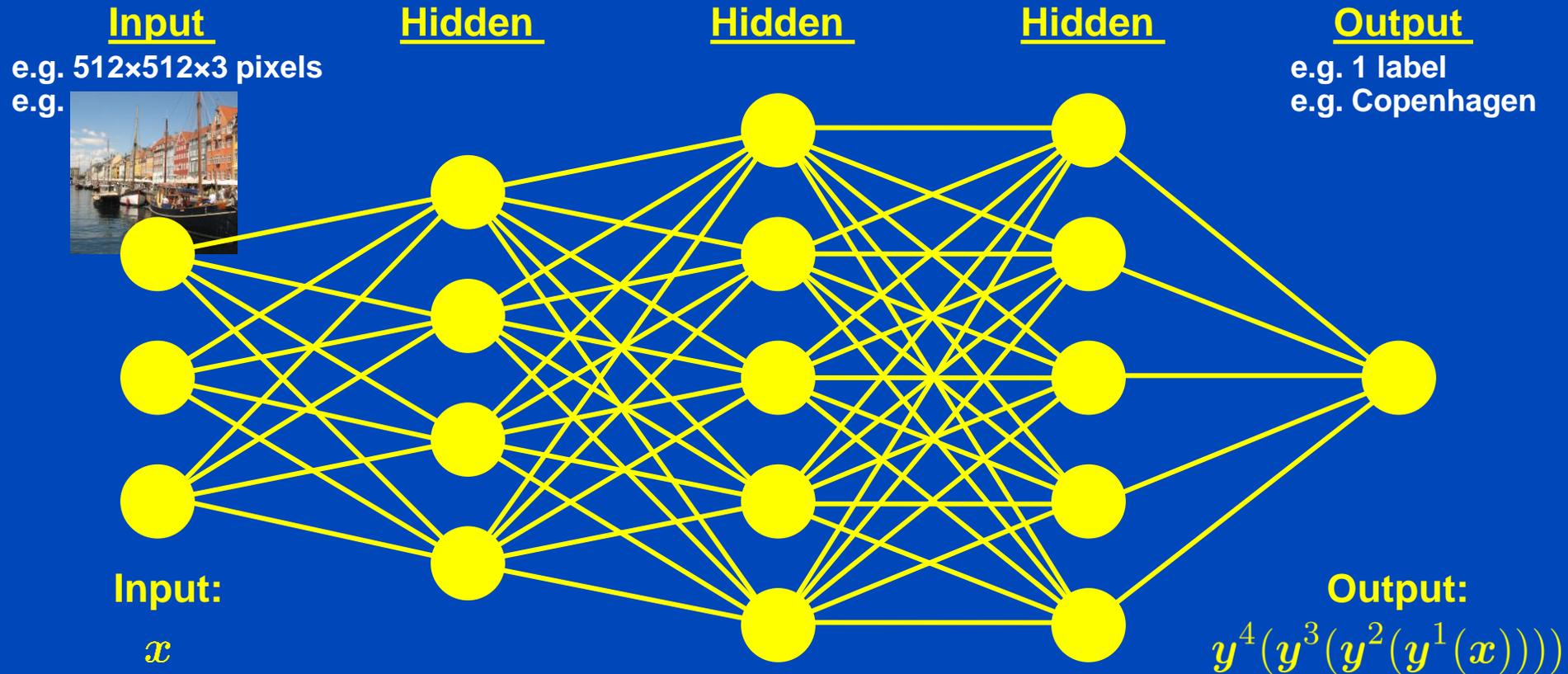
with \mathbf{x}_n being the training data input, $\mathbf{y}(\mathbf{c}, \mathbf{x}_n)$ being the network output, and \mathbf{y}_n being the so-called labels, i.e. the training target, and N being the number of training samples.

- An example for such a loss function is the MSE loss

$$L(\mathbf{c}, \mathbf{x}_n, \mathbf{y}_n) = \left(\mathbf{y}(\mathbf{c}, \mathbf{x}_n) - \mathbf{y}_n \right)^2$$

Fully-Connected Neural Network

- Each layer fully connects to previous layer
- Difficult to train (many parameters in W and b)
- Spatial relations not necessarily preserved



$y(x) = f(W \cdot x + b)$ with $f(x) = (f(x_1), f(x_2), \dots)$ point-wise scalar, e.g. $f(x) = x \vee 0 = \text{ReLU}$

Universal Approximation Theorem¹

- A fully-connected network with at least three layers is also called a multi layer perceptron (MLP).

If σ is continuous, bounded and nonconstant, then $\mathfrak{R}_M(\sigma)$ is dense in the subset \mathcal{C} of continuous and compact functions of \mathbb{R}^M .

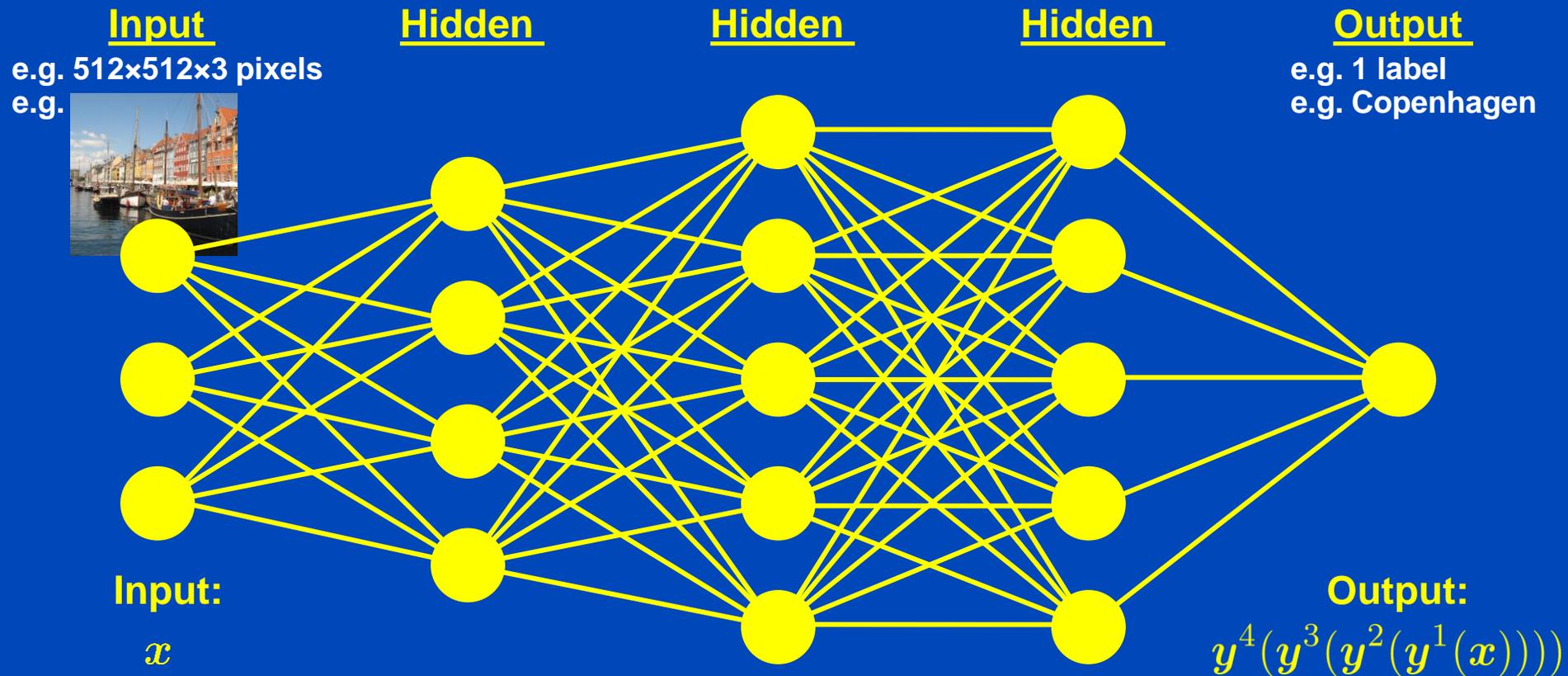
$$\begin{aligned}\mathfrak{R}_M(\sigma) &= \bigcup_N \mathfrak{R}_M^{(N)}(\sigma) \\ &= \bigcup_N \left\{ g : \mathbb{R}^M \rightarrow \mathbb{R} \mid \overbrace{g(x) = \sum_{j=0}^N w_{k,j}^{(3)} \sigma \left(\sum_{i=0}^M w_{j,i}^{(2)} x_i \right)}^{\text{3-layer MLP with final activation linear and } N \text{ hidden neurons}} \right\}\end{aligned}$$

- For any function $h(x) \in \mathcal{C} \subseteq \mathbb{R}^M$ we can find a function $g(x) \in \mathfrak{R}_M(\sigma)$ for which $\|h(x) - g(x)\|_p < \epsilon$.
- Any 3-layer MLP with appropriately chosen layer sizes and activation function, e.g. the sigmoid function, is a universal function approximator.
- **This theorem does not provide any insight into how to find the unknowns!**

¹Hornik, Kurt; Stinchcombe, Maxwell; White, Halbert (1989). Multilayer Feedforward Networks are Universal Approximators. Neural Networks. Vol. 2. Pergamon Press. pp. 359–366; Theorem 2

Fully-Connected Neural Network

- Each layer fully connects to previous layer
- Difficult to train (many parameters in W and b)
- Spatial relations not necessarily preserved



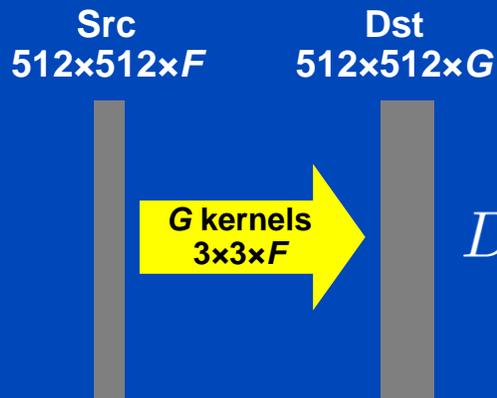
$y(x) = f(W \cdot x + b)$ with $f(x) = (f(x_1), f(x_2), \dots)$ point-wise scalar, e.g. $f(x) = x \vee 0 = \text{ReLU}$

Convolutional Neural Network (CNN)

- Replace dense W in $y(x) = f(W \cdot x + b)$ by a sparse matrix W with sparsity being of convolutional type (band diagonal of Toeplitz type).
- CNNs consist (mainly) of convolutional layers.
- Convolutional layers are not fully connected.
- Convolutional layers are small, say 3×3 , convolution kernels whose entries need to be found by training.
- CNNs preserve spatial relations to some extent.

$$W_{1D} = \begin{pmatrix} 1 & 8 & 0 & 0 & 0 & 0 \\ 2 & 1 & 8 & 0 & 0 & 0 \\ 0 & 2 & 1 & 8 & 0 & 0 \\ 0 & 0 & 2 & 1 & 8 & 0 \\ 0 & 0 & 0 & 2 & 1 & 8 \\ 0 & 0 & 0 & 0 & 2 & 1 \end{pmatrix}$$

Only three unknowns!



$$D_{i,j,g} = \sum_f S_{i,j,f} * K_{i,j,f}^g = \sum_{a,b,f} S_{i-a,j-b,f} K_{a,b,f}^g$$

Attention: No convolution in depth direction!

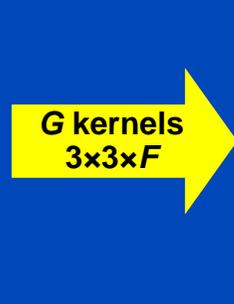
Here, a 2D example is shown. Conv layers also exist in 3D and higher dimensions.

Convolution Layers

- Input layer S
 - vector of size I with F features: $I \times F$
 - image of size I by J with F features: $I \times J \times F$
 - volume of size I by J by K with F features: $I \times J \times K \times F$
 - ...
- Convolution kernel K
 - G kernels of size $(2A+1) \times (2B+1) \times F$ with or without padding*
- Output layer D
 - same spatial dimensions as input layer*
 - G features (depth G)

Src
 $512 \times 512 \times F$

Dst
 $512 \times 512 \times G$

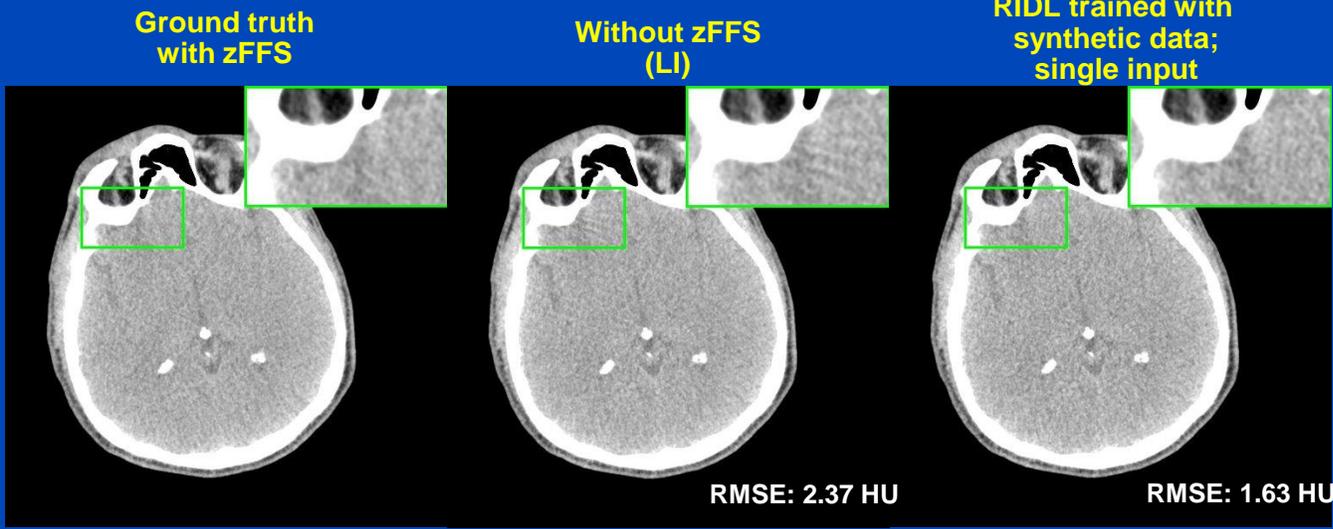
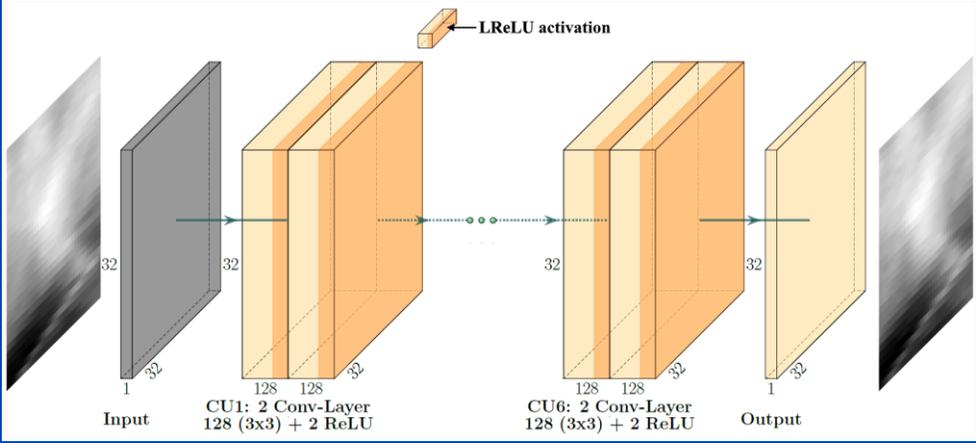
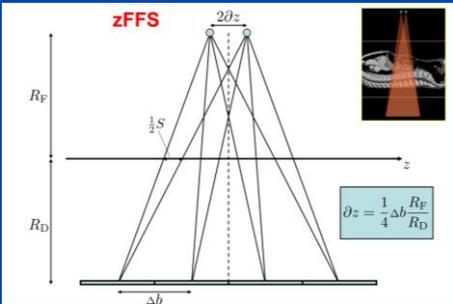


$$D_{i,j,g} = \sum_f S_{i,j,f} * K_{i,j,f}^g = \sum_{a,b,f} S_{i-a,j-b,f} K_{a,b,f}^g$$

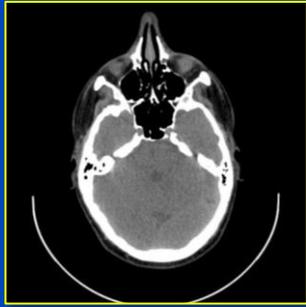
Attention: No convolution in depth direction!

