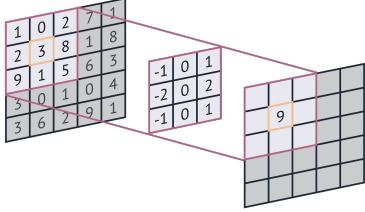
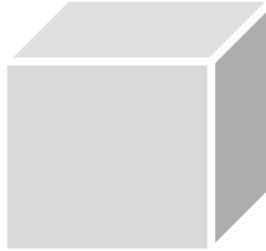


Image segmentation with deep learning

Current state of the art



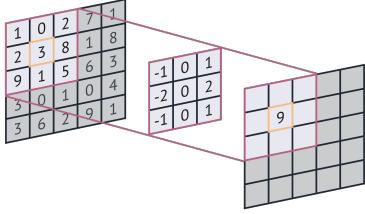
Recap on deep learning for image segmentation



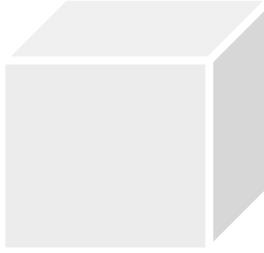
3D image segmentation



nnU-Net



Recap on deep learning for image segmentation



3D image segmentation



nnU-Net

**In deep learning, we hope to find functions that utilise
the structure of our data to perform a given task**

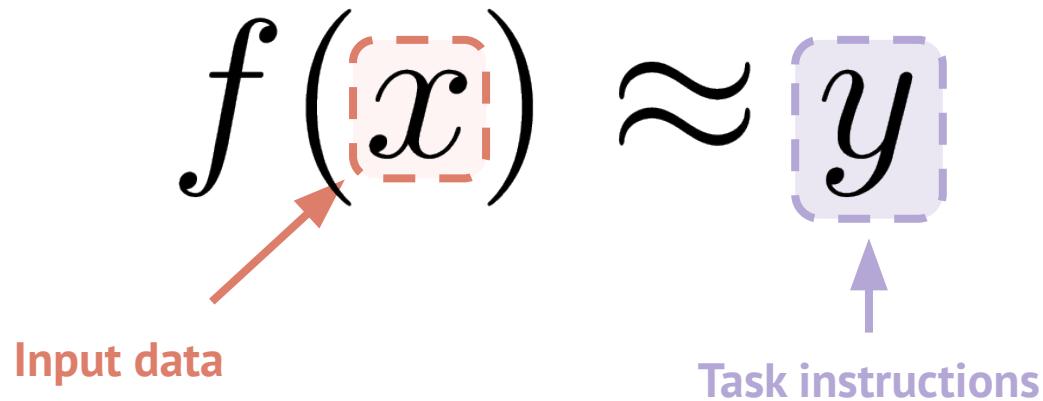
$$f(x) \approx y$$

In deep learning, we hope to find functions that utilise the structure of our data to perform a given task

$$f(\tilde{x}) \approx y$$

Input data

In deep learning, we hope to find functions that utilise the structure of our data to perform a given task



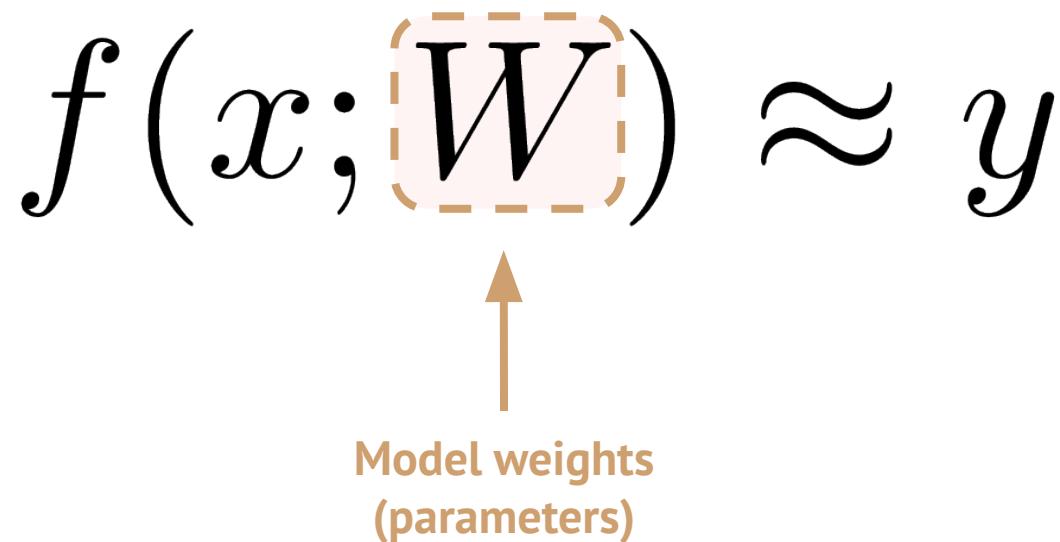
In deep learning, we hope to find functions that utilise the structure of our data to perform a given task



To discover this function, we parametrise it by a set of weights and try to find the optimal weights

$$f(x; W) \approx y$$

To discover this function, we parametrise it by a set of weights and try to find the optimal weights

$$f(x; \mathbb{W}) \approx y$$


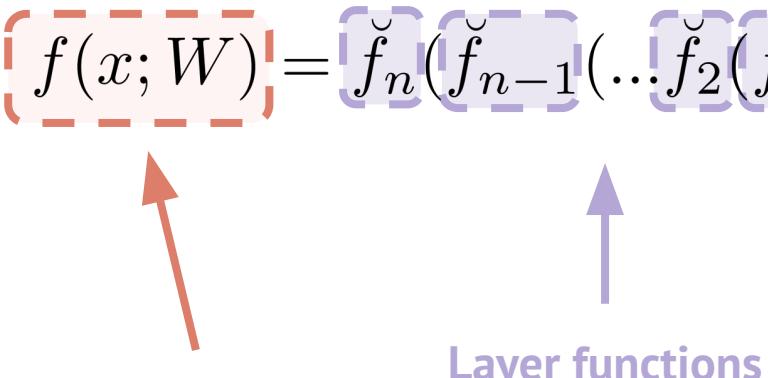
Model weights
(parameters)

Neural networks are described as a cascade of function compositions

$$f(x; W) = \check{f}_n(\check{f}_{n-1}(\dots \check{f}_2(\check{f}_1(x; W_1), W_2)\dots, W_{n-1}), W_n)$$

Neural networks are described as a cascade of function compositions

$$f(x; W) = \check{f}_n(\check{f}_{n-1}(\dots \check{f}_2(\check{f}_1(x; W_1), W_2), \dots, W_{n-1}), W_n)$$



Neural network

Layer functions

Each layer is generally an affine map (linear plus constant) followed by an element-wise nonlinear activation function

$$f(x; W) = \check{f}_n(\check{f}_{n-1}(\dots \check{f}_2(\check{f}_1(x; W_1), W_2), \dots, W_{n-1}), W_n)$$

$$\check{f}_k(x; A, b) = \phi(Ax + b)$$

Each layer is generally an affine map (linear plus constant) followed by an element-wise nonlinear activation function

$$f(x; W) = \check{f}_n(\check{f}_{n-1}(\dots \check{f}_2(\check{f}_1(x; W_1), W_2), \dots, W_{n-1}), W_n)$$

$$\check{f}_k(x; A, b) = \phi(Ax + b)$$



Layer weights (W_k above)

Each layer is generally an affine map (linear plus constant) followed by an element-wise nonlinear activation function

$$f(x; W) = \check{f}_n(\check{f}_{n-1}(\dots \check{f}_2(\check{f}_1(x; W_1), W_2), \dots, W_{n-1}), W_n)$$

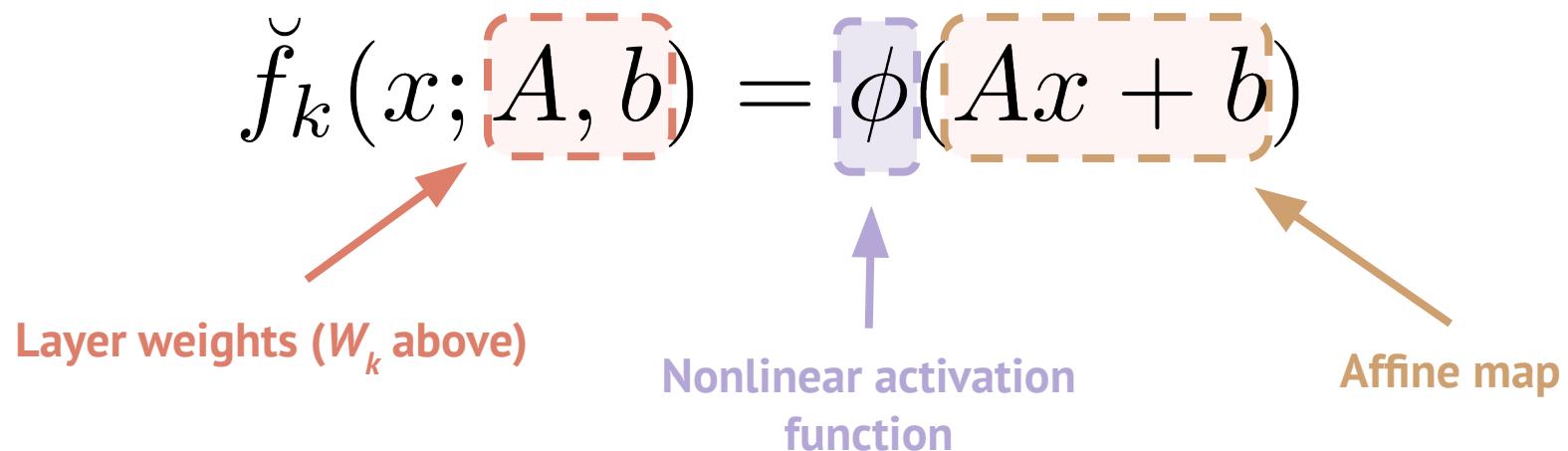
$$\check{f}_k(x; A, b) = \phi(Ax + b)$$

Layer weights (W_k above)