*Convolution may include a stride (step size) > 1 . Similar to convolution with stride 1 followed by pooling.

Row Interpolation to Mimick zFFS



$C = 0 \text{ HU}, W = 200 \text{ HU}$



$C = 60 \text{ HU}, W = 360 \text{ HU}$



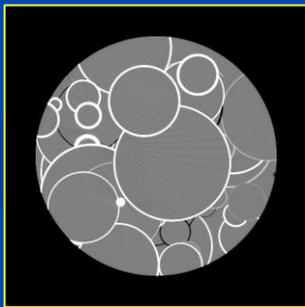
$C = -400 \text{ HU}, W = 1500 \text{ HU}$



$C = 60 \text{ HU}, W = 400 \text{ HU}$



Random synthetic phantom



$C = 0 \text{ HU}, W = 400 \text{ HU}$

Pooling

Downsampling

- **Input layer S**
 - image of size I by J with F features: $I \times J \times F$
 - ...
- **Pooling kernel**
 - pooling function, e.g. max, mean, stochastic, ...
 - size and strides
- **Output layer D**
 - reduced spatial size
 - same depth

1	1	1	3	2	3	1	2
2	3	0	3	1	9	6	9
1	8	0	4	0	8	9	9
1	1	2	3	9	2	3	1
0	5	1	3	2	1	1	3
1	1	1	1	0	0	1	1
2	5	0	7	1	9	7	9
2	0	0	8	2	4	0	1

2x2 stride 2x2
max pool

3	3	9	9
8	4	9	9
5	3	2	3
5	8	9	9

Src
64x64xF

Dst
32x32xF

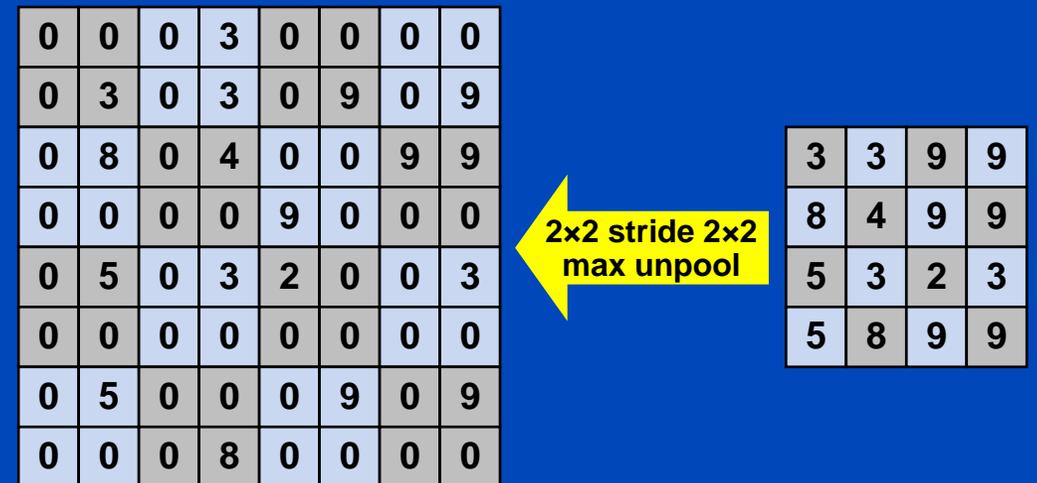
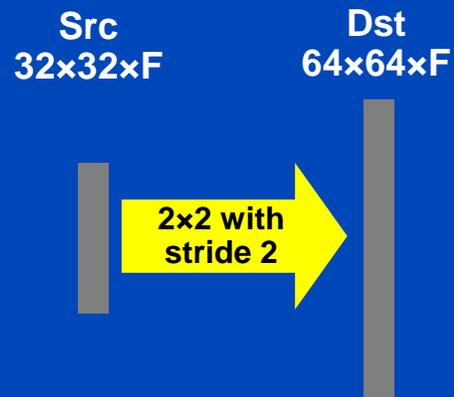
2x2 with
stride 2

$$D_{i,j,f} = \max_{b,d} S_{ai+b,cj+d,f}$$

Unpooling

Upsampling

- **Input layer S**
 - image of size I by J with F features: $I \times J \times F$
 - ...
- **Unpooling kernel**
 - pooling function, e.g. max, mean, stochastic, ...
 - size and strides
- **Output layer D**
 - increased spatial size
 - same depth

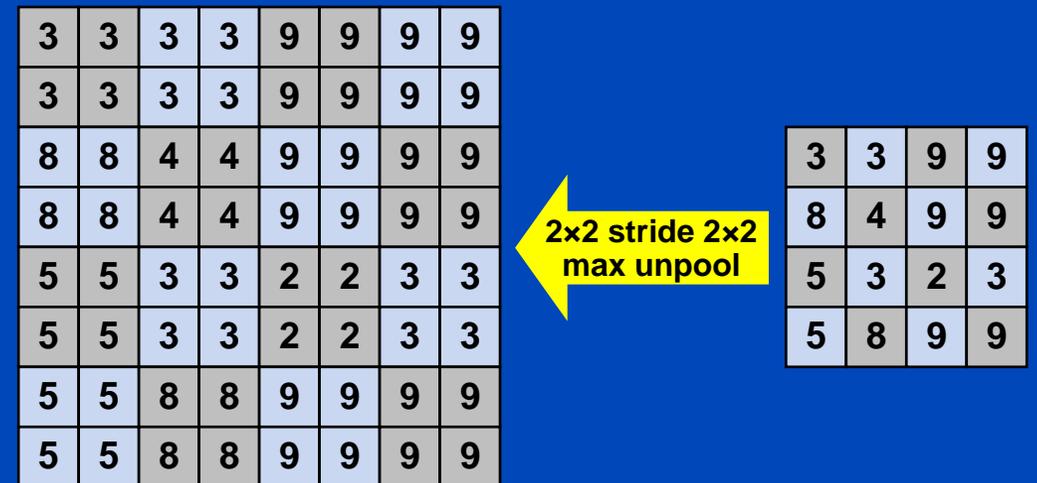


Max values at max positions that were originally found during pooling. Zeroes at non-max positions.

Unpooling

Upsampling

- **Input layer S**
 - image of size I by J with F features: $I \times J \times F$
 - ...
- **Unpooling kernel**
 - pooling function, e.g. max, mean, stochastic, ...
 - size and strides
- **Output layer D**
 - increased spatial size
 - same depth



Max values at all positions.

Dilated Convolutions

- Convolution

$$D_{i,j,g} = \sum_f S_{i,j,f} * K_{i,j,f}^g = \sum_{a,b,f} S_{i-a,j-b,f} K_{a,b,f}^g$$

- 8-dilated convolution

$$D_{i,j,g} = \sum_f S_{i,j,f} *_8 K_{i,j,f}^g = \sum_{a,b,f} S_{i-8a,j-8b,f} K_{a,b,f}^g$$

- Dilation helps to increase the receptive field of the kernel without increasing the number of unknowns in the kernel.
- Similar effect as pooling followed by convolution.

Transposed Convolution

- Sometimes also called fractionally-strided convolution layer or deconvolution layer
- Deconvolution layer is a very unfortunate name and should rather be called a transposed convolutional layer.

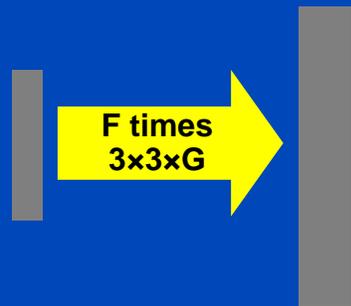
$$D_{i,j,g} = \sum_f S_{2i,2j,f} * K_{i,j,f}^g = \sum_{a,b,f} S_{2i-a,2j-b,f} K_{a,b,f}^g$$

Convolution with stride 2

Src
32x32xG

Dst
64x64xF

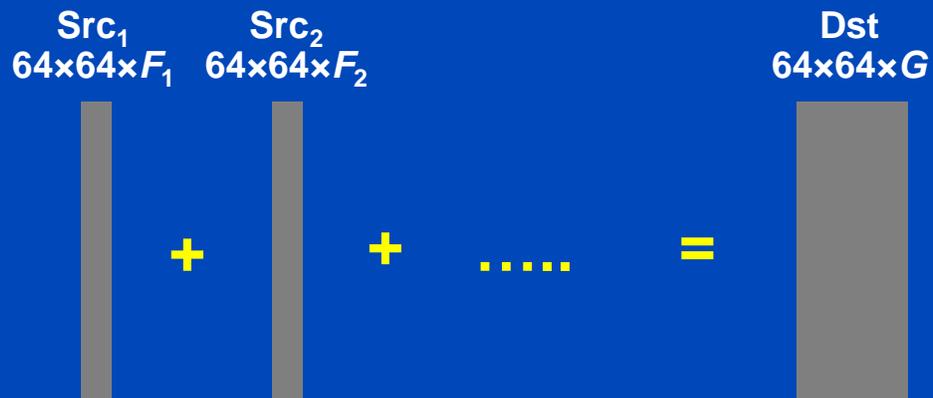
Transposed convolution with stride 2



$$D_{i,j,f} = \sum_{a,b,g} S_{i/2+a,j/2+b,g} K_{a,b,g}^f$$

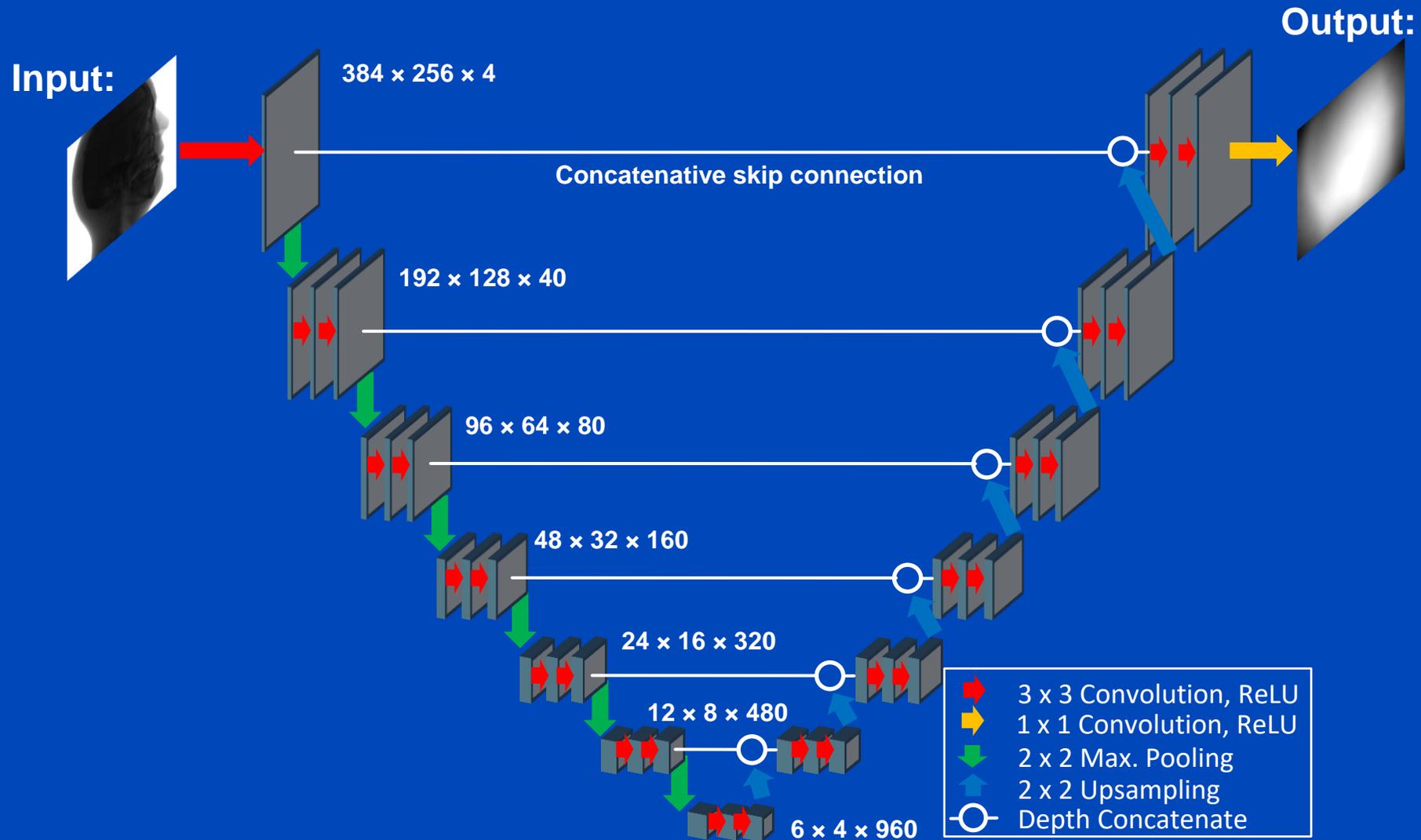
Depth Concatenation

- N input layers S_n
 - vector of size I with F_n features: $I \times F_n$
 - image of size I by J with F_n features: $I \times J \times F_n$
 - volume of size I by J by K with F_n features: $I \times J \times K \times F_n$
 - ...
- Output layer D
 - same spatial dimensions as input layer
 - $G = F_1 + F_2 + \dots + F_N$ features



$$G = \sum_n F_n$$

U-Net¹

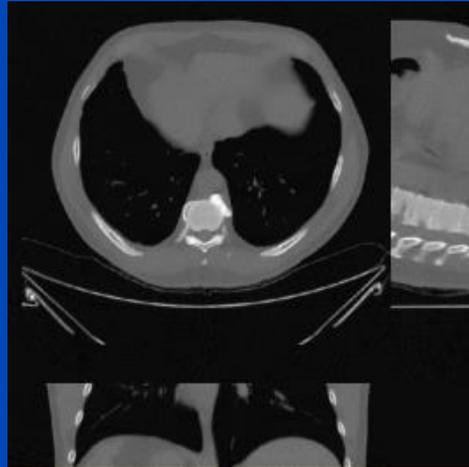


¹O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. Proc. MICCAI:234-241, 2015.

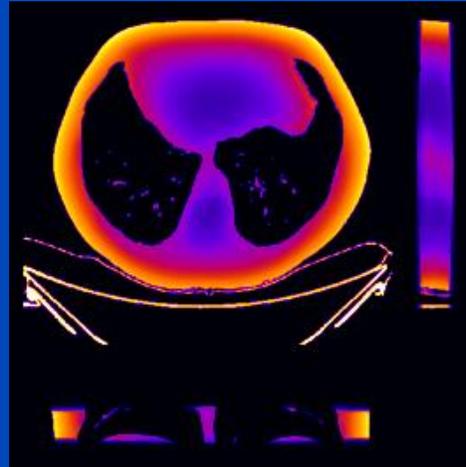
U-Net Example: The Deep Dose Estimation (DDE)

Thorax, tube A, 120 kV, no bowtie

CT image



First order dose

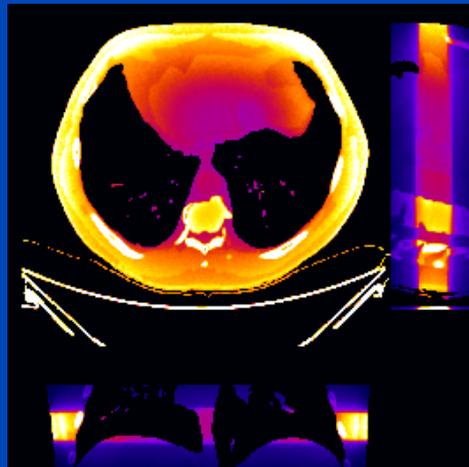


	MC	DDE
48 slices	1 h	0.25 s
whole body	20 h	5 s

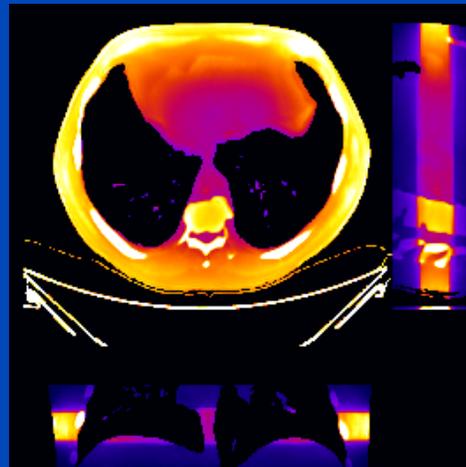
MC uses 16 CPU kernels
DDE uses one Nvidia Quadro P600 GPU

DDE training took 74 h for 300 epochs,
1440 samples, 48 slices per sample

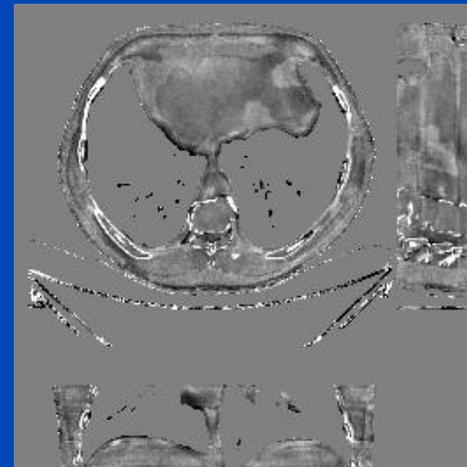
MC ground truth



DDE



Relative error



C = 0%
W = 40%

Monte Carlo (180 min)

Compute times as of 2021

Deep Dose Estimation (2 s)

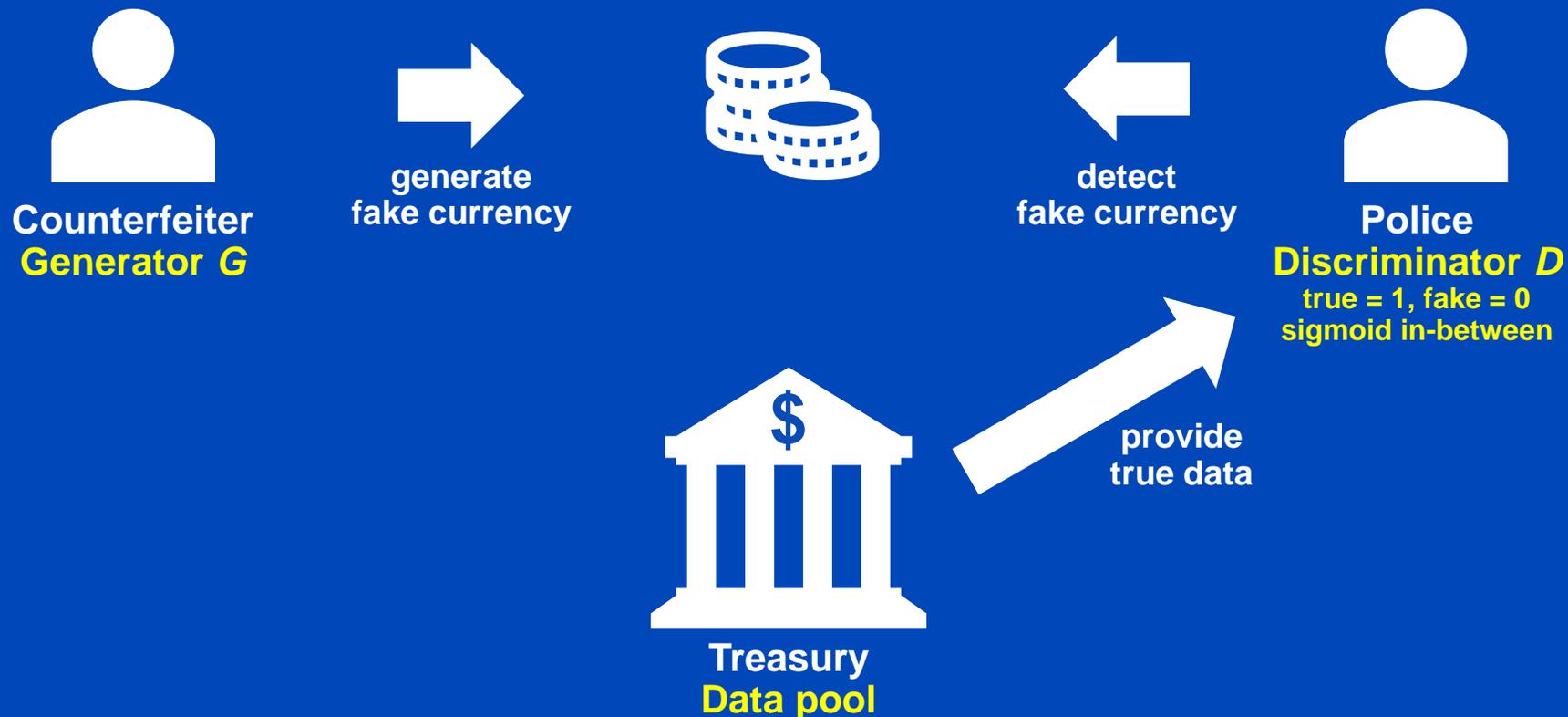
Percentage Error



FIGURE 5 Sagittal and coronal view of the dose distribution of a 100 kV whole-body spiral computed tomography (CT) scan including a bowtie filter and an angular tube current modulation. Here, the two left columns show the ground truth, the middle columns show the deep dose estimation (DDE) prediction and the two right columns the corresponding percentage error. Note that dose to air is neglected for computational reasons, and therefore, displayed as zero

Generative Adversarial Network¹ (GAN)

- Useful, if no direct ground truth (GT) is available, the training data are unpaired, unsupervised learning



¹1. Goodfellow et al. Generative Adversarial Nets, arXiv 2014

Generative Adversarial Network (GAN)

- Typical loss function and minimax game:

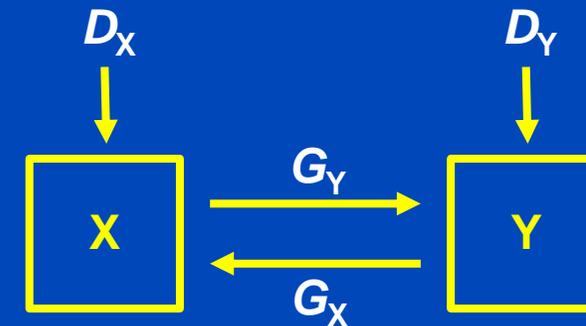
$$\min_G \max_D L(D, G) := E_x \ln (1 - D(G(x))) + E_y \ln D(y)$$

- **Conditional GAN¹**

- Conditional GANs sample the generator input x not from a uniform distribution but from a conditional distribution, e.g. noisy CT images.
- Need some measure to ensure similarity to input distribution (e.g. pixelwise loss added to the minimax loss function)

- **Cycle GAN²**

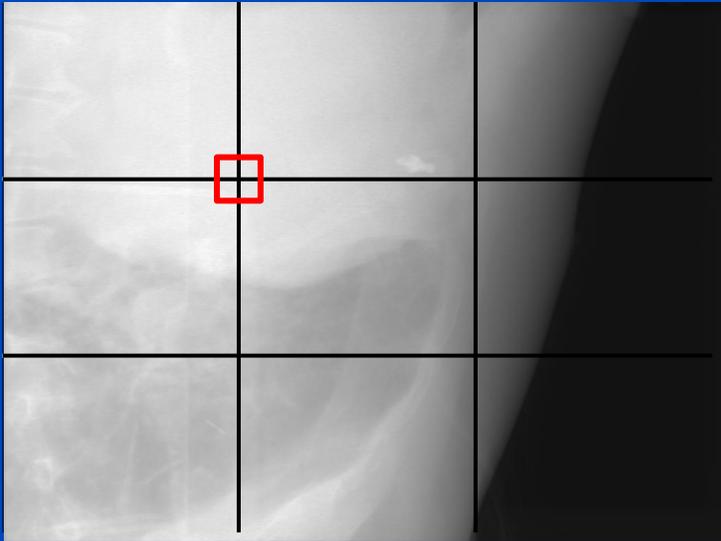
- Two GANs ($X \rightarrow Y$ and $Y \rightarrow X$)
- Demand cyclic consistency, i.e. $x = G_X(G_Y(x))$ and $y = G_Y(G_X(y))$



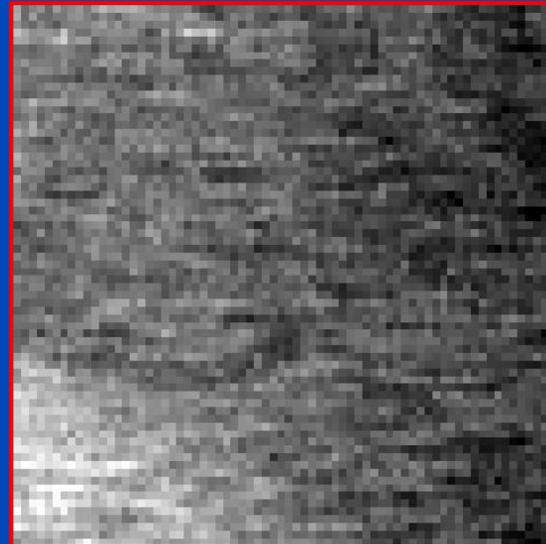
¹Isola et al. 2017

²Zhu et al., 2017

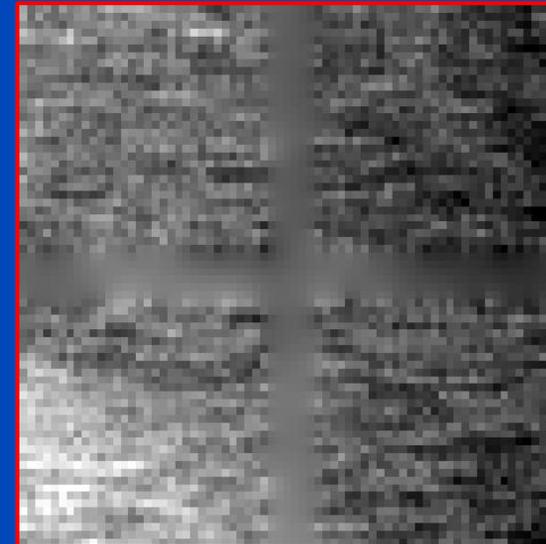
Example: Inpainting



Original

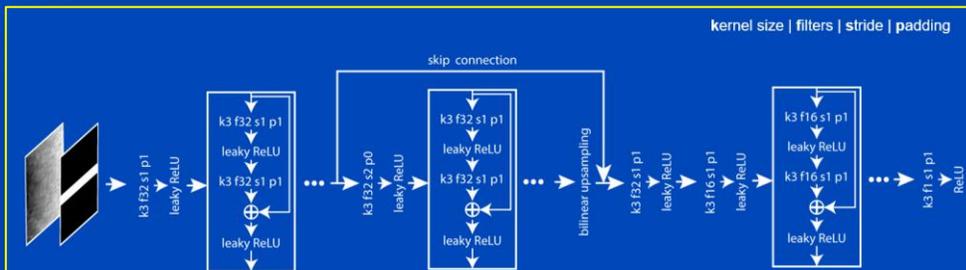


CNN with MSE only



CNN with GAN

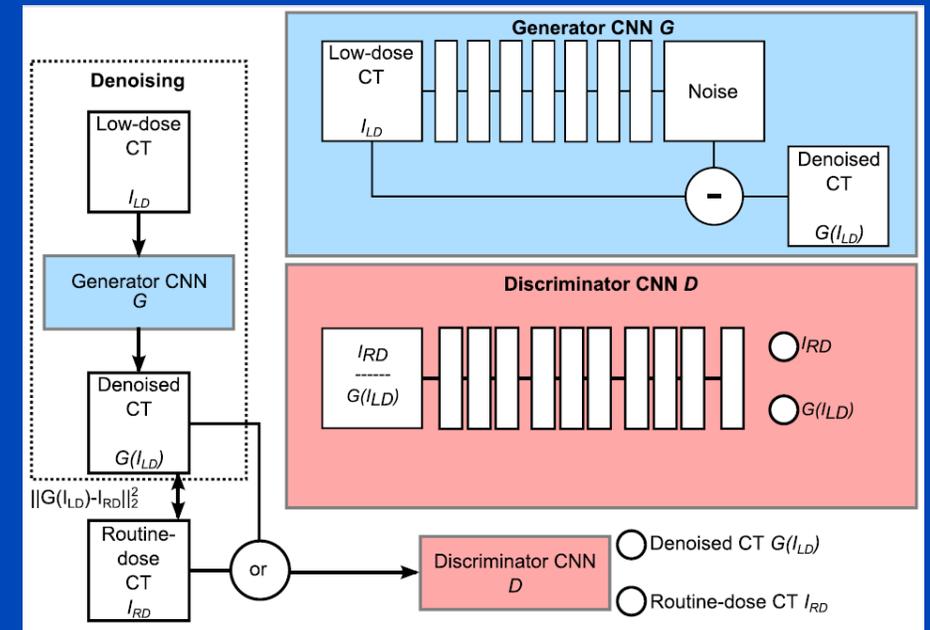
$$G(x) =$$



- Receives patch and mask as input
- Fully convolutional
- Leaky ReLUs as nonlinearities to further training stability
- Consists of several residual blocks
- One downsampling with skip connection to upsampled image to increase the receptive field.

GAN Example: Noise Removal with GAN

- Task: Reduce noise from low dose CT images.
- A conditional generative adversarial networks (GAN) is used
- Generator G :
 - 3D CNN that operates on small cardiac CT sub volumes
 - Seven $3 \times 3 \times 3$ convolutional layers yielding a receptive field of $15 \times 15 \times 15$ voxels for each destination voxel
 - Depths (features) from 32 to 128
 - Batch norm only in the hidden layers
 - Subtracting skip connection
- Discriminator D :
 - Sees either routine dose image or a generator-denoised low dose image
 - Two $3 \times 3 \times 3$ layers followed by several 3×3 layers with varying strides
 - Feedback from D prevents smoothing.
- Training data:
 - 120 kV
 - Unenhanced (why?) patient data acquired with Philips Brilliance iCT 256.
 - Two scans (why?) per patient, one with 0.2 mSv and one with 0.9 mSv effective dose.



GAN Example: Noise Removal with GAN



Low dose image (0.2 mSv)

GAN Example: Noise Removal with GAN



iDose level 3 reconstruction (0.2 mSv)

GAN Example: Noise Removal with GAN



Denoised low dose image (0.2 mSv)

GAN Example: Noise Removal with GAN



Normal dose image (0.9 mSv)

Loss Function

- The neural network coefficients (weights and biases) \mathbf{c} are chosen by minimizing a loss function (cost function)

$$\mathbf{c} = \arg \min_{\mathbf{c}} \sum_{n=1}^N L(\mathbf{c}, \mathbf{x}_n, \mathbf{y}_n)$$

with \mathbf{x}_n being the training data input, $\mathbf{y}(\mathbf{c}, \mathbf{x}_n)$ being the network output, and \mathbf{y}_n being the so-called labels, i.e. the training target, and N being the number of training samples.