Affine map

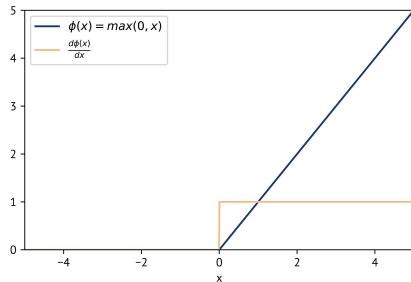
Each layer is generally an affine map (linear plus constant) followed by an element-wise nonlinear activation function

$$f(x; W) = \check{f}_n(\check{f}_{n-1}(\dots \check{f}_2(\check{f}_1(x; W_1), W_2), \dots, W_{n-1}), W_n)$$

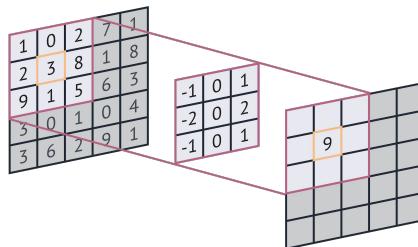


This parametrisation gives rise to a large number of choices when designing a neural network

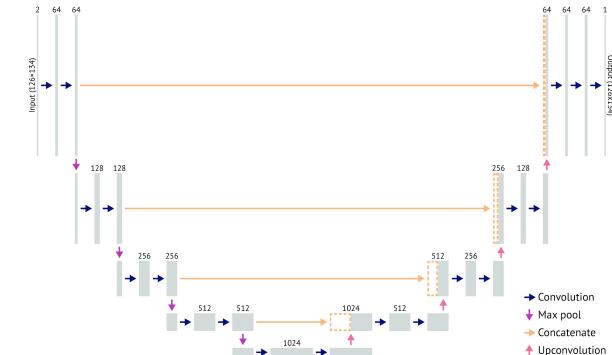
Activation function



Structure of linear map

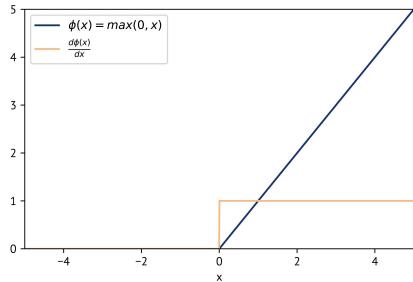


Structure of function composition

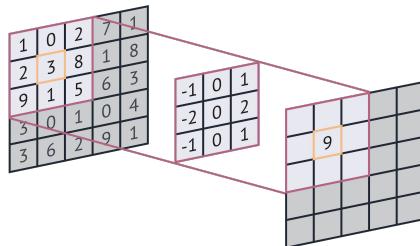


The exact choice of model is called an *architecture*

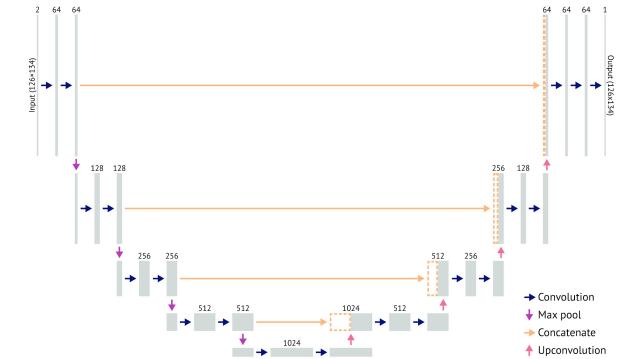
Activation function



Structure of linear map



Structure of function composition



Model architecture

The most commonly used activation function is the ReLU function

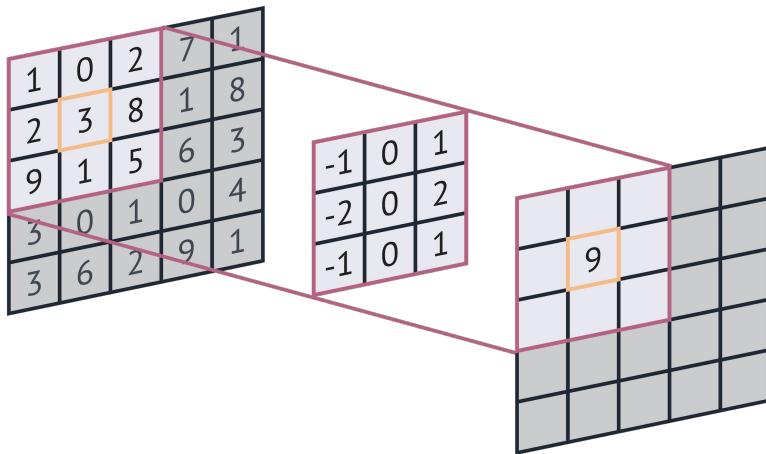
$$\phi(x) = \max(x, 0)$$

However, because of some numerical benefits, the *leaky ReLU* activation function is also often used

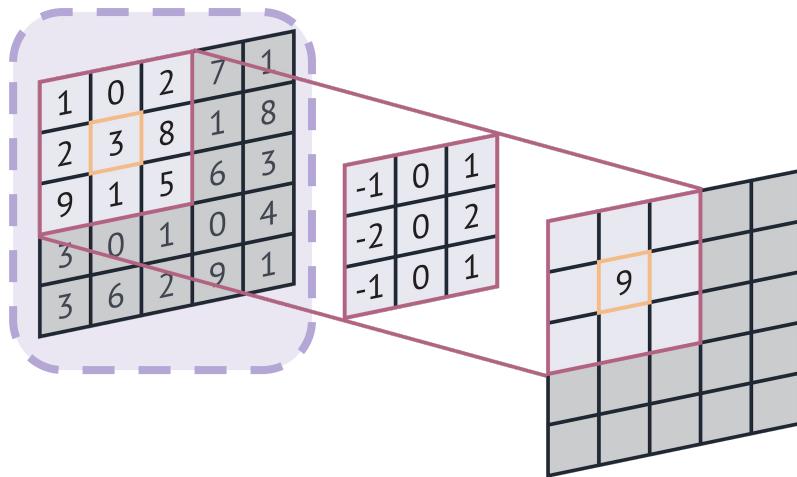
$$\phi(x) = \max(x, 0)$$

$$\phi_a(x) = \max(x, ax)$$

For image segmentation, the linear map is a convolution

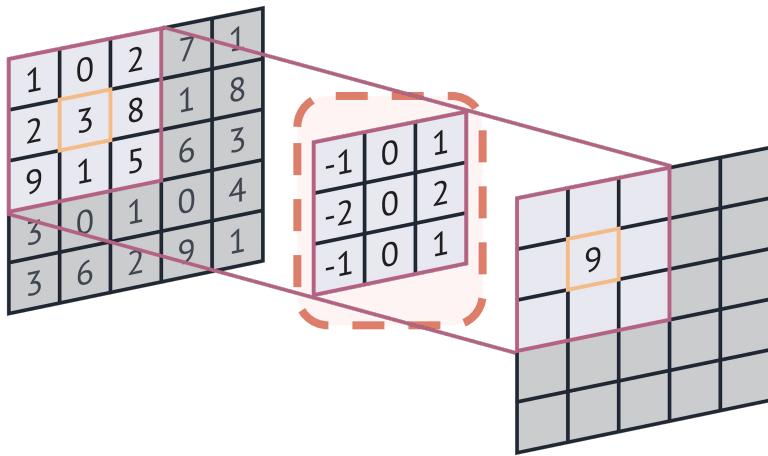


For image segmentation, the linear map is a convolution



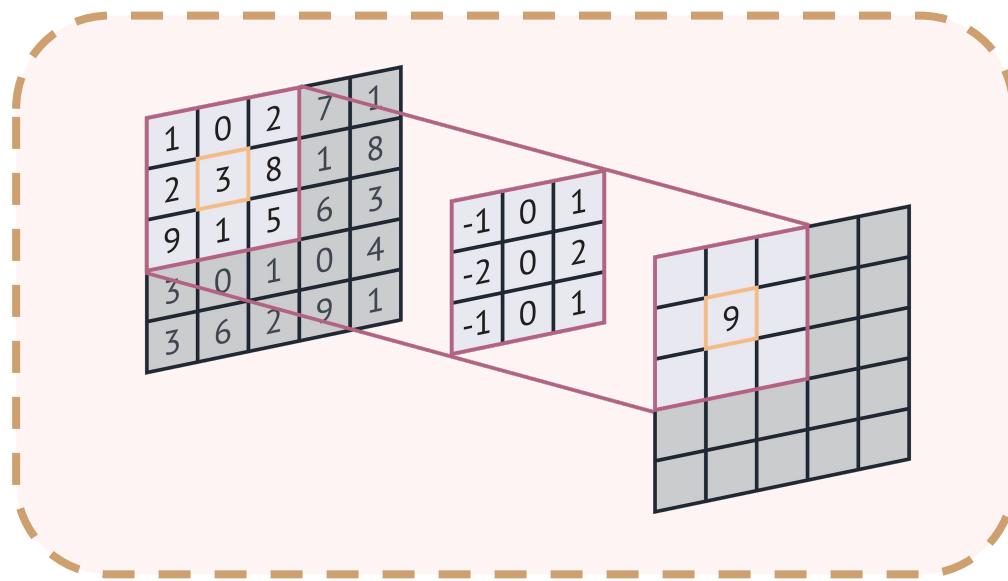
$$\check{f}_k(x; A, b) = \phi(A\underline{x} + b)$$

For image segmentation, the linear map is a convolution



$$\check{f}_k(x; A, b) = \phi(\underline{A}x + b)$$

For image segmentation, the linear map is a convolution



$$\check{f}_k(x; A, b) = \phi(\underline{Ax} + b)$$

To illustrate convolution concepts, we will look at 1D convolutions

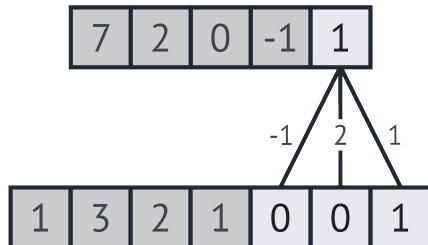
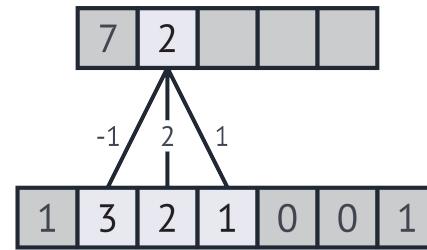
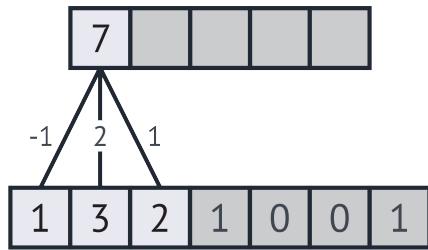
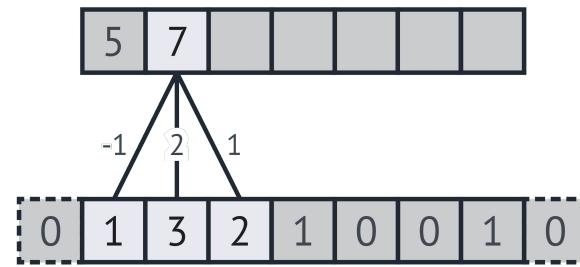
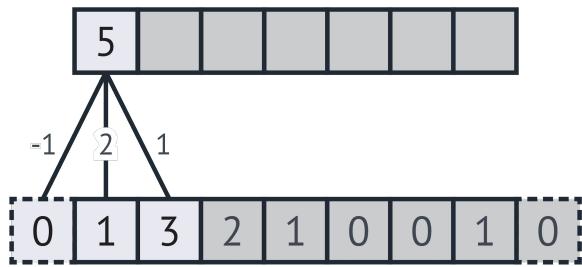
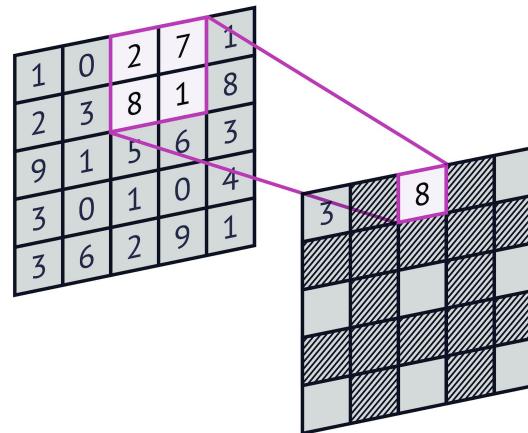
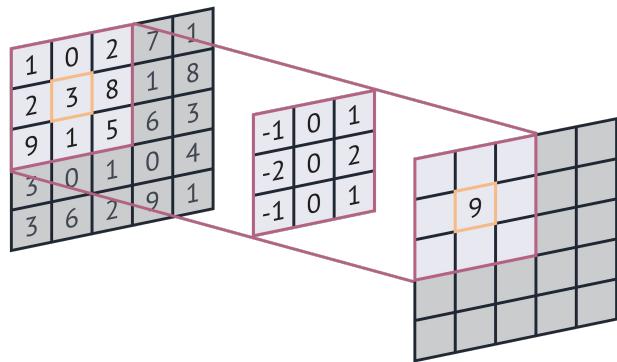


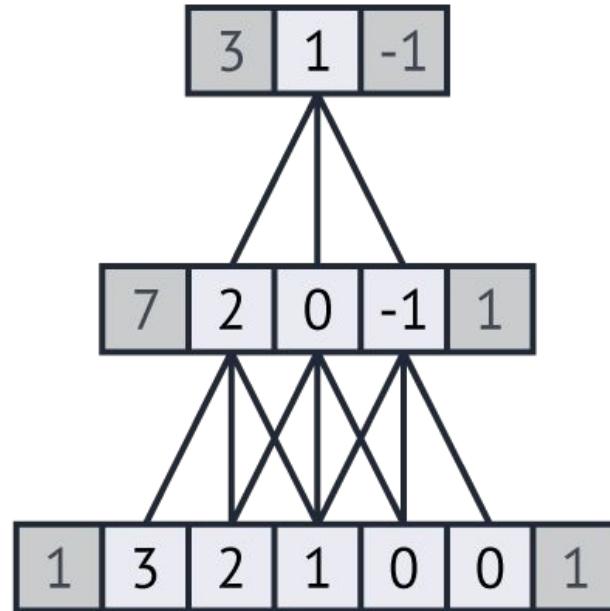
Image boundaries are often padded by zeros to keep image size constant



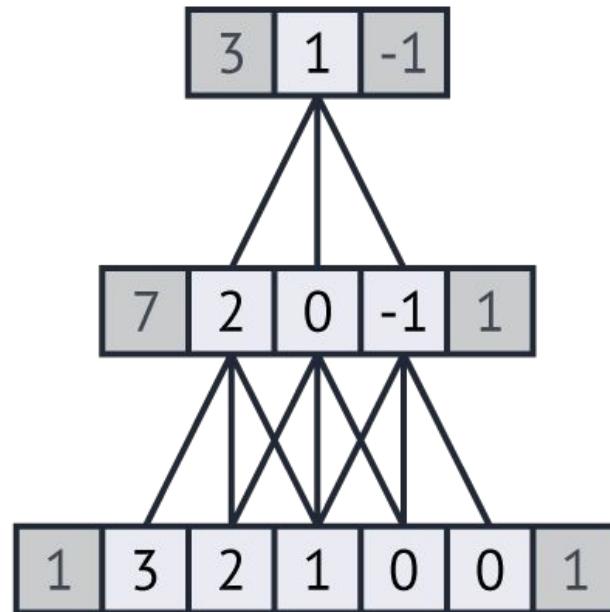
A convolutional neural network is a neural network where all layers can be sliding window operators



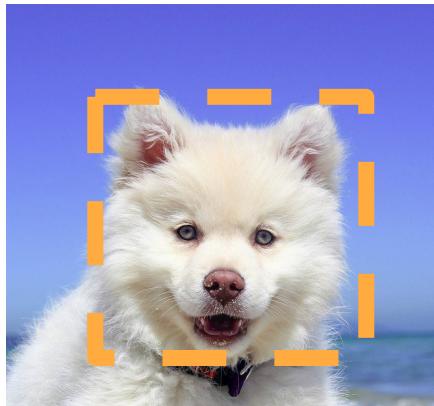
An important concept for convolutional neural networks is *receptive fields*



The receptive field of a layer is the size of the input image that affects it

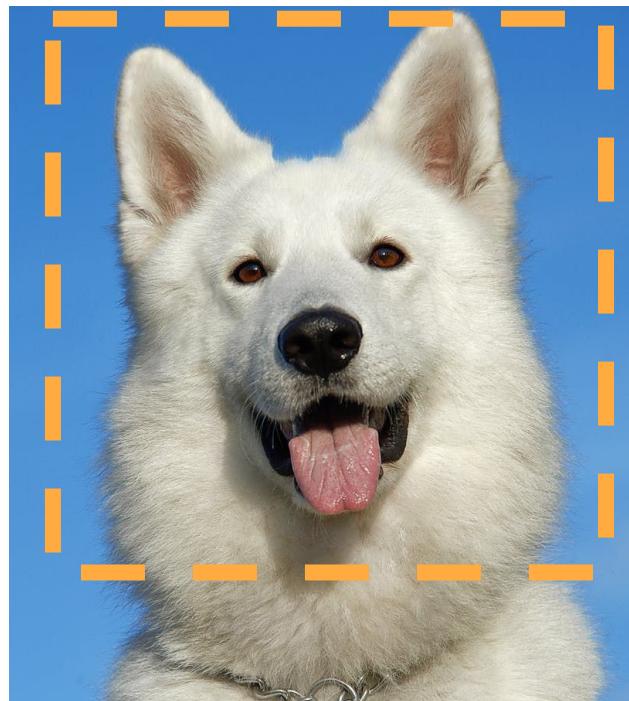


The receptive field describes how large features we are able to discover in the input image

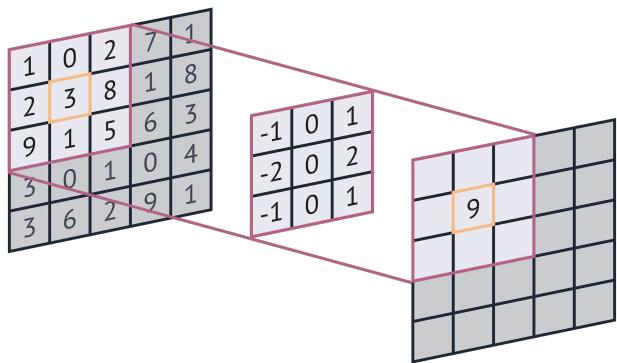


Small receptive
field

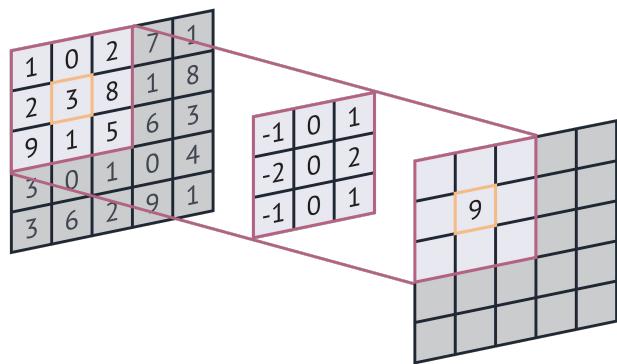
Large receptive
field



A single 3×3 convolution has a receptive field of 3×3



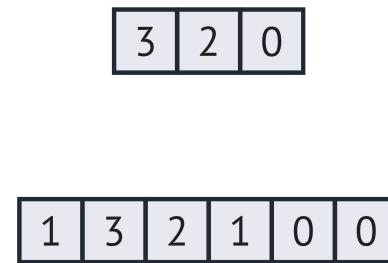
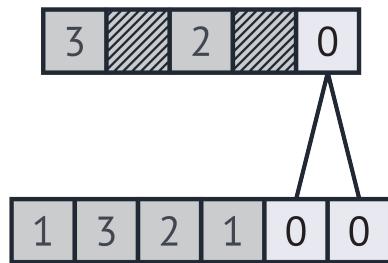
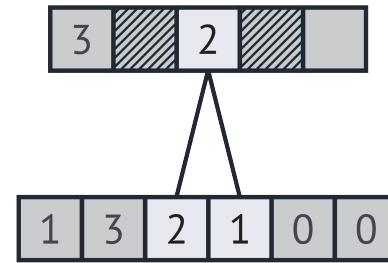
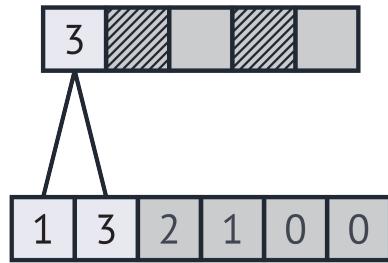
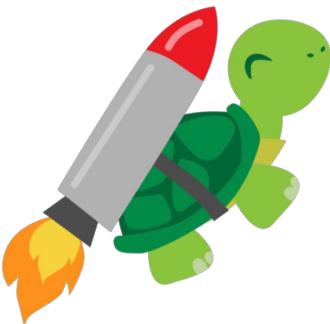
Two stacked 3x3 convolutions has a receptive field of 5x5