- An example for such a loss function is the MSE loss

$$L(\mathbf{c}, \mathbf{x}_n, \mathbf{y}_n) = (\mathbf{y}(\mathbf{c}, \mathbf{x}_n) - \mathbf{y}_n)^2$$

Gradient Descent – A Method to find L 's Minimum

- Walk along the direction of the negative gradient
- Steepest descent
- Learning rate η

$$\mathbf{c}^{\text{new}} = \mathbf{c}^{\text{old}} - \eta \nabla_{\mathbf{c}} L(\mathbf{c}, \mathbf{x}_n, \mathbf{y}_n)$$

- Easy to understand, but not optimal
- Methods in use
 - Batch gradient descent
 - Stochastic gradient descent
 - Mini-batch gradient descent
 - Conjugate gradient descent
 - Quasi Newton methods
 - Momentum methods

Toy Example

Nested scalar functions $f(c, x)$ with unknown coefficients c

Loss
Function

$$L(c_3, c_2, c_1, x) = (f_3(c_3, f_2(c_2, f_1(c_1, x))) - y)^2$$

last layer output
2nd layer output
1st layer output
input
training target

Intermediate
Values

$$L_3 = \frac{dL}{df_3} = 2(f_3 - y)$$
$$L_2 = \frac{dL}{df_2} = \frac{dL}{df_3} \frac{df_3}{df_2} = L_3 \frac{df_3}{df_2}$$

...

$$L_n = \frac{dL}{df_n} = L_{n+1} \frac{df_{n+1}}{df_n}$$

Desired
Gradients

$$\frac{dL}{dc_3} = \frac{dL}{df_3} \frac{df_3}{dc_3} = L_3 \frac{df_3}{dc_3}$$
$$\frac{dL}{dc_2} = \frac{dL}{df_3} \frac{df_3}{df_2} \frac{df_2}{dc_2} = L_2 \frac{df_2}{dc_2}$$

...

$$\frac{dL}{dc_n} = L_n \frac{df_n}{dc_n}$$

Toy Example

Nested scalar functions $f(c, x)$ with unknown coefficients c

Loss
Function

$$L(c_3, c_2, c_1, x) = (f_3(c_3, f_2(c_2, f_1(c_1, x))) - y)^2$$

last layer output
2nd layer output
1st layer output
input
training target

Intermediate
Values

$$\begin{aligned} L_3 &= \frac{dL}{df_3} = 2(f_3 - y) \\ L_2 &= \frac{dL}{df_2} = \frac{dL}{df_3} \frac{df_3}{df_2} = L_3 \frac{df_3}{df_2} \\ &\dots \\ L_n &= \frac{dL}{df_n} = L_{n+1} \frac{df_{n+1}}{df_n} \end{aligned}$$

Backpropagation

Desired
Gradients

$$\begin{aligned} \frac{dL}{dc_3} &= \frac{dL}{df_3} \frac{df_3}{dc_3} = L_3 \frac{df_3}{dc_3} \\ \frac{dL}{dc_2} &= \frac{dL}{df_3} \frac{df_3}{df_2} \frac{df_2}{dc_2} = L_2 \frac{df_2}{dc_2} \\ &\dots \\ \frac{dL}{dc_n} &= L_n \frac{df_n}{dc_n} \end{aligned}$$

Toy Example 2

Nested vector-valued functions $f(\mathbf{c}, \mathbf{x})$ with unknown coefficients vectors \mathbf{c}

Loss
Function

$$L(\mathbf{c}_3, \mathbf{c}_2, \mathbf{c}_1, \mathbf{x}) = \left(\underset{n}{f_3}(\underset{\nu}{\mathbf{c}_3}, \underset{m}{f_2}(\underset{\mu}{\mathbf{c}_2}, \underset{l}{f_1}(\underset{\lambda}{\mathbf{c}_1}, \mathbf{x}))) - \underset{\text{training target}}{\mathbf{y}} \right)^2$$

last layer output
2nd layer output
1st layer output
input

Intermediate
Values

$$L_{3n} = \frac{dL}{df_{3n}} = 2(f_3 - \mathbf{y})_n$$

$$L_{2m} = \frac{dL}{df_{2m}} = \sum_n \frac{dL}{df_{3n}} \frac{df_{3n}}{df_{2m}} = \sum_n L_{3n} \frac{df_{3n}}{df_{2m}}$$

$$L_{1l} = \frac{dL}{df_{1l}} = \dots = \sum_m L_{2m} \frac{df_{2m}}{df_{1l}}$$

Desired
Gradients

$$\frac{dL}{dc_{3\nu}} = \sum_n \frac{dL}{df_{3n}} \frac{df_{3n}}{dc_{3\nu}} = \sum_n L_{3n} \frac{df_{3n}}{dc_{3\nu}} = 0$$

$$\frac{dL}{dc_{2\mu}} = \sum_{nm} \frac{dL}{df_{3n}} \frac{df_{3n}}{df_{2m}} \frac{df_{2m}}{dc_{2\mu}} = \sum_m L_{2m} \frac{df_{2m}}{dc_{2\mu}} = 0$$

$$\frac{dL}{dc_{1\lambda}} = \dots = \sum_l L_{1l} \frac{df_{1l}}{dc_{1\lambda}}$$

Toy Example 3

Backprojection of convolved CNN-preprocessed rawdata with unknown \mathbf{c}_1

Loss Function

$$L(\mathbf{c}_3, \mathbf{c}_2, \mathbf{c}_1, \mathbf{x}) = \left(\underset{n}{f_3}(\underset{\nu}{\mathbf{c}_3}, \underset{m}{f_2}(\underset{\mu}{\mathbf{c}_2}, \underset{l}{f_1}(\underset{\lambda}{\mathbf{c}_1}, \mathbf{x}))) - \mathbf{y} \right)^2$$

backprojection
convolution
CNN
rawdata
target image

Intermediate Values

$$L_{3n} = \frac{dL}{df_{3n}} = \text{image} = 2 \cdot \text{reconstructed} - \text{target image}$$

backprojection derived wrt the m -th sinogram value, i.e. image where pixel n holds the intersection of ray m with pixel n

$$L_{2m} = \frac{dL}{df_{2m}} = \sum_n \frac{dL}{df_{3n}} \frac{df_{3n}}{df_{2m}} = \sum_n L_{3n} \frac{df_{3n}}{df_{2m}} = \text{forward projection of image } L_3$$

$$L_{1l} = \frac{dL}{df_{1l}} = \dots = \sum_m L_{2m} \frac{df_{2m}}{df_{1l}} = \text{convolution of } L_2 \text{ with RamLak kernel}$$

convolution derived wrt the l -th sinogram value, i.e. sinogram where pixel m holds the convolution kernel entry $m-l$

Desired Gradients

$$\frac{dL}{dc_{3\nu}} = \sum_n \frac{dL}{df_{3n}} \frac{df_{3n}}{dc_{3\nu}} = \sum_n L_{3n} \frac{df_{3n}}{dc_{3\nu}} = 0$$

$$\frac{dL}{dc_{2\mu}} = \sum_{nm} \frac{dL}{df_{3n}} \frac{df_{3n}}{df_{2m}} \frac{df_{2m}}{dc_{2\mu}} = \sum_m L_{2m} \frac{df_{2m}}{dc_{2\mu}} = 0$$

$$\frac{dL}{dc_{1\lambda}} = \dots = \sum_l L_{1l} \frac{df_{1l}}{dc_{1\lambda}}$$

Gradient Descent

For each epoch:

Shuffle the N data samples

For each batch (batch size B):

For each data sample \mathbf{x}_b of the current batch:

Calculate the loss function's gradient

$$\frac{dL}{d\mathbf{c}} = \nabla_{\mathbf{c}} L(\mathbf{c}, \mathbf{x}_b)$$

wrt the unknowns \mathbf{c} by backpropagation.

Now, update the network parameters (weights, biases, etc.):

$$\mathbf{c}^{\text{new}} = \mathbf{c}^{\text{old}} - \frac{\eta}{B} \sum_{b=1}^B \nabla_{\mathbf{c}} L(\mathbf{c}, \mathbf{x}_b)$$

$B = 1$: stochastic gradient descent
$B = N$: gradient descent
else	: batch gradient descent

Optimization Algorithms

SGD with Momentum: Add fraction of past update vector to current update vector

$$v' = \gamma v + \eta \cdot \nabla_{\theta} C(\theta)$$

$$\theta' = \theta - v'$$

Nesterov Accelerated Gradient (NAG): Compute gradients with respect to the approximate future position of parameters

$$v' = \gamma v + \eta \cdot \nabla_{\theta} C(\theta - \gamma v)$$

$$\theta' = \theta - v'$$

Adagrad: Adapt updates to each individual parameter

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t$$

$$g_{t,i} = \nabla_{\theta} C(\theta_{t,i})$$

$$G_t = \begin{pmatrix} \sum_{t' < t} g_{t',1}^2 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sum_{t' < t} g_{t',N}^2 \end{pmatrix}$$

Optimization Algorithms

Adadelta: Improve on Adagrad by restricting the window of past time steps to some fixed window → No vanishing gradients

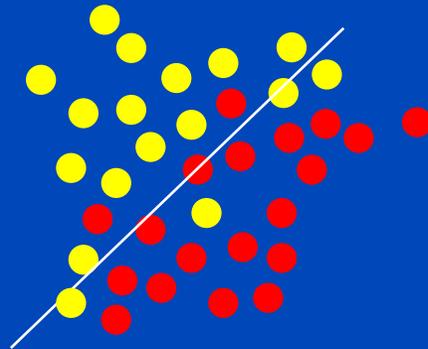
Additionally, perform normalization of gradients → No learning rate needed

RMSprop: Identical to Adadelta, without normalization. Adapt learning rate with exponentially decaying average of past squared gradients.

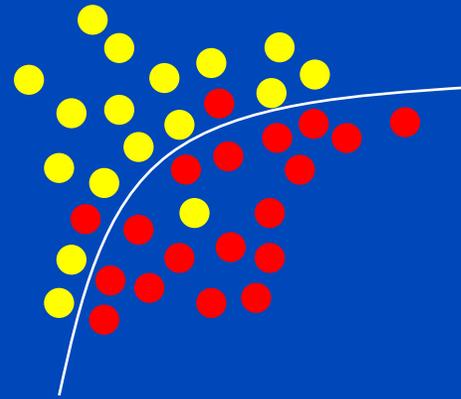
Adaptive Moment Estimation (Adam): Momentum for first, and second order moment (mean and variance) of gradients

$$\begin{aligned} m_t^{(1)} &= \beta_1 m_{t-1}^{(1)} + (1 - \beta_1) g_t & \hat{m}_t^{(1)} &= \frac{m_t^{(1)}}{1 - \beta_1} \\ m_t^{(2)} &= \beta_2 m_{t-1}^{(2)} + (1 - \beta_2) g_t^2 & \hat{m}_t^{(2)} &= \frac{m_t^{(2)}}{1 - \beta_2} \end{aligned} \quad \theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{m}_t^{(2)} + \epsilon}} \odot \hat{m}_t^{(1)}$$

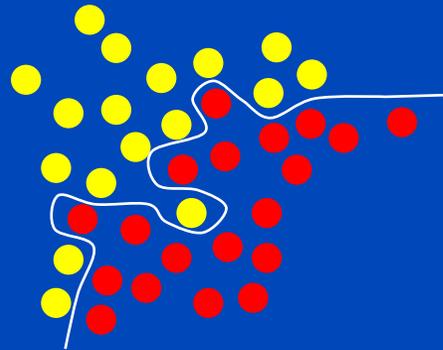
Underfitting



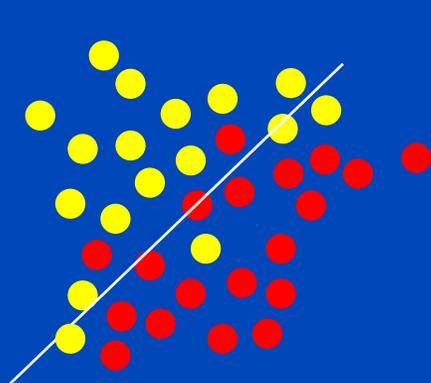
Fitting



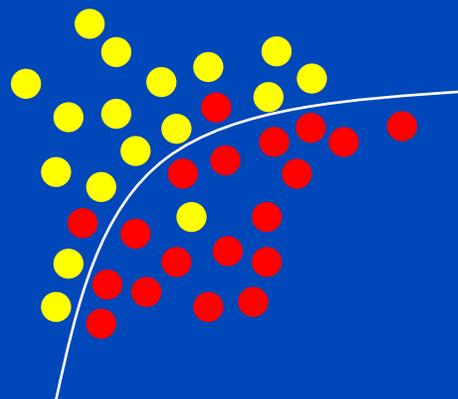
Overfitting



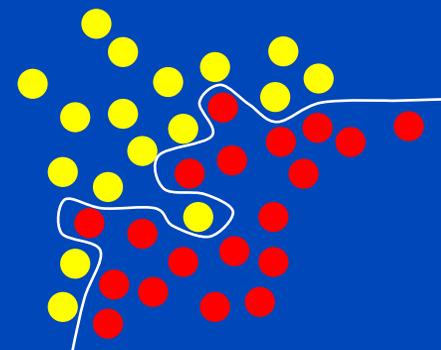
Fitting



underfit



reasonable



overfit

Overfitting

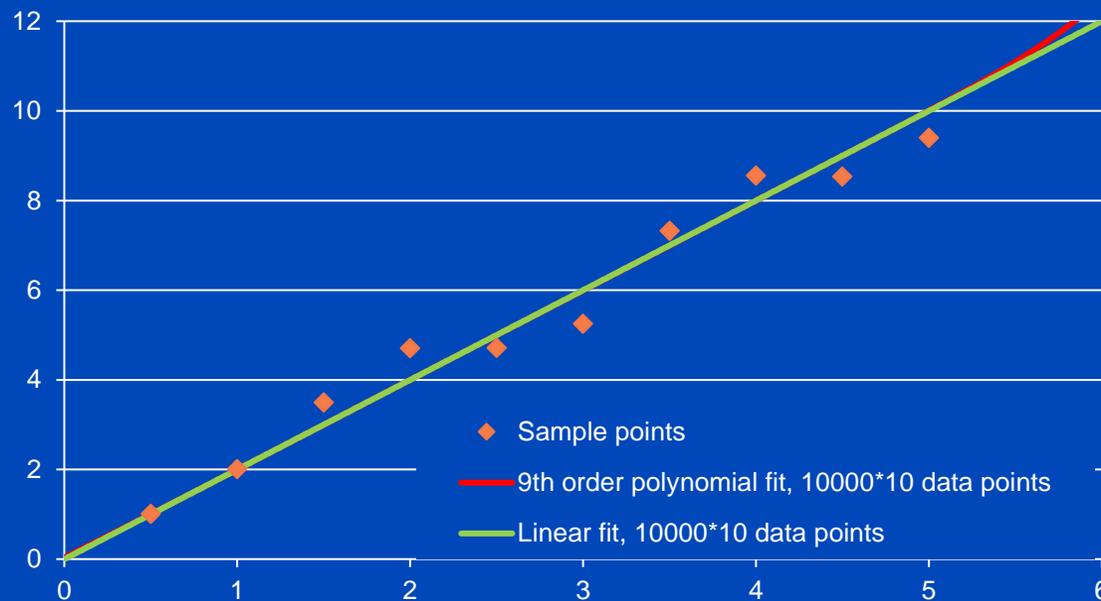
- Assume our training data result from sampling the function $f(x) = 2x$ at a given number of points.
- Since the sampling might include some random noise, the samples slightly deviate from the function $f(x) = 2x$.
- A 9th order polynomial perfectly fits the training data, but fails to appropriately predict test data such as $x = 0.25$ for instance.



Regularization

Increase of training data

- The increase of the amount of training data makes the network more robust against single deviations.
- The training data can also be increased artificially.
- Similar results can be observed if the polynomial is fitted to 100000 samples.



Coefficients	Linear	9th order
C_0	-0.00295	0.03343
C_1	2.000325	1.904762
C_2		0.079125
C_3		-0.02262
C_4		0.000435
C_5		-4.96E-05
C_6		0.000339
C_7		-4.25E-05
C_8		-9.19E-06
C_9		1.43E-06

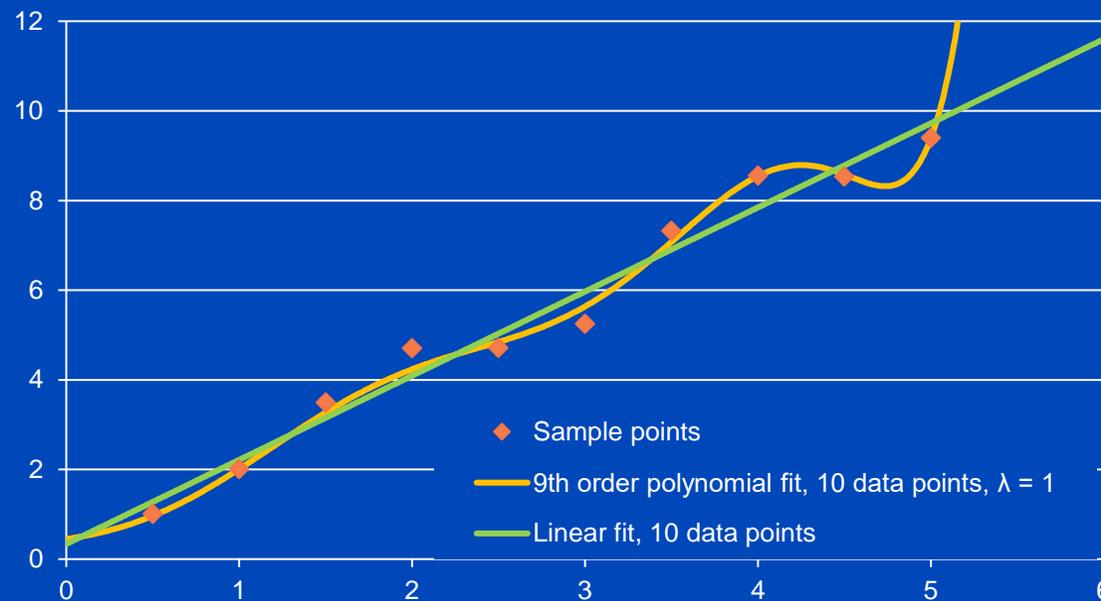
Regularization

Penalizing large weights

- Modification of the cost function to penalize large weight (i.e. quadratic penalty):

$$C = C_0 + \lambda \sum_w w^2$$

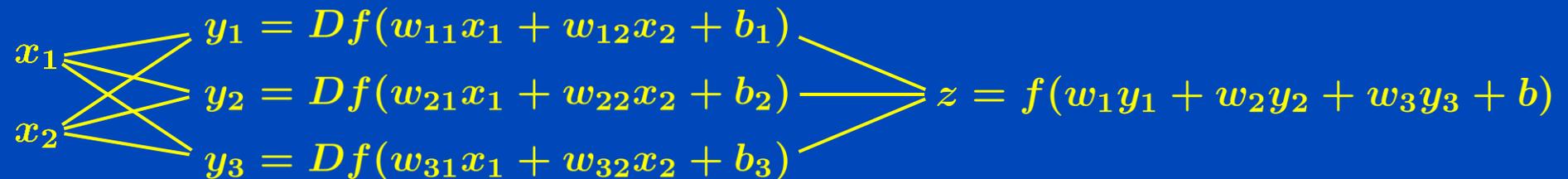
- If a certain weight is large, the output strongly depends on the input of that weight.