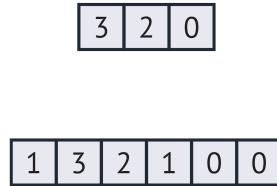
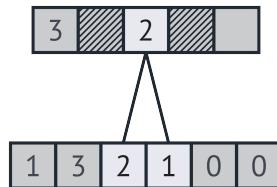
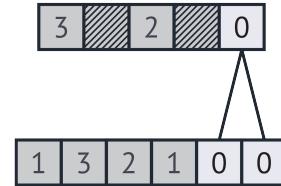
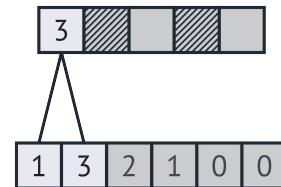
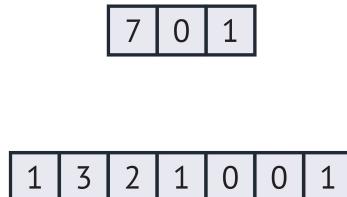
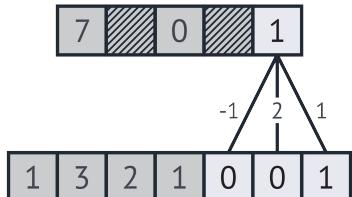
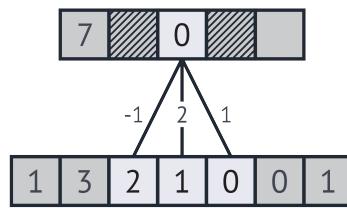
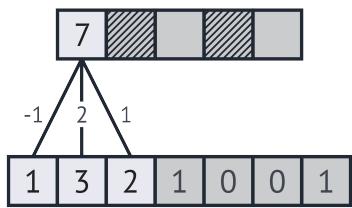
If we just stack convolutions, then the receptive field grows slowly



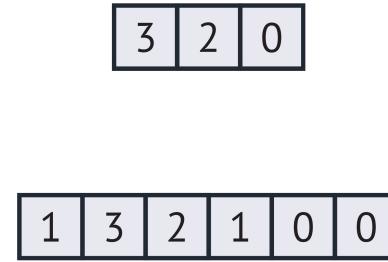
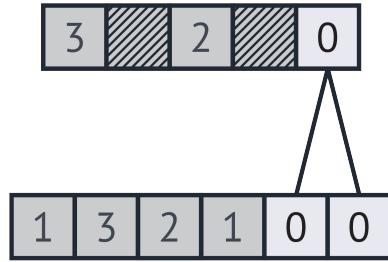
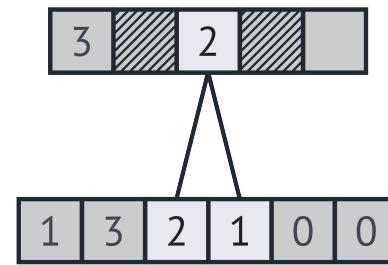
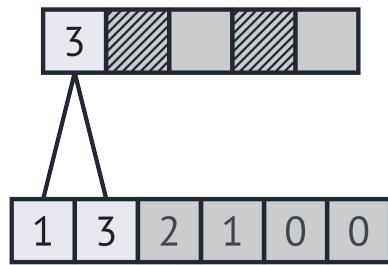
To combat this, we generally use downsampling, or *pooling* layers



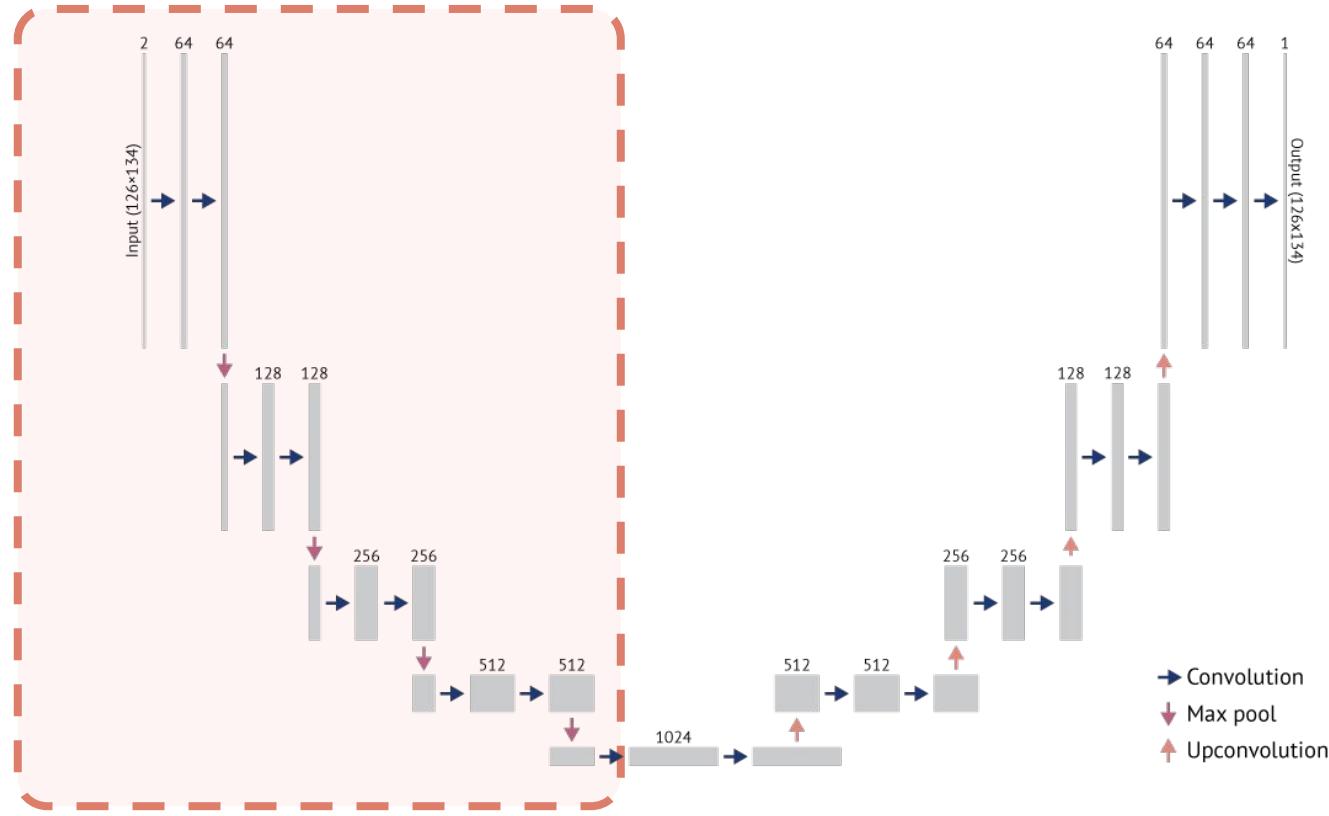
Downsampling is generally either performed via a max-pooling or a strided convolution



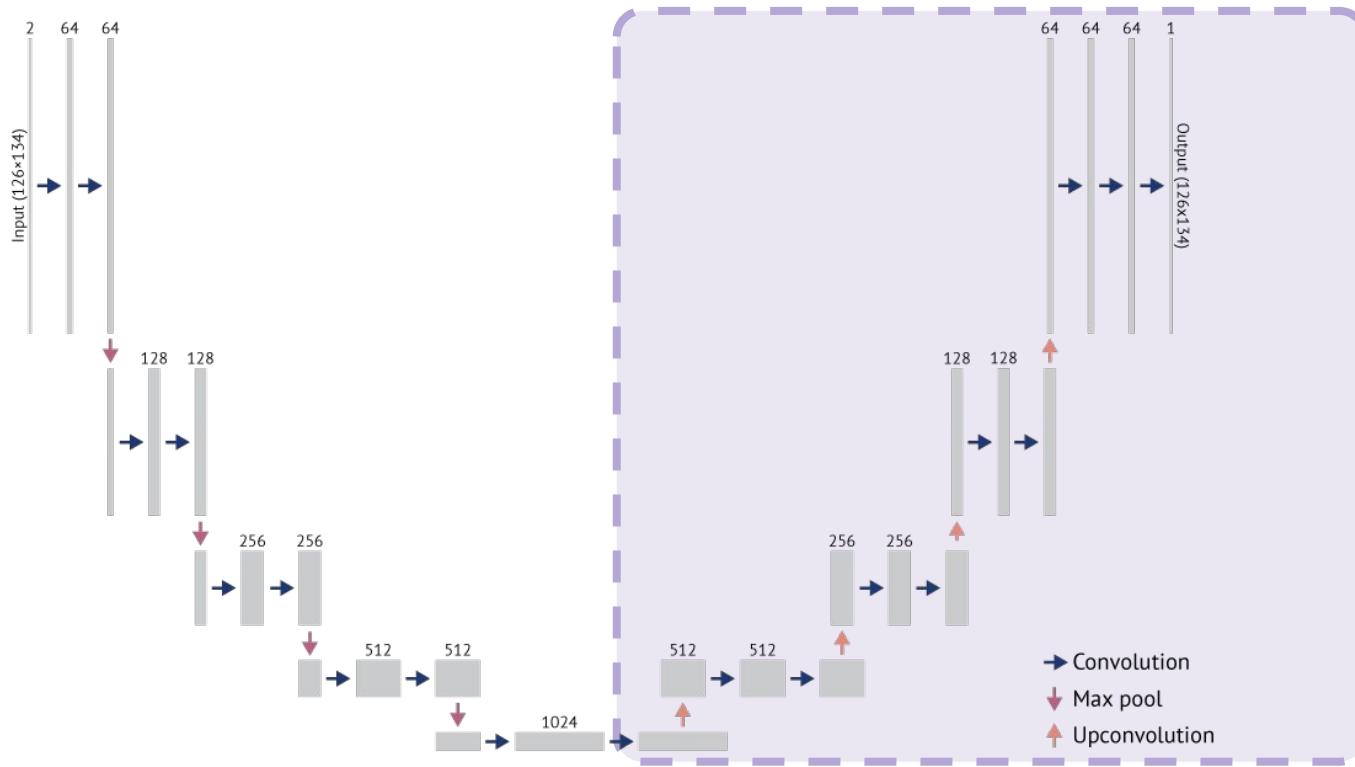
Downsampling reduces the image size, so the output of the network will be smaller than the input



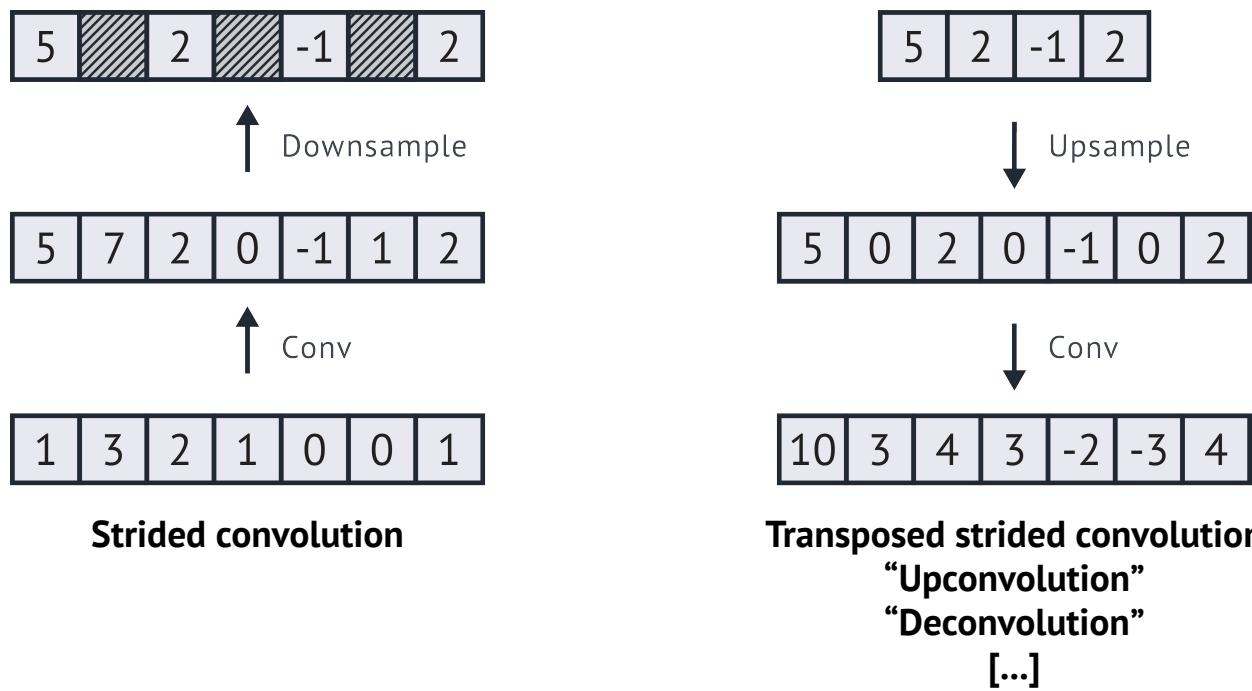
The encoder down-samples the image, allowing us to discover large features in the input image



The decoder upscales the output from the encoder, allowing us to locate where each object is in the image



This upscaling is normally done via a transposed strided convolution



Upconvolution can lead to checkerboard artefacts, so some authors prefer to use interpolation instead

5		2		-1		2
---	--	---	--	----	--	---

↑ Downsample

5	7	2	0	-1	1	2
---	---	---	---	----	---	---

↑ Conv

1	3	2	1	0	0	1
---	---	---	---	---	---	---

Strided convolution

5	2	-1	2
---	---	----	---

↓ Upsample

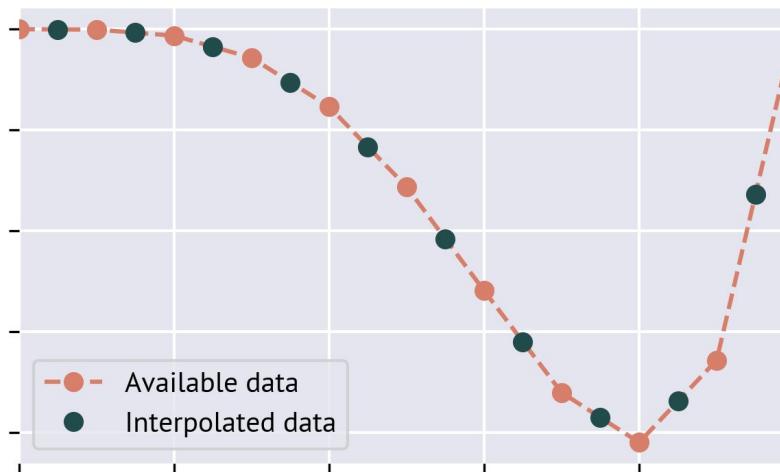
5	0	2	0	-1	0	2
---	---	---	---	----	---	---

↓ Conv

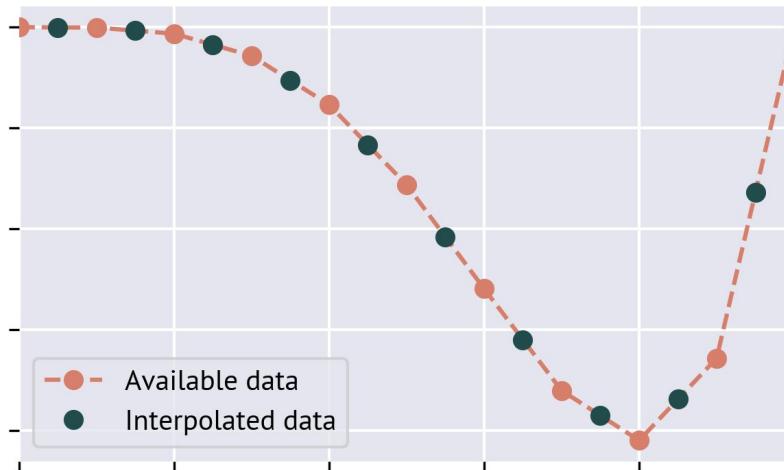
10	3	4	3	-2	-3	4
----	---	---	---	----	----	---

Upconvolution

If you image has odd dimensions before downscaling, then a linear interpolation is needed



If you image has odd dimensions before downscaling, then a linear interpolation is needed



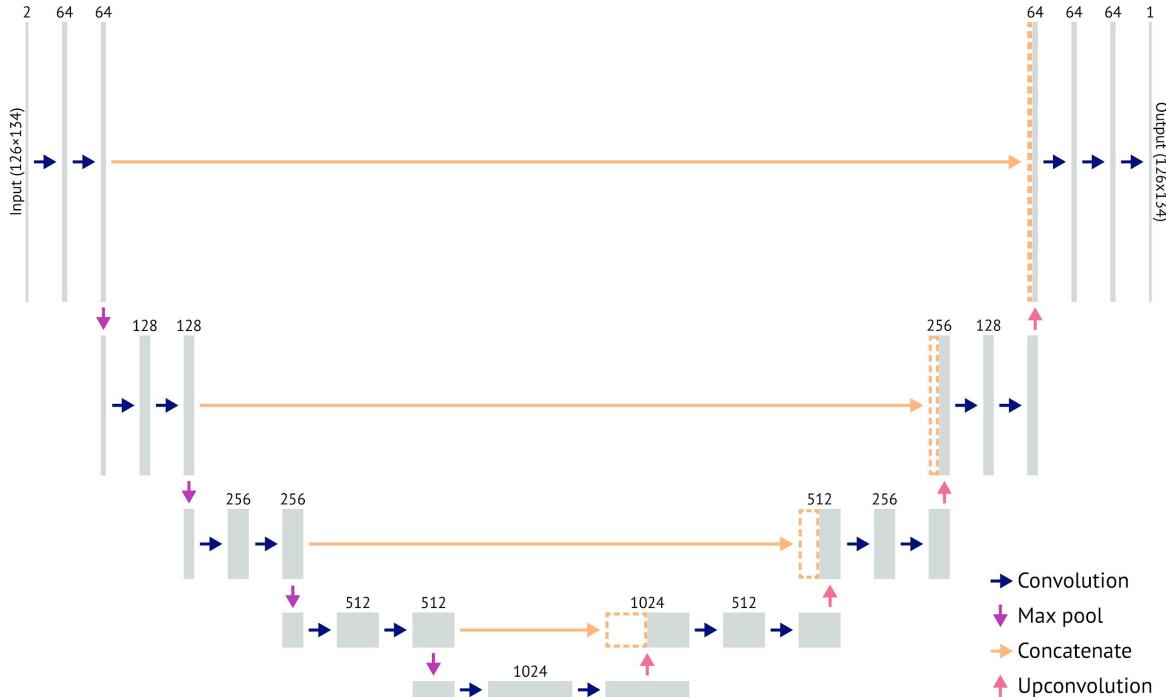
Cubic interpolation cannot run on GPU (at least in TensorFlow)

Unfortunately, the downsampling procedure discards high-frequency information

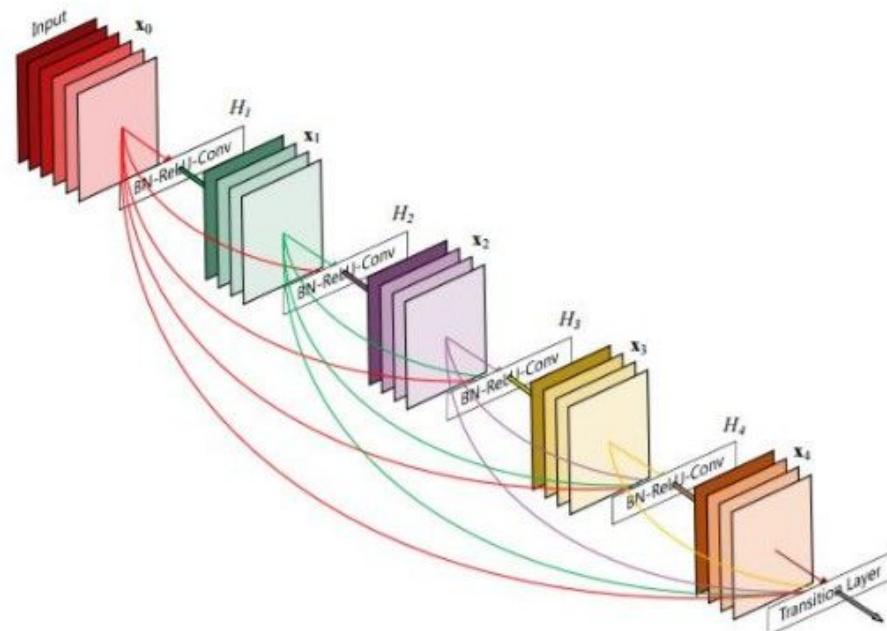
Where is the tumour edge?