Coefficients	Linear	9th order
C_0	-0.00295	0.447558
C_1	2.000325	0.575279
C_2		0.665781
C_3		0.562606
C_4		0.049884
C_5		-0.45894
C_6		0.186099
C_7		-0.01496
C_8		-0.00342
C_9		0.000471

Regularization: Dropout

- Dropout is intended to reduce overfitting.
- Dropout randomly (before each mini batch) zeroes activations with probability $1-p$ during training.
- Let D be a Bernoulli random variate with $P(D = 1) = p$ and regard the following toy example:



- After training, on inference, the outputs y_1 , y_2 , and y_3 need to be multiplied by p to be equivalent to the training situation, on expectation (alternatively divide by p during training).
- Dropout can be realized using dropout layers.

Batch Normalization

Batch normalization

- normalizes each activation to have zero expectation and unit variance within the batch
- introduces trainable scale and offset for each activation (or for each feature map) to, potentially, denormalize again
- is part of the model architecture
- reduces the need for dropout
- reduces internal covariate shift and thus accelerates training
- fixes the means and variances of layer inputs
- improves gradient flow through the network
- allows for higher learning rates without the risk of divergence
- prevents the net from getting trapped in saturated modes
- makes it possible to use saturating nonlinearities

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$;
Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

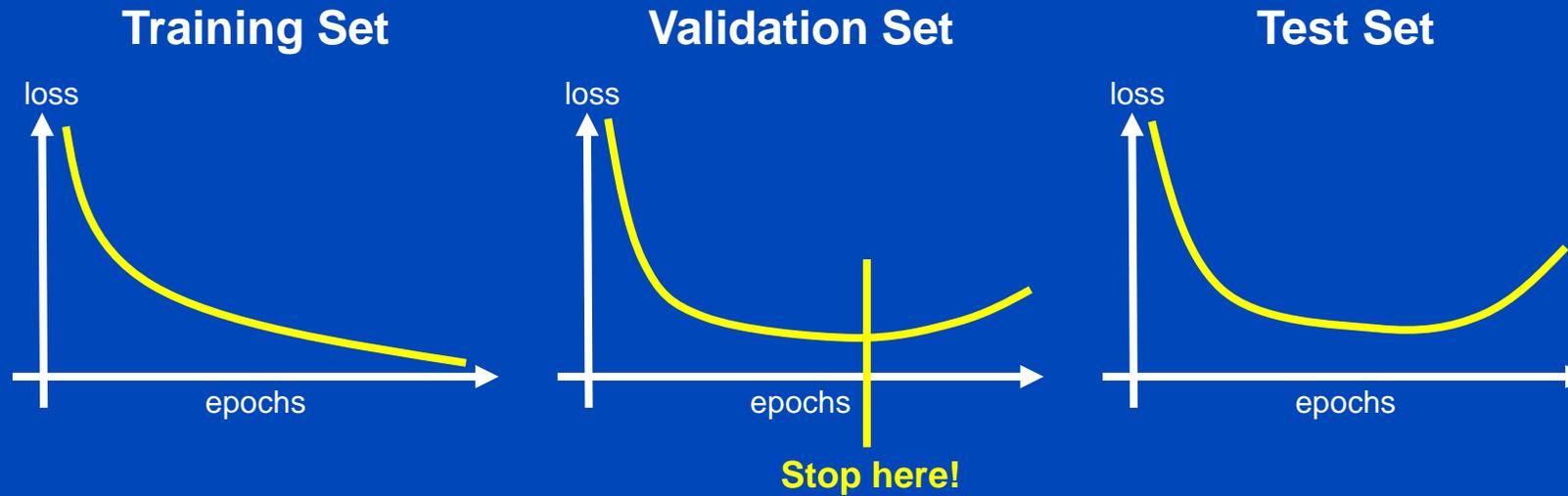
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

Further Means to Avoid Overfitting

- **Choose adequate network architecture**
- **Preprocess data**
 - Normalize data (mean, var, ...)
 - Add prior knowledge (e.g. $\exp(-x)$)
- **Data augmentation**
 - Random transformations (mirror, affine, deformable, ...)
 - Gray value distribution
 - Change spatial resolution
 - Add noise
 - ...
- **Penalize loss function**
 - Enforce small weights
 - Enforce sparse weights
 - ...

Learning Curve



- Training and validation set are part of the training
- Do not use test set for training
- Early stopping (at minimum validation loss)
- Training : Validation : Test \approx 70 : 20 : 10

Weight Initialization

- Weights in neural networks should be initialized such that the neurons are not saturated (since saturation often decreases the learning rate).
- Assume we have a fully-connected network with 1000 input neurons.
- Let us further assume that half of the input equals 1 and the other half equals 0.
- If the weights and the bias are initialized with Gaussian random numbers with zero mean and a standard deviation of 1, the weighted sum $z = \sum w_j x_j + b$ to the first hidden neuron is zero mean Gaussian with standard deviation $\sigma = \sqrt{501} \approx 22.4$.
- Thus, it is very likely that $z \ll 0$ or $z \gg 0$ and the neuron saturates.
- Therefore, if we have n_{in} inputs, an initialization with Gaussian random numbers with zero mean and a standard deviation of $1/\sqrt{n_{\text{in}}}$ would be a better choice.

What was not Discussed Here

- Attention mechanism
- Transformer networks
- Diffusion networks
- ...

Thank You!



This presentation will soon be available at www.dkfz.de/ct.

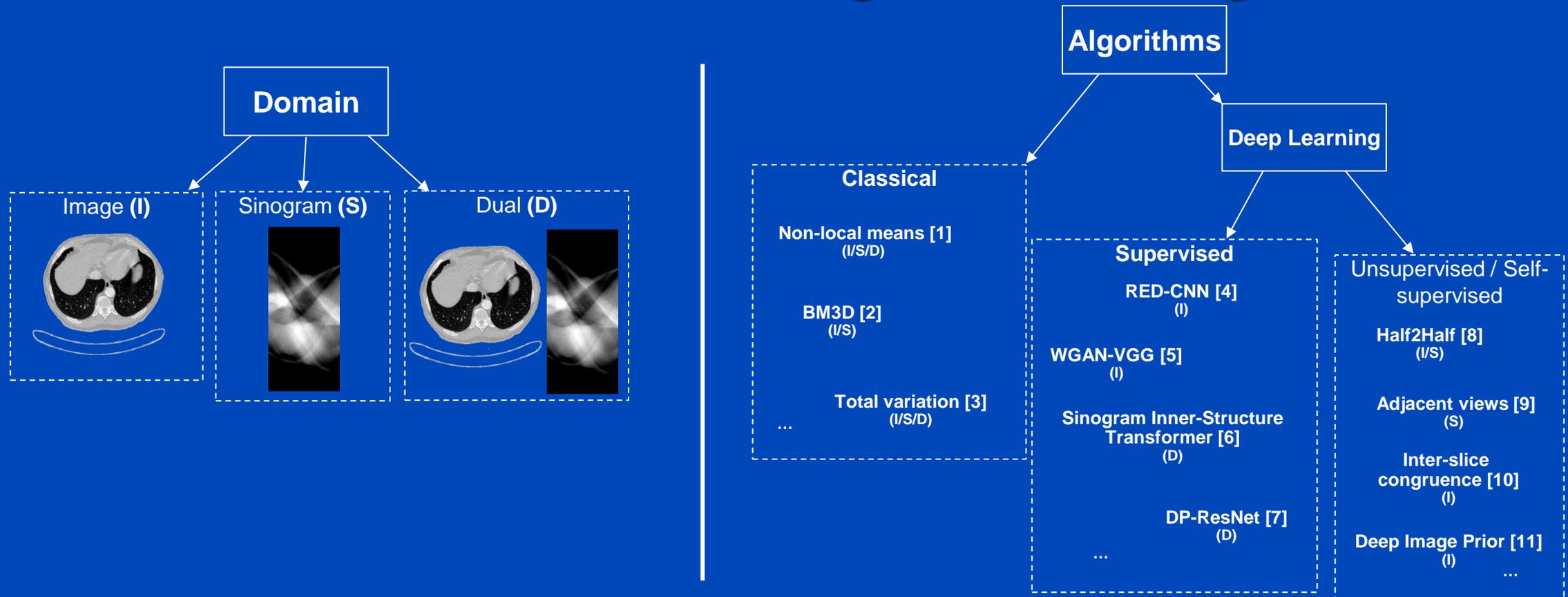
Job opportunities through DKFZ's international PhD or Postdoctoral Fellowship programs (marc.kachelriess@dkfz.de).

Parts of the reconstruction software were provided by RayConStruct® GmbH, Nürnberg, Germany.

Denoising benchmark with surprising results

PART 10: IS NEWER ALWAYS BETTER?

Low Dose CT Image Denoising



[1] Yuan Y, Zhang YB, Yu HY (2018) "Adaptive nonlocal means method for denoising basis material images [...]". *J Comput Assist Tomogr* 42:972–981.

[2] Feruglio, P Fumene, C Vinegoni, J Gros, A Sbarbati, and R Weissleder (2010) "Block Matching 3D Random Noise Filtering for Absorption Optical Projection Tomography." *Physics in Medicine and Biology* 55 (18): 5401–15.

[3] Jia L, Zhang Q, Shang Y, et al. (2018) "Denoising for low-dose CT image by discriminative weighted nuclear norm minimization". IEEE Access

[4] Chen, Hu, Yi Zhang, Mannudeep K. Kalra, Feng Lin, Yang Chen, Peixi Liao, Jiliu Zhou, and Ge Wang. 2017. "Low-Dose CT with a Residual Encoder-Decoder Convolutional Neural Network." *IEEE Transactions on Medical Imaging* 36 (12): 2524–35.

[5] Yang, Qingsong, Pingkun Yan, Yanbo Zhang, et al.. 2018. "Low-Dose CT Image Denoising Using a Generative Adversarial Network with Wasserstein Distance and Perceptual Loss." *IEEE Transactions on Medical Imaging* 37 (6): 1348–57.

[6] L. Yang, Z. Li, R. Ge, J. Zhao, H. Si and D. Zhang, "Low-Dose CT Denoising via Sinogram Inner-Structure Transformer," in *IEEE Transactions on Medical Imaging*, vol. 42, no. 4, pp. 910-921, 2023.

[7] Yin, Xiangrui, Qianlong Zhao, Jin Liu, Wei Yang, Jian Yang, Guotao Quan, Yang Chen, Huazhong Shu, Limin Luo, and Jean-Louis Coatrieux. 2019. "Domain Progressive 3D Residual Convolution Network to Improve Low-Dose CT Imaging." *IEEE TMI* 38 (12).

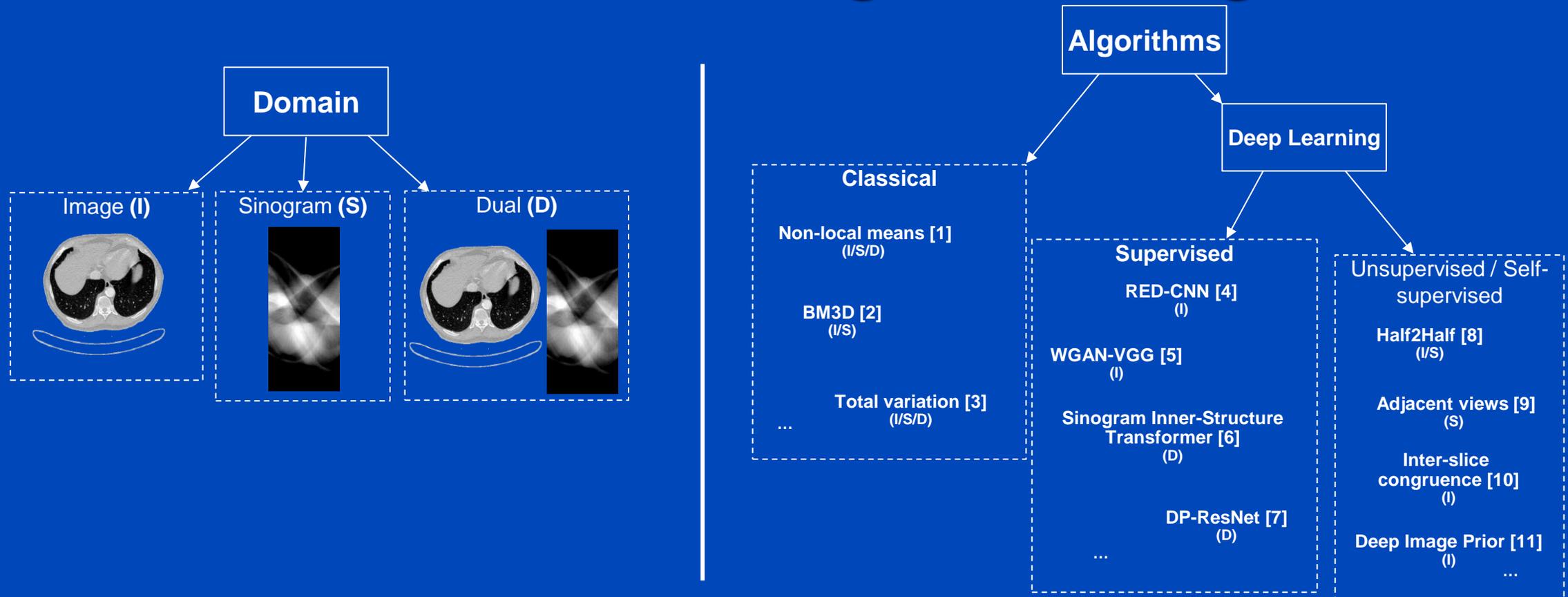
[8] Yuan, N., Zhou, J., & Qi, J. (2020). Half2Half: deep neural network based CT image denoising without independent reference data. *Physics in Medicine & Biology*, 65(21), 215020.

[9] Hong, Zixuan, Dong Zeng, Xi Tao, and Jianhua Ma. 2023. "Learning CT Projection Denoising from Adjacent Views." *Medical Physics* 50 (3): 1367–77.

[10] Bera, Sutanu, and Prabir Kumar Biswas. 2023. "Self Supervised Low Dose Computed Tomography Image Denoising Using Invertible Network Exploiting Inter Slice Congruence." In *2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 5603–12.

[11] Baguer, Daniel Otero, Johannes Leuschner, and Maximilian Schmidt. 2020. "Computed Tomography Reconstruction Using Deep Image Prior and Learned Reconstruction Methods." *Inverse Problems* 36 (9).

Low Dose CT Image Denoising



[1] Yuan Y, Zhang YB, Yu HY (2018) "Adaptive nonlocal means method for denoising basis material images [...]". *J Comput Assist Tomogr* 42:972–981.

[2] Feruglio, P Fumene, C Vinegoni, J Gros, A Sbarbati, and R Weissleder (2010) "Block Matching 3D Random Noise Filtering for Absorption Optical Projection Tomography." *Physics in Medicine and Biology* 55 (18): 5401–15.

[3] Jia L, Zhang Q, Shang Y, et al. (2018) "Denoising for low-dose CT image by discriminative weighted nuclear norm minimization". *IEEE Access*

[4] Chen, Hu, Yi Zhang, Mannudeep K. Kalra, Feng Lin, Yang Chen, Peixi Liao, Jiliu Zhou, and Ge Wang. 2017. "Low-Dose CT with a Residual Encoder-Decoder Convolutional Neural Network." *IEEE Transactions on Medical Imaging* 36 (12): 2524–35.

[5] Yang, Qingsong, Pingkun Yan, Yanbo Zhang, et al.. 2018. "Low-Dose CT Image Denoising Using a Generative Adversarial Network with Wasserstein Distance and Perceptual Loss." *IEEE Transactions on Medical Imaging* 37 (6): 1348–57.

[6] L. Yang, Z. Li, R. Ge, J. Zhao, H. Si and D. Zhang, "Low-Dose CT Denoising via Sinogram Inner-Structure Transformer," in *IEEE Transactions on Medical Imaging*, vol. 42, no. 4, pp. 910-921, 2023.