The U-Net and SegNet architectures were designed to preserve this information



Other architectures also exist, however, U-Net is regarded as a good choice



Now that we have decided on an architecture, how do we fit these models?



We formulate a loss function that measures the severity of our prediction errors

$$L(W; \mathcal{T}) = \sum_{\substack{(x^{(i)}, y^{(i)}) \in \mathcal{T}}} l(f(x^{(i)}; W), y^{(i)})$$

Diagram illustrating the components of the loss function:

- Training example**: Points to the term $(x^{(i)}, y^{(i)})$ in the sum.
- Training set**: Points to the set \mathcal{T} in the sum.
- Loss function**: Points to the function $l(\cdot)$ enclosed in a dashed orange box.

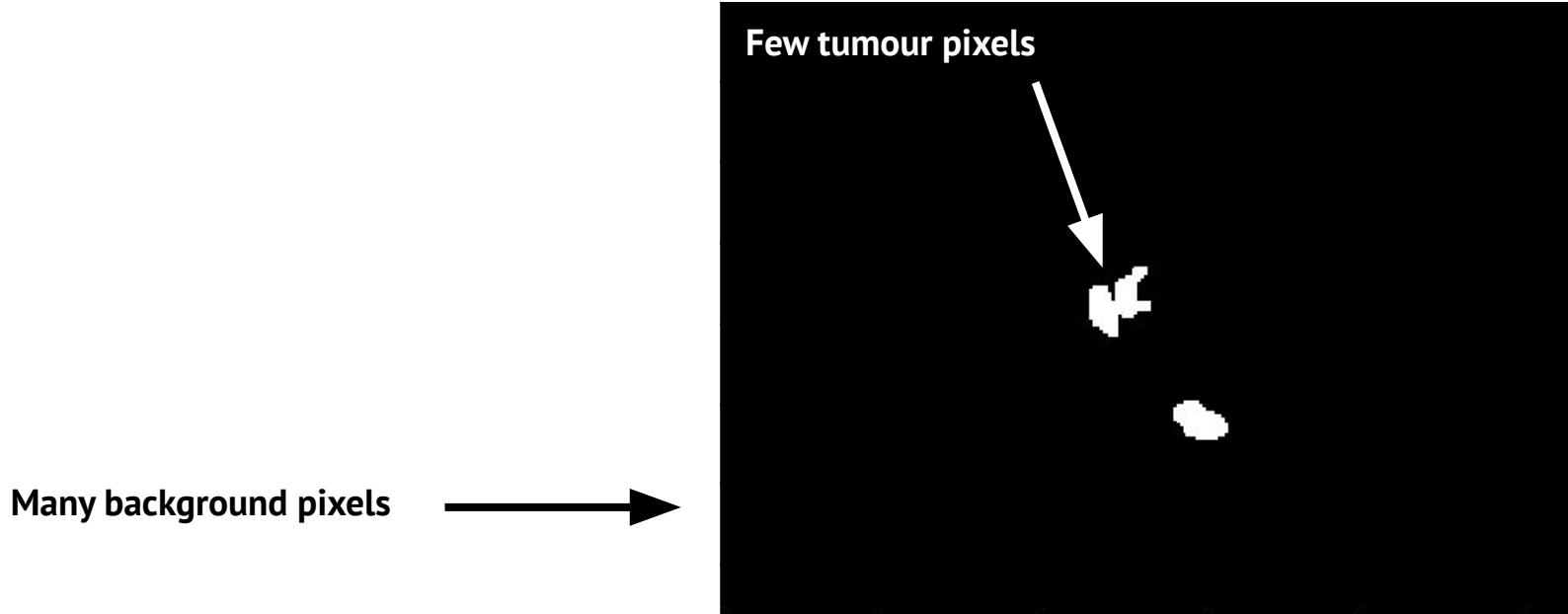
For classification problems, we often use the cross entropy error

$$l_{\text{CE}}(f(x^{(i)}; W), y^{(i)}) = - \sum_p \left[y_p^{(i)} \log \left(f \left(x_p^{(i)} \right) \right) + \left(1 - y_p^{(i)} \right) \log \left(1 - f \left(x_p^{(i)} \right) \right) \right]$$

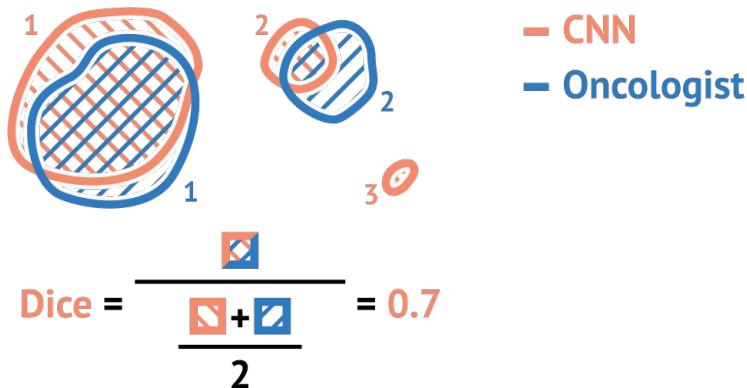
True pixel class

Predicted pixel class

However, this loss function puts too much weight into correctly classifying background pixels

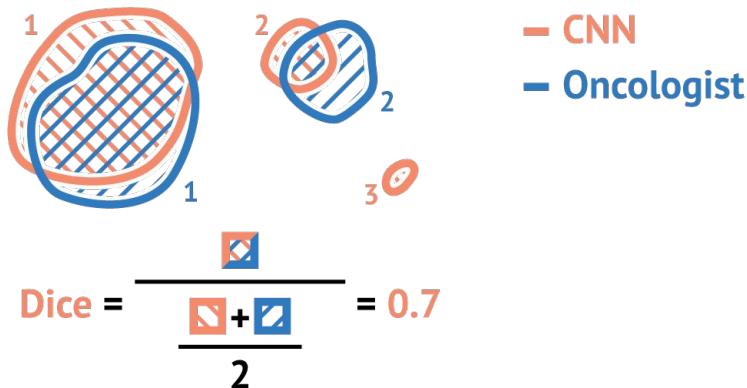


Therefore, we often use the Dice loss instead



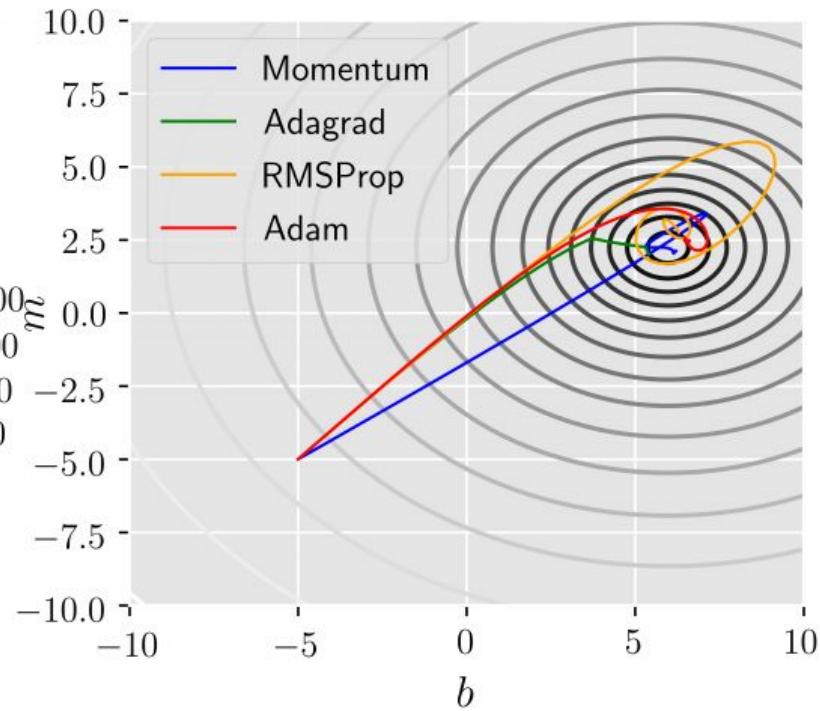
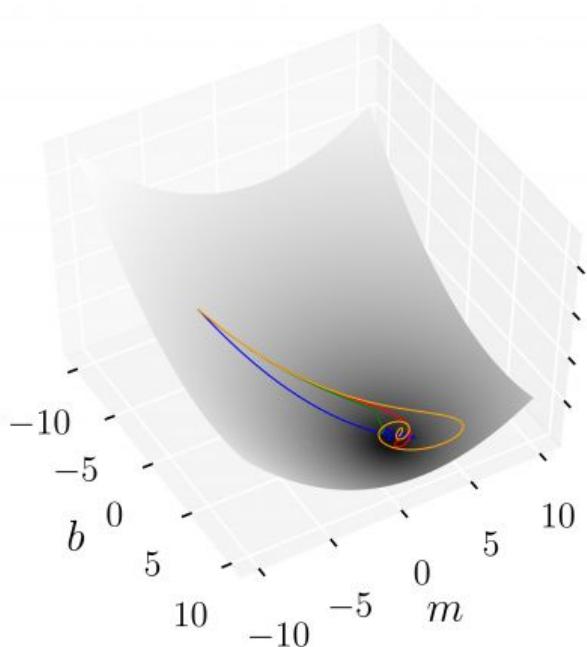
$$l_{\text{dice}}(f(x^{(i)}; W), y^{(i)}) = \frac{2 \sum_p \left[f(x_p^{(i)}) y_p^{(i)} \right]}{\sum_p f(x_p^{(i)}) + \sum_p y_p^{(i)}}$$