[7] Yin, Xiangrui, Qianlong Zhao, Jin Liu, Wei Yang, Jian Yang, Guotao Quan, Yang Chen, Huazhong Shu, Limin Luo, and Jean-Louis Coatrieux. 2019. "Domain Progressive 3D Residual Convolution Network to Improve Low-Dose CT Imaging." *IEEE TMI* 38 (12).

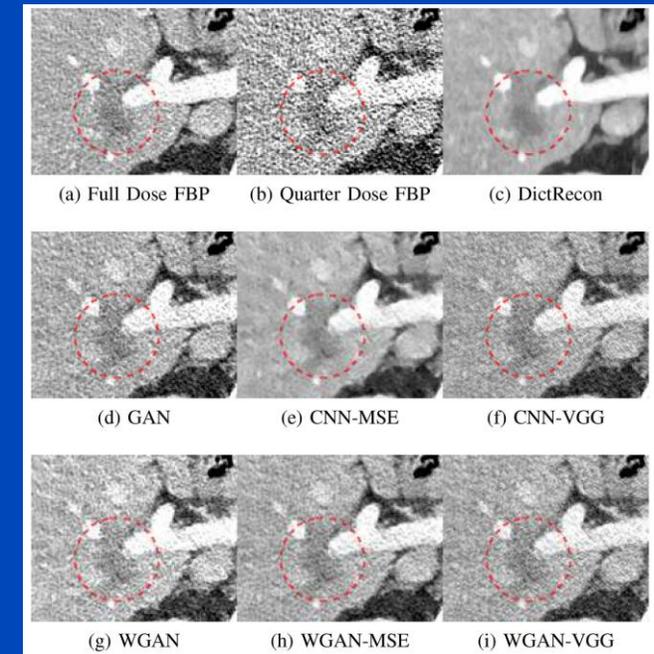
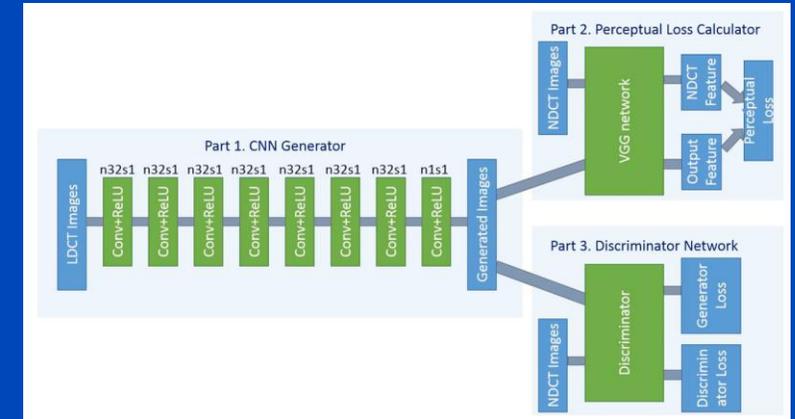
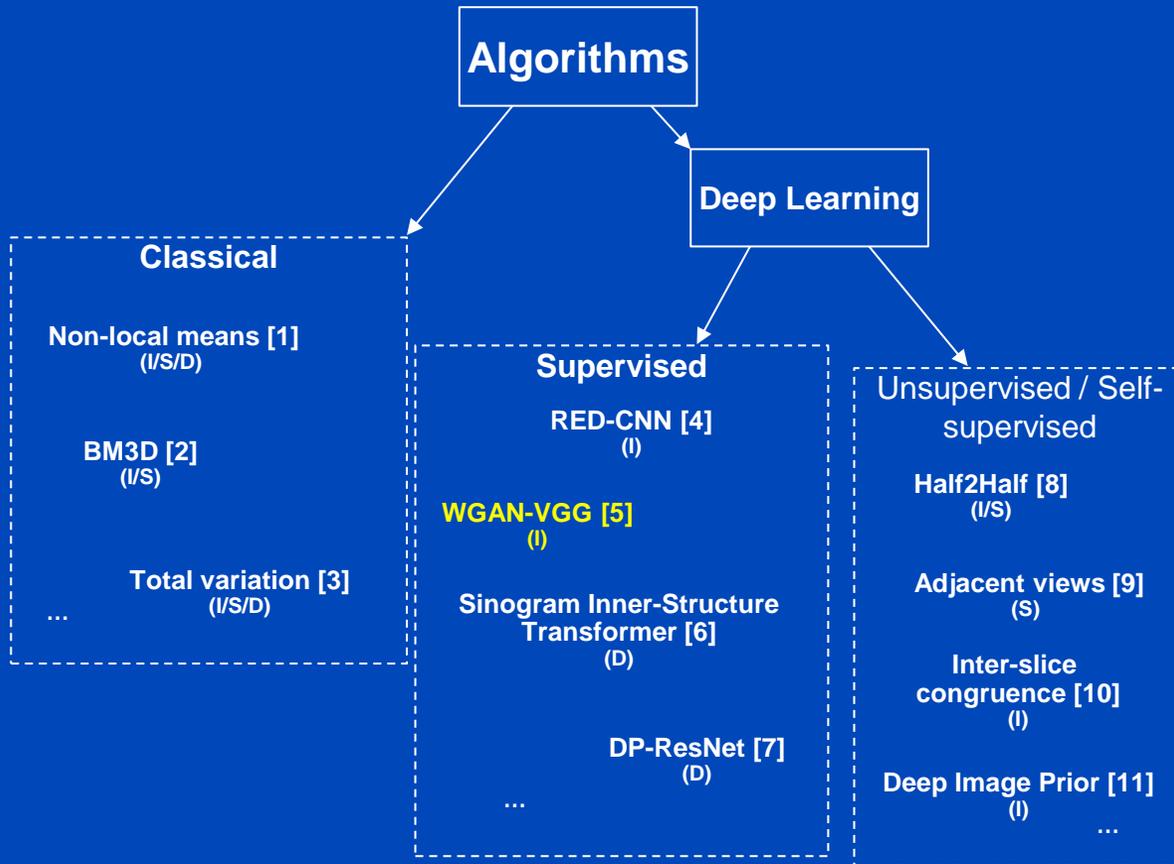
[8] Yuan, N., Zhou, J., & Qi, J. (2020). Half2Half: deep neural network based CT image denoising without independent reference data. *Physics in Medicine & Biology*, 65(21), 215020.

[9] Hong, Zixuan, Dong Zeng, Xi Tao, and Jianhua Ma. 2023. "Learning CT Projection Denoising from Adjacent Views." *Medical Physics* 50 (3): 1367–77.

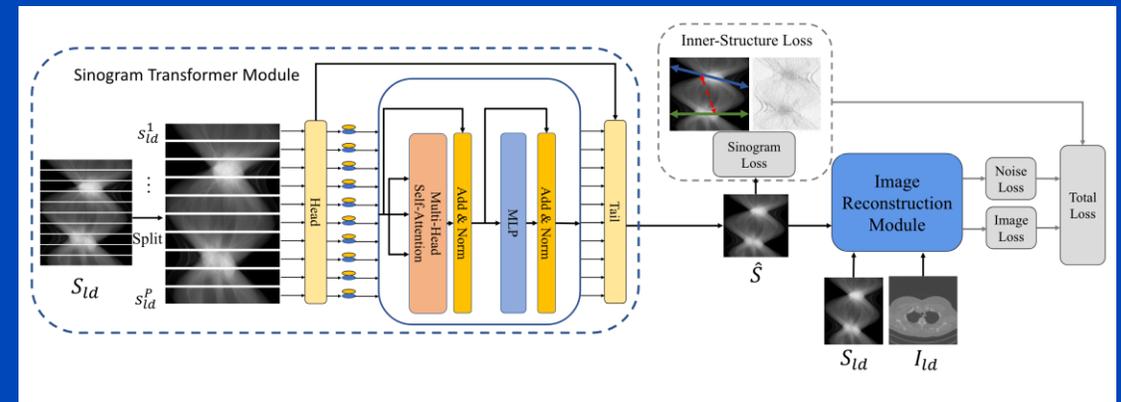
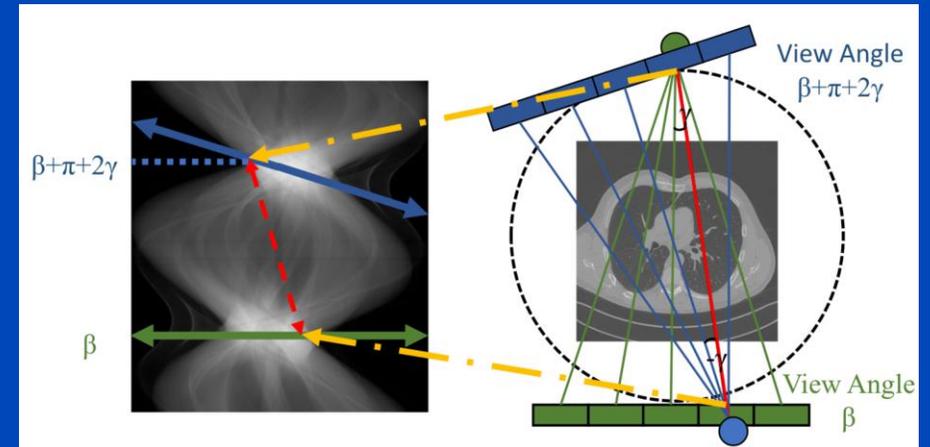
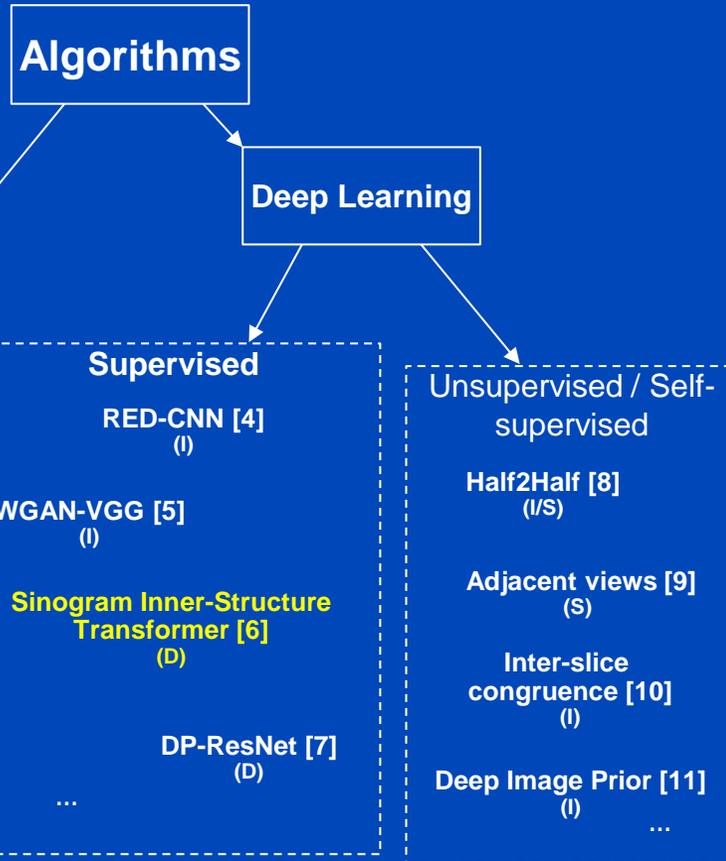
[10] Bera, Sutanu, and Prabir Kumar Biswas. 2023. "Self Supervised Low Dose Computed Tomography Image Denoising Using Invertible Network Exploiting Inter Slice Congruence." In *2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 5603–12.

[11] Baguer, Daniel Otero, Johannes Leuschner, and Maximilian Schmidt. 2020. "Computed Tomography Reconstruction Using Deep Image Prior and Learned Reconstruction Methods." *Inverse Problems* 36 (9).

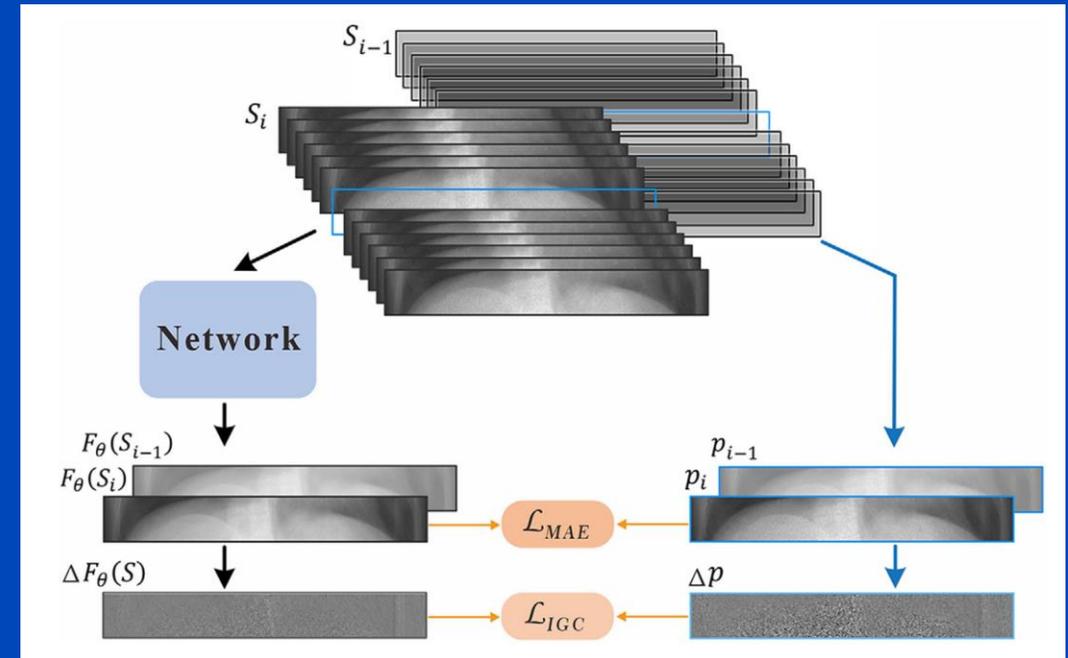
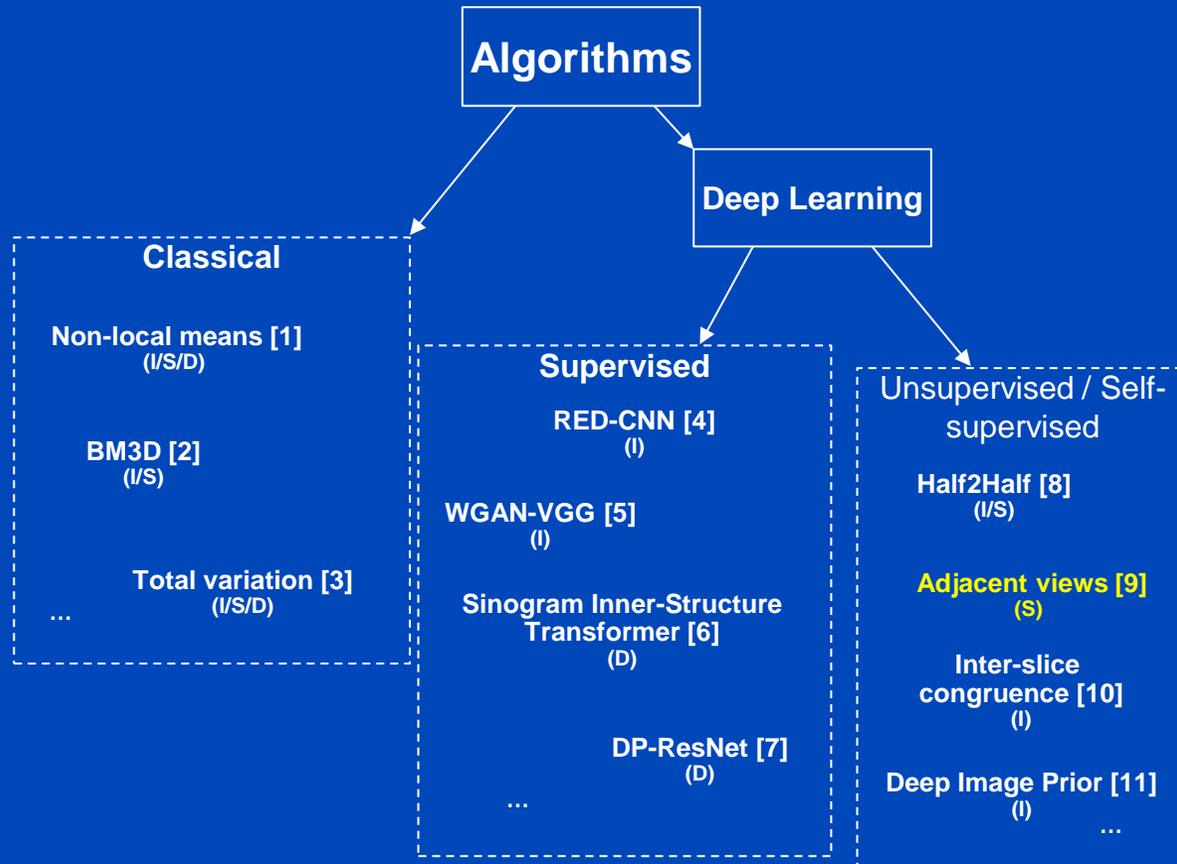
Low Dose CT Image Denoising



Low Dose CT Image Denoising



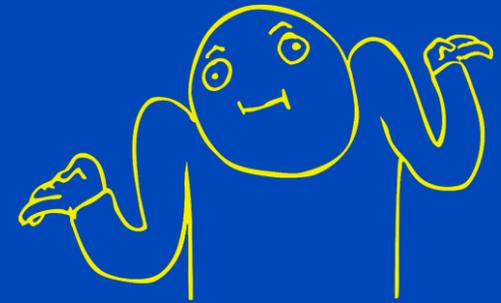
Low Dose CT Image Denoising



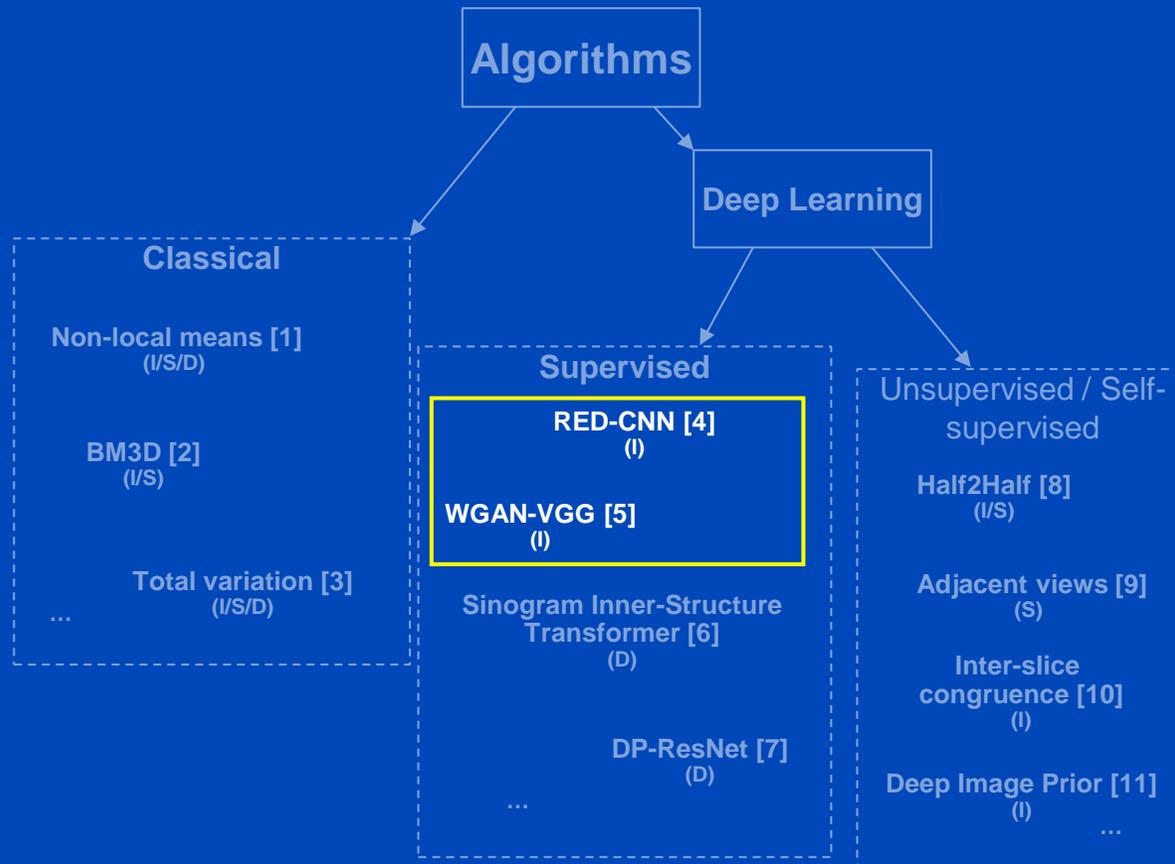
Low Dose CT Image Denoising

Q: Which algorithm performs best?

A:



Low Dose CT Image Denoising



$$\theta^* = \arg \min_{\theta} \mathbb{E}_{x,y \sim \mathcal{D}^{\text{train}}} \|f_{\theta}(x) - y\|$$

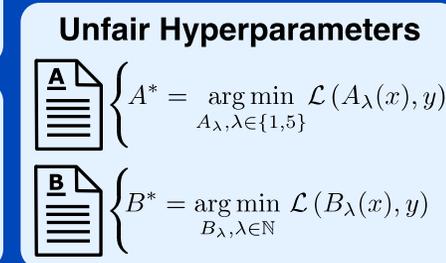
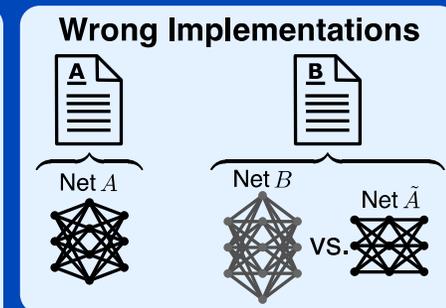
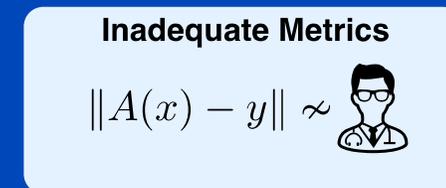
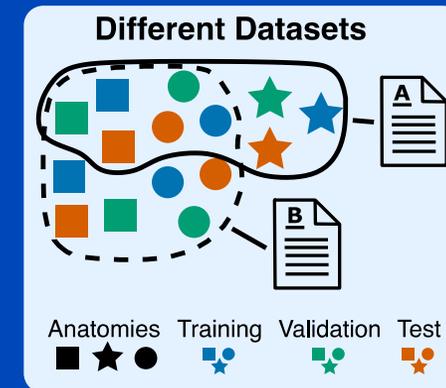
Flaws of Current Evaluation Protocols

No open-source implementations: Increases chance of (un)intentional errors

Meanwhile, to keep the reasonable model complexity, we reduced 96 filters to 32 filters in each layer. [1]

Inadequate metrics: Standard IQA metrics (MSE, SSIM, ...) do not correlate well with human reader ratings [2]

Unfair hyperparameters: Either no HP optimization or limited to subset of parameters / methods. Often authors use HPs reported in reference publications
→ Problematic since no consensus dataset exists



Benchmark Setup

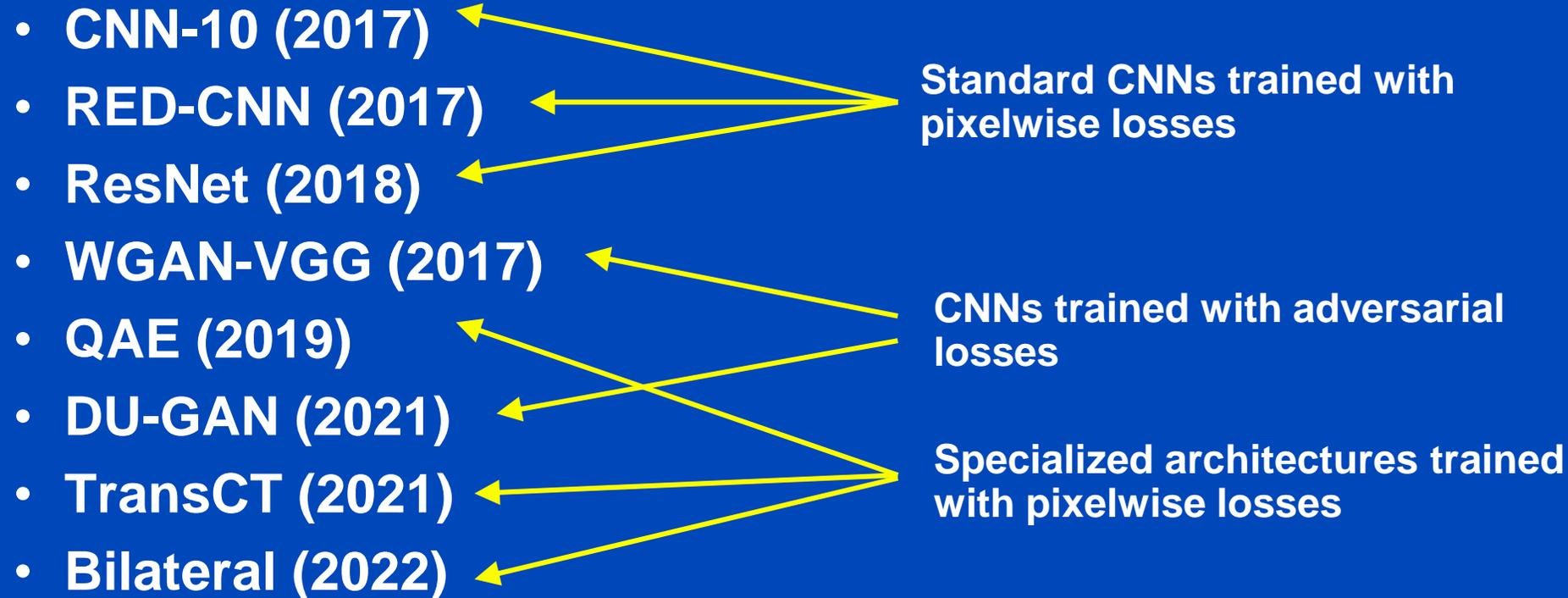
Dataset

Open-source *LDCT and Projection Dataset* [1] with 150 CT scans of abdomen, head, and chest at routine dose levels. Low dose images were simulated at 25% dose for abdomen/head and 10% dose for chest

All tested methods use the same train/validation set and were evaluated on the same test set

Benchmark Setup

LDCT Denoising Algorithms



Benchmark Setup

- **Hyperparameter Optimization**

- Rigorous HP optimization including the weighting factors of loss function terms
- 50 iterations of sequential model-based optimization (SMBO) using Gaussian processes and expected improvement as acquisition function
- As metric to optimize we use SSIM on validation dataset

- **Retrain each method 10 times with**

- optimal HPs and
- different random seeds

	Parameter	Prior
All algorithms	Learning rate	$\log \mathcal{U}(1 \times 10^{-5}, 0.01)$
	Maximum iterations	$\mathcal{U}(1 \times 10^3, 1 \times 10^5)$
	Mini-batch size	$\mathcal{U}(2, 128)$
CNN-10 (2017)	Patchsize	$\mathcal{U}(32, 128)$
RED-CNN (2017)	Patchsize	$\mathcal{U}(32, 128)$
WGAN-VGG (2017)	β_1 of Adam	$\mathcal{U}(0.3, 0.9)$
	Loss weight: $\lambda_{\text{perceptual}}$	$\mathcal{U}(0, 1)$
	Critic updates	$\mathcal{U}(1, 5)$
	Patchsize	$\mathcal{U}(32, 128)$
ResNet (2018)	Patchsize	$\mathcal{U}(32, 128)$
QAE (2019)	Patchsize	$\mathcal{U}(32, 128)$
DU-GAN (2021)	β_1 of Adam	$\mathcal{U}(0.3, 0.9)$
	Cutmix warmup	$\mathcal{U}(0, 1 \times 10^4)$
	Loss weight: λ_{adv}	$\mathcal{U}(0, 1)$
	Loss weight: λ_{CM}	$\mathcal{U}(0, 10)$
	Loss weight: $\lambda_{\text{px,grad}}$	$\mathcal{U}(0, 40)$
TransCT (2021)	Critic updates	$\mathcal{U}(1, 5)$
	Patchsize	$\mathcal{U}(32, 128)$
	-	-
	-	-
Bilateral (2022)	Learning rate for σ_r	$\log \mathcal{U}(1 \times 10^{-5}, 0.01)$
	Patchsize	$\mathcal{U}(32, 128)$
	Initialization for σ_r	$\mathcal{U}(0, 1)$
	Initialization for $\sigma_{x,y}$	$\mathcal{U}(0, 1)$

\mathcal{U} : Uniform distribution; $\log \mathcal{U}$: Log-uniform distribution

Benchmark Setup

Metrics

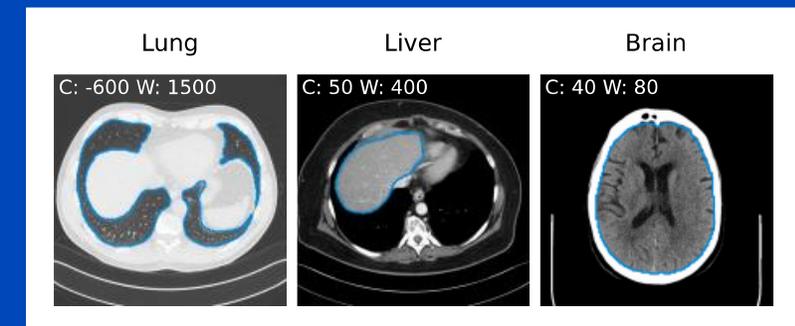
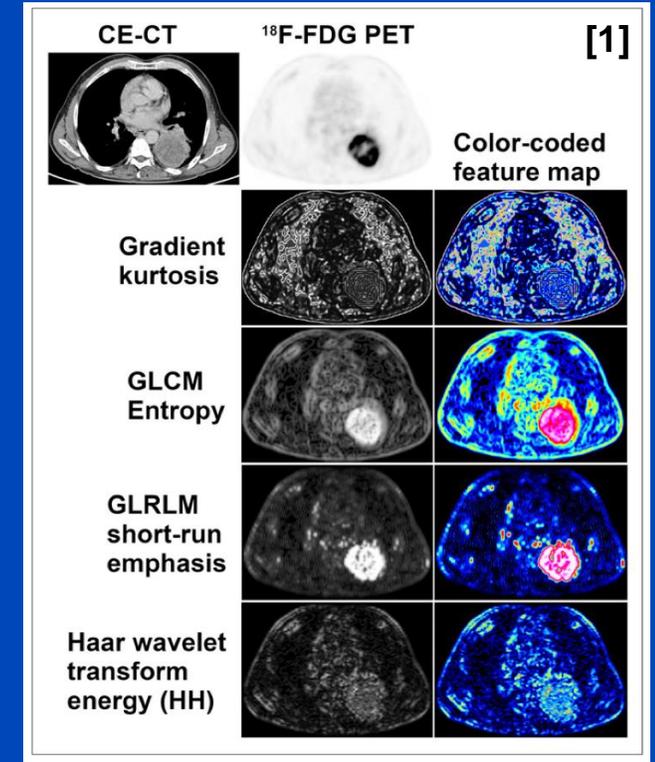
- Standard metrics: **SSIM, PSNR, Visual information fidelity (VIF)**
- Clinically relevant image properties: **Radiomic feature similarity (RFS)**
 1. Automatically segment organs in high-dose scan (s)
 2. Compare features on high-dose scan with those on denoised scans

$$\text{RFS}_i^{(s)} = \cos \left(r_i^{(s)}, r_{\text{GT}}^{(s)} \right), \quad r_i^{(s)} = \left(R_{i,1}^{(s)}, \dots, R_{i,J}^{(s)} \right)$$