Therefore, we often use the Dice loss instead

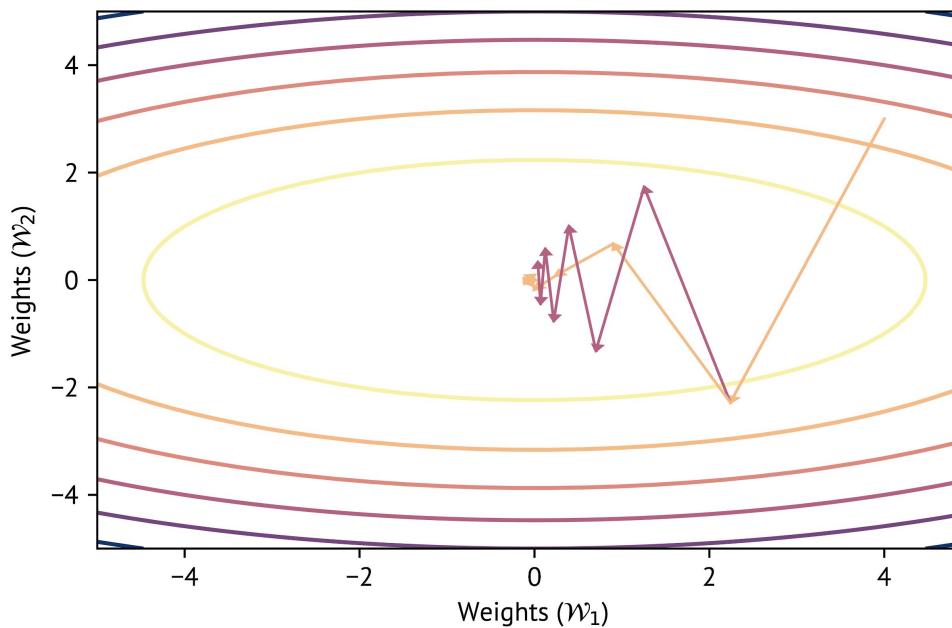


$$l_{\text{dice}}(f(x^{(i)}; W), y^{(i)}) = \frac{2 \sum_p \left[f(x_p^{(i)}) y_p^{(i)} \right]}{\sum_p f(x_p^{(i)}) + \sum_p y_p^{(i)}}$$

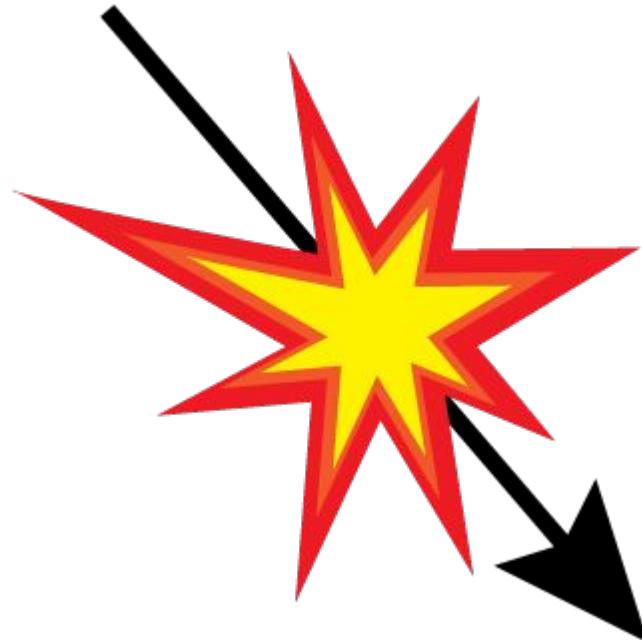
First minimise the loss using an adaptive momentum method such as Adam

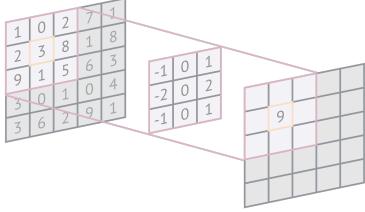


If the generalisation error is large, try using momentum SGD optimiser instead of an adaptive momentum method

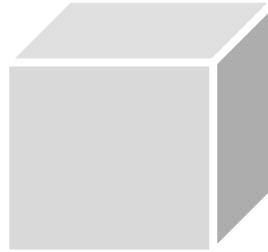


And always remember to normalise your layers to avoid exploding gradients!





Recap on deep learning for image segmentation



3D image segmentation



nnU-Net

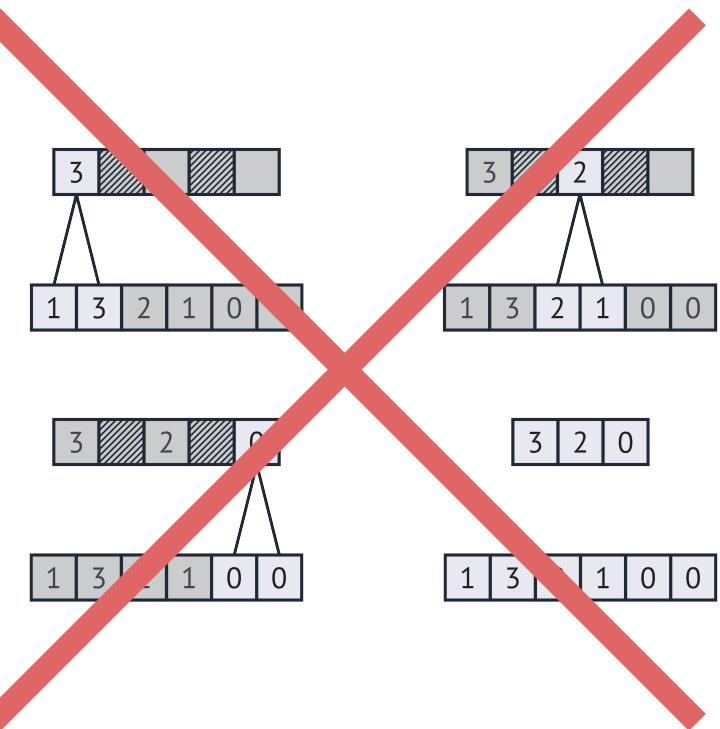
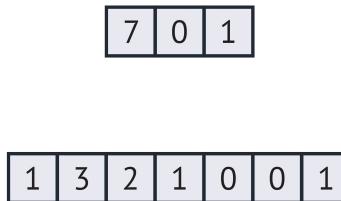
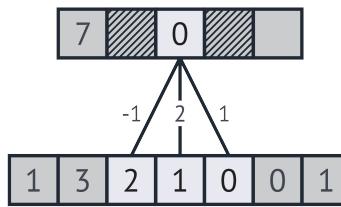
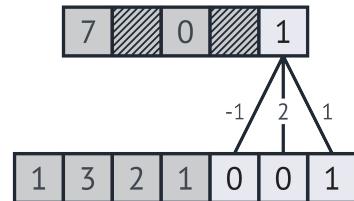
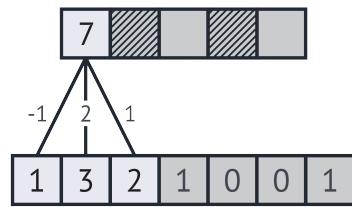
With 3D images, the memory footprint grows significantly



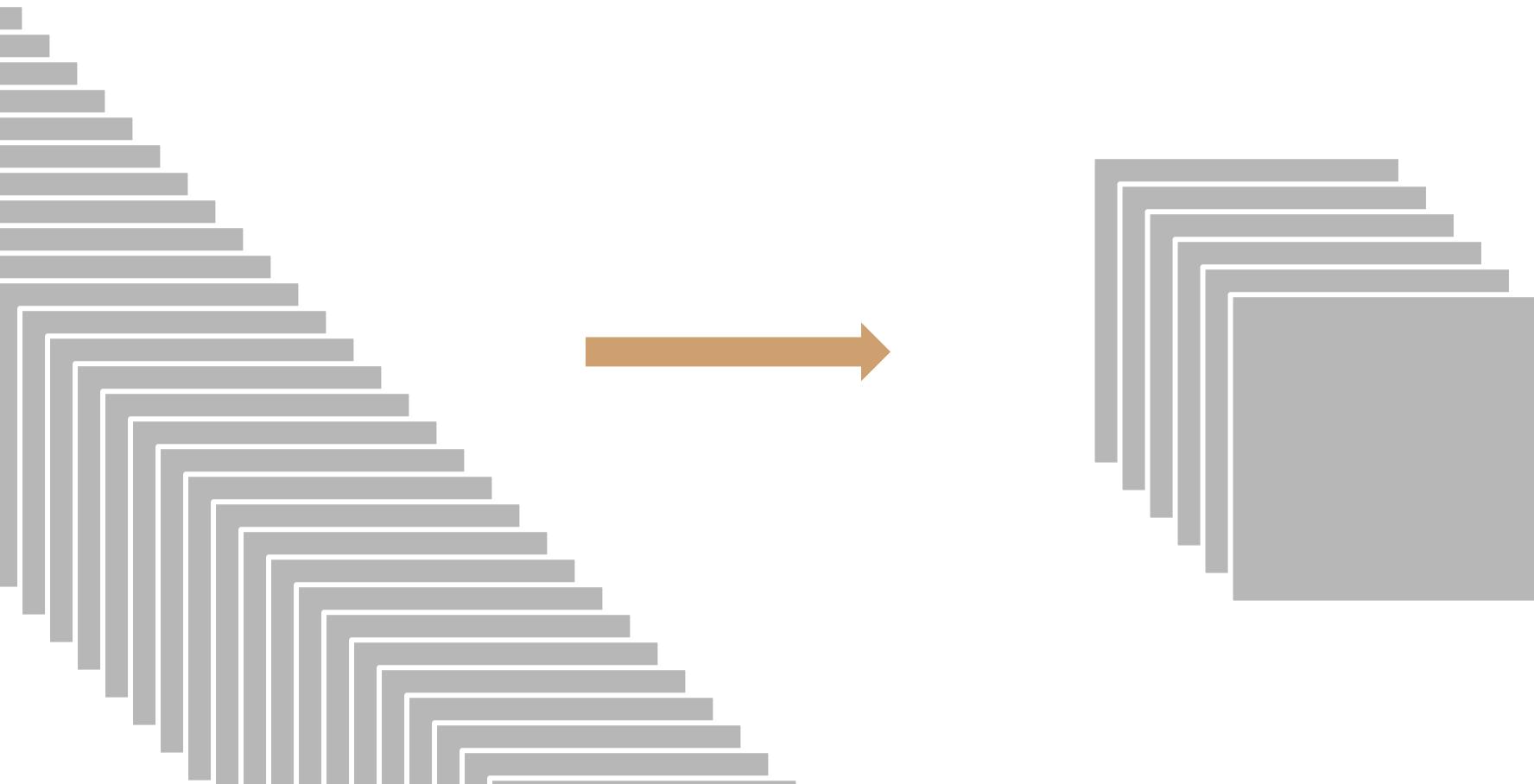
**Therefore, we need to use a shallower network than in 3D
to fit the model in memory**