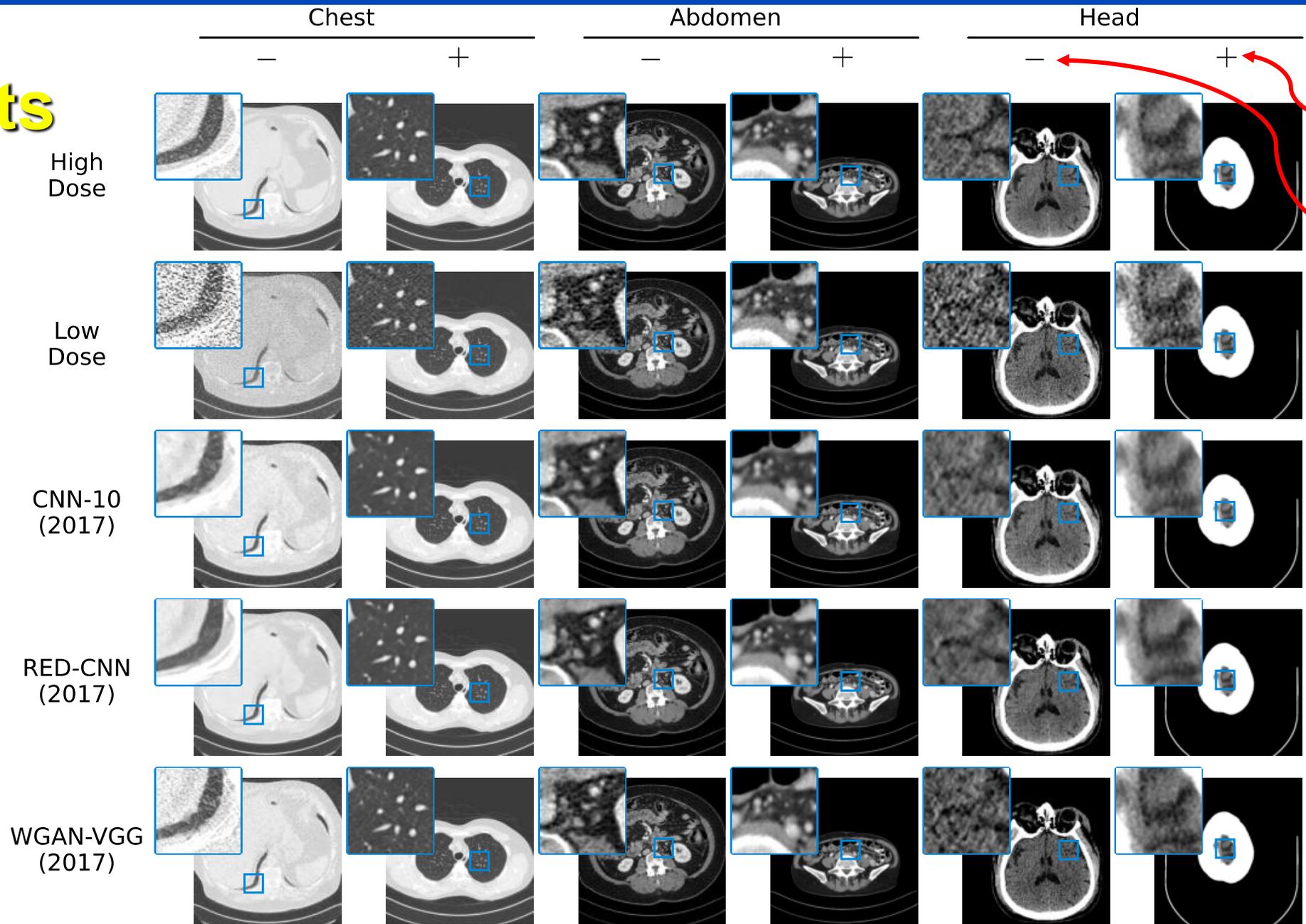
$j \in \{1, 2, \dots, J\}$: Radiomic features

$i \in \{1, 2, \dots, N\}$: Algorithms

Normalized radiomic features for some algorithm i on scan s



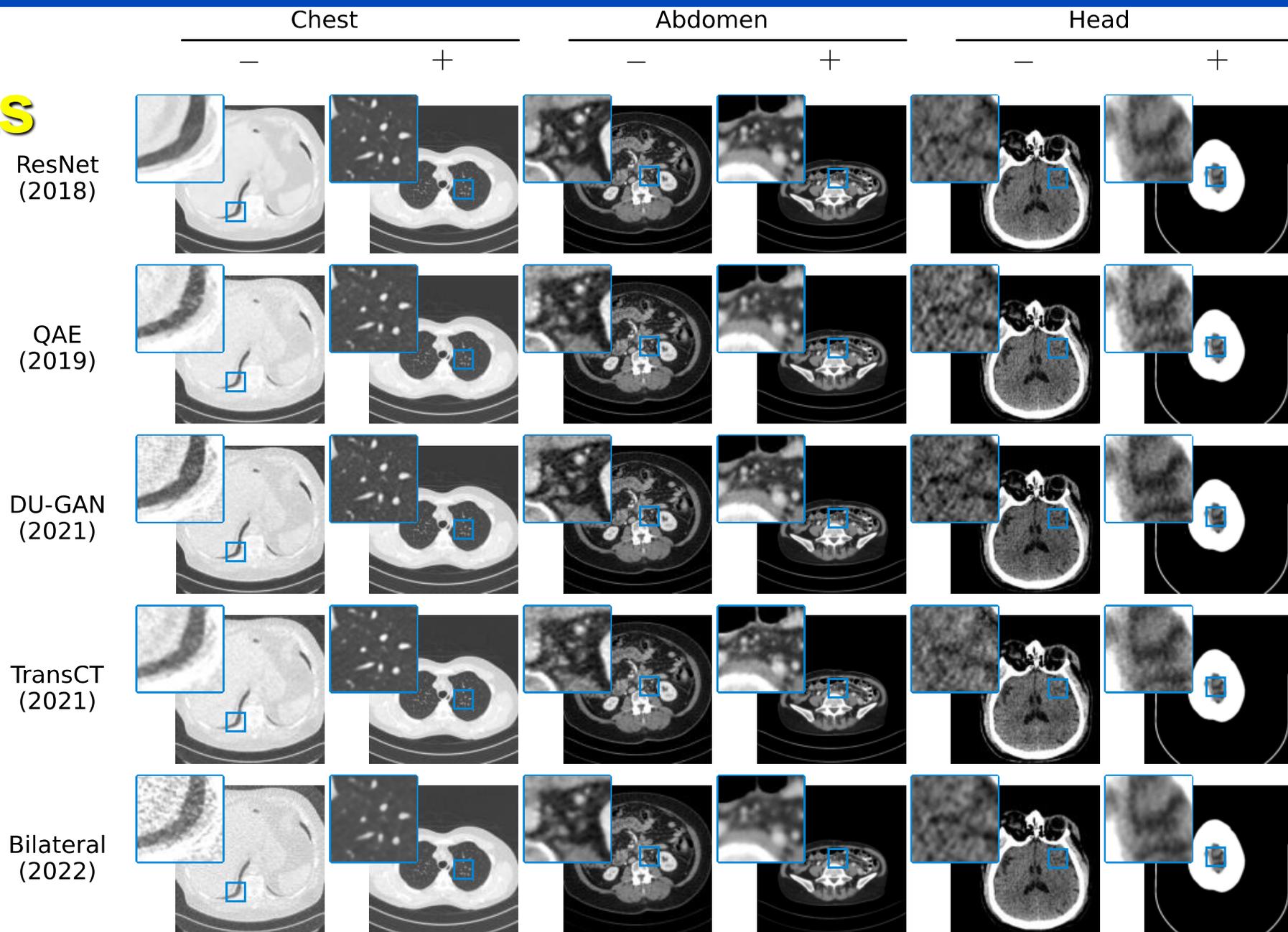
Results



Highest SSIM across all head slices and methods

Lowest SSIM across all head slices and methods

Results



Results

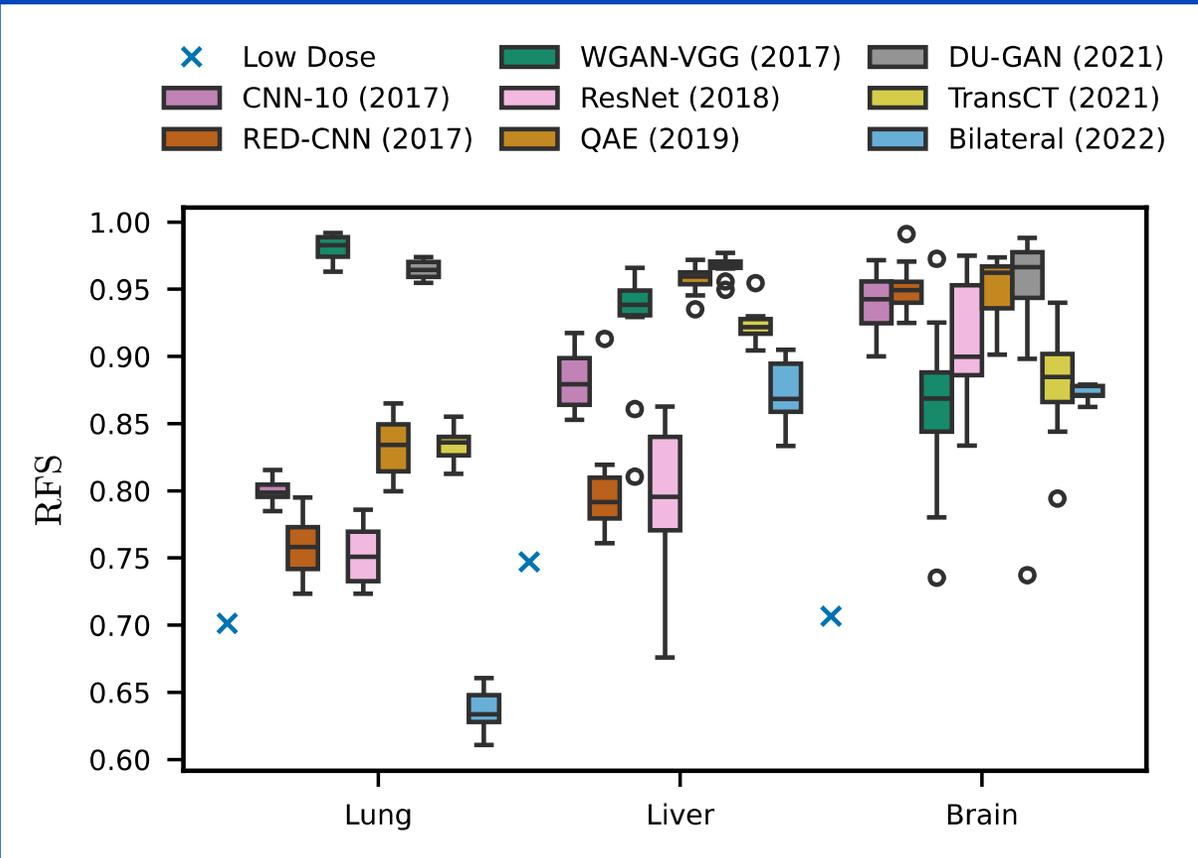
SSIM, PSNR, VIF

	Chest (10% dose)			Abdomen (25% dose)			Head (25% dose)		
	SSIM	PSNR (dB)	VIF	SSIM	PSNR (dB)	VIF	SSIM	PSNR (dB)	VIF
LD	0.34	18.77	0.09	0.84	28.67	0.34	0.88	26.4	0.55
CNN-10 (2017)	0.5867 ± 0.0006	27.71 ± 0.02	0.1915 ± 0.0008	0.896 ± 0.001	32.4 ± 0.1	0.449 ± 0.003	0.896 ± 0.004	28.9 ± 0.6	0.620 ± 0.006
RED-CNN (2017)	0.609 ± 0.002	28.36 ± 0.03	0.221 ± 0.003	0.9028 ± 0.0007	33.22 ± 0.07	0.491 ± 0.008	0.904 ± 0.001	30.4 ± 0.2	0.69 ± 0.01
WGAN-VGG (2017)	<i>0.51 ± 0.03</i>	<i>25.5 ± 0.2</i>	<i>0.148 ± 0.004</i>	<i>0.882 ± 0.002</i>	<i>30.5 ± 0.9</i>	<i>0.38 ± 0.01</i>	<i>0.88 ± 0.02</i>	<i>25 ± 3</i>	<i>0.53 ± 0.02</i>
ResNet (2018)	0.610 ± 0.001	28.42 ± 0.03	0.224 ± 0.002	<i>0.901 ± 0.002</i>	<i>33.15 ± 0.08</i>	<i>0.487 ± 0.006</i>	<i>0.901 ± 0.005</i>	<i>29.6 ± 0.8</i>	<i>0.67 ± 0.02</i>
QAE (2019)	<i>0.584 ± 0.003</i>	<i>27.62 ± 0.09</i>	<i>0.186 ± 0.003</i>	<i>0.894 ± 0.002</i>	<i>32.0 ± 0.2</i>	<i>0.418 ± 0.007</i>	<i>0.899 ± 0.001</i>	<i>28.5 ± 0.3</i>	<i>0.594 ± 0.008</i>
DU-GAN (2021)	<i>0.565 ± 0.004</i>	<i>26.7 ± 0.1</i>	<i>0.168 ± 0.002</i>	<i>0.894 ± 0.002</i>	<i>32.1 ± 0.3</i>	<i>0.427 ± 0.005</i>	<i>0.903 ± 0.003</i>	<i>29 ± 1</i>	<i>0.622 ± 0.005</i>
TransCT (2021)	<i>0.563 ± 0.002</i>	<i>26.99 ± 0.05</i>	<i>0.167 ± 0.002</i>	<i>0.877 ± 0.003</i>	<i>30.5 ± 0.2</i>	<i>0.372 ± 0.007</i>	<i>0.849 ± 0.005</i>	<i>24.7 ± 0.4</i>	<i>0.44 ± 0.01</i>
Bilateral (2022)	<i>0.555 ± 0.001</i>	<i>25.59 ± 0.04</i>	<i>0.159 ± 0.002</i>	<i>0.859 ± 0.003</i>	<i>27.1 ± 0.1</i>	<i>0.361 ± 0.003</i>	<i>0.873 ± 0.002</i>	<i>26.6 ± 0.1</i>	<i>0.500 ± 0.004</i>

Bold: Significantly better than previously best method
Italics: Significantly worse than previously best method

Results

Radiomic Feature Similarity (RFS)



	Lung	Liver	Brain
LD	0.7	0.75	0.71
CNN-10 (2017)	0.800 ± 0.009	0.88 ± 0.02	0.94 ± 0.02
RED-CNN (2017)	0.76 ± 0.02	0.80 ± 0.04	0.95 ± 0.02
WGAN-VGG (2017)	0.98 ± 0.01	0.92 ± 0.05	0.86 ± 0.07
ResNet (2018)	0.75 ± 0.02	0.79 ± 0.06	0.91 ± 0.05
QAE (2019)	0.83 ± 0.02	0.96 ± 0.01	0.95 ± 0.02
DU-GAN (2021)	0.965 ± 0.007	0.967 ± 0.008	0.94 ± 0.08
TransCT (2021)	0.83 ± 0.01	0.92 ± 0.01	0.88 ± 0.04
Bilateral (2022)	0.64 ± 0.01	0.87 ± 0.02	0.873 ± 0.006

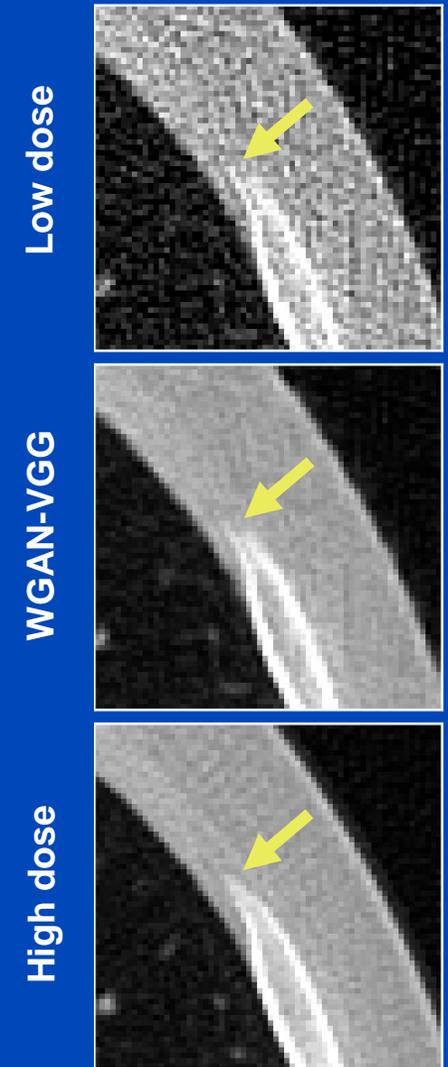
Bold: Significantly better than previously best method
Italics: Significantly worse than previously best method

Summary & Conclusions

- Revisited some of the deep learning-based methods for low dose CT image denoising
- Newer algorithms do not consistently outperform earlier ones both in terms of standard IQA metrics and the proposed radiomic feature similarity

→ Highlights the need for more rigorous and fair evaluation of novel deep learning based denoising methods for LDCT image denoising*

Important research direction: Develop metrics that capture the robustness of algorithms wrt anatomical details



*Similar to reality checks in related fields [1, 2]

[1] G. Melis, C. Dyer, and P. Blunsom, "On the state of the art of evaluation in neural language models," in *ICLR*, 2018.

[2] K. Musgrave, S. Belongie, and S.-N. Lim, "A metric learning reality check," in *ECCV*, 2020.