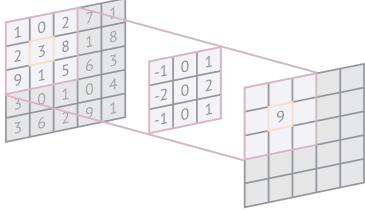


You should consider strided convolutions for 3D instead of max-pooling to reduce memory footprint

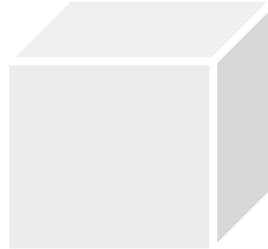


You also need to reduce the number of images per batch





Recap on deep learning for image segmentation



3D image segmentation



nnU-Net

nnU-Net: Self-adapting Framework for U-Net-Based Medical Image Segmentation



<https://arxiv.org/pdf/1809.10486.pdf>

The idea is that the U-Net architecture is good enough, only small tweaks to the original architecture is needed

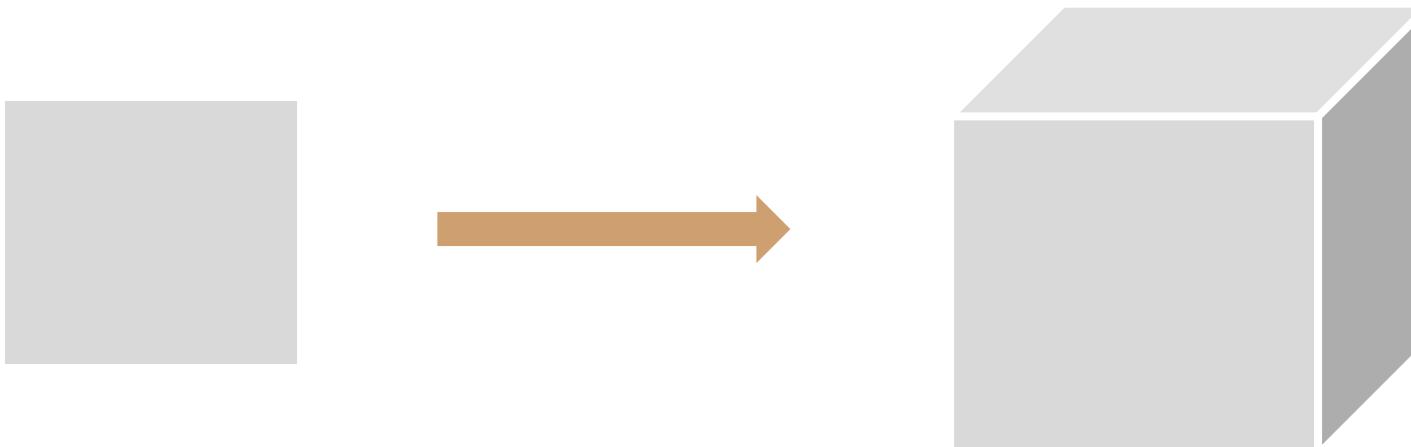
ReLU → LeakyReLU ($\alpha=0.01$)

BatchNormalization → InstanceNormalization

Dice loss → Dice + cross entropy loss

Raw image input → Random data augmentation

For 3D images, they replaced 2D convolutions with 3D convolutions



The depth of the models in the nnU-Net architecture is adaptive to the dataset size

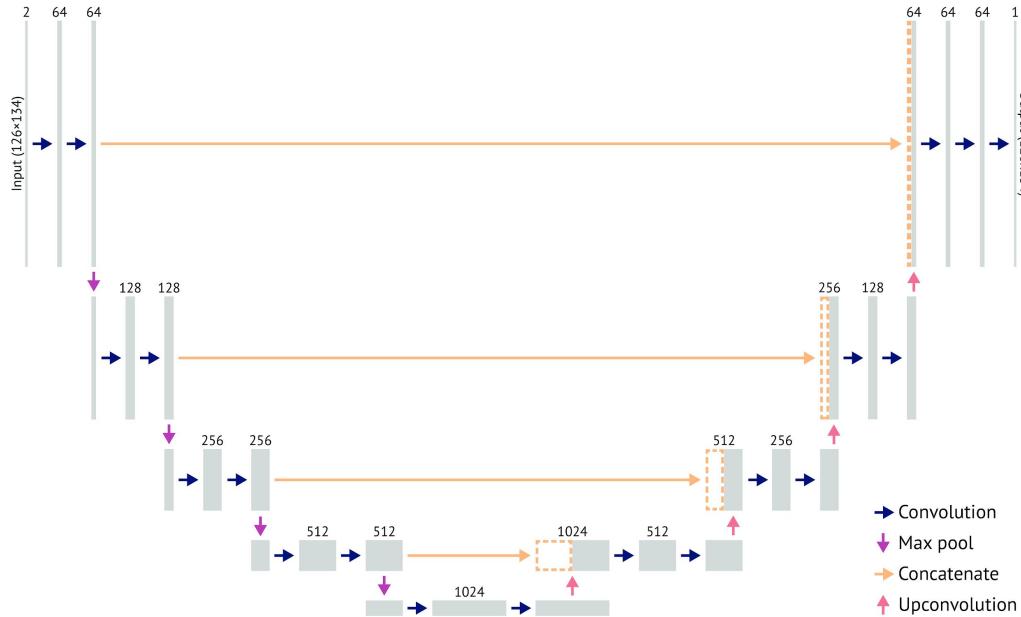
Architecture specification:

- Two convolutions between each up/down-sample operation
- Downsample until image has size 8 in each dimension

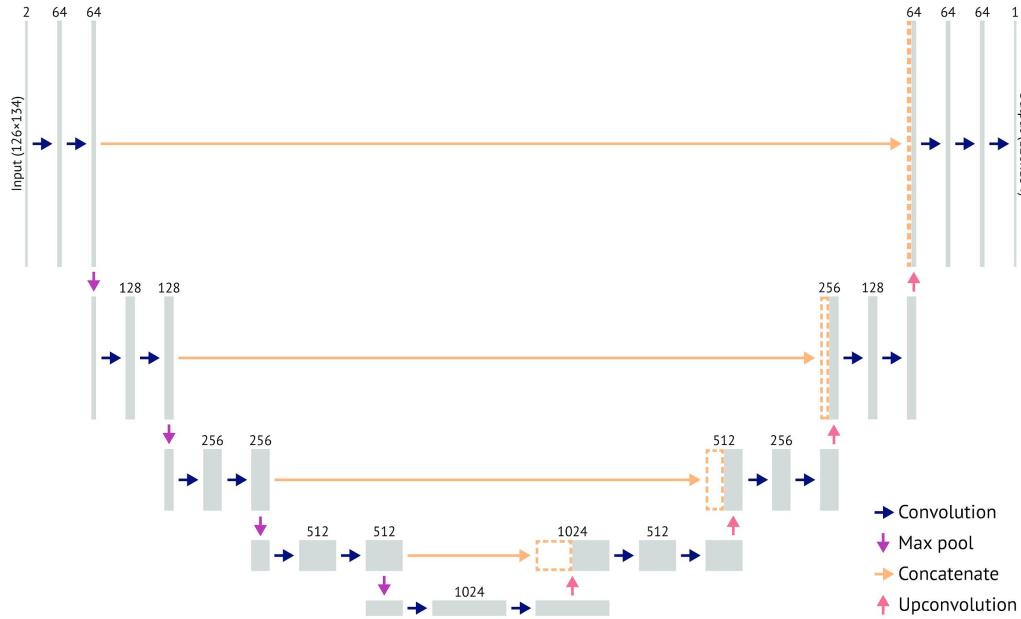
In a later version of the preprint, they show that there is generally not a large effect of any hyperparameter choice

	BraTS	Liver lowres	Liver fullres	Hippocampus	Prostate	Lung nodule	Pancreas
Vanilla nnU-Net	0.72	0.79	0.78	0.89	0.77	0.65	0.65
Batch norm instead of Inst. norm	1.0%	-0.1%	2.9%	-0.1%	-1.3%	-14.2%	-3.7%
No feature map normalization	1.1%	-4.6%	-22.8%	-0.2%	-4.2%	3.0%	-100.0%
ReLU instead of LeakyReLU	0.6%	0.0%	1.0%	-0.1%	-0.2%	-0.4%	0.5%
No data augmentation	-0.8%	-4.9%	1.5%	-1.5%	-0.4%	4.2%	-11.3%
Only cross-entropy loss	-0.6%	-12.0%	-6.3%	0.0%	-1.4%	-25.4%	-8.8%
Only dice loss	0.9%	-2.5%	-10.1%	-0.3%	-3.0%	-11.5%	1.6%

In conclusion, don't use anything fancier than U-Net for segmentation tasks



In conclusion, don't use anything fancier than U-Net for segmentation tasks



Except possibly using CRFs as a post-processing step