

Deep Learning in Computer Vision

June 7th, 2019

Janus Nørtoft Jensen, inje@dtu.dk

Topics of this lecture

- Better optimization
-

Optimization

Learning rate annealing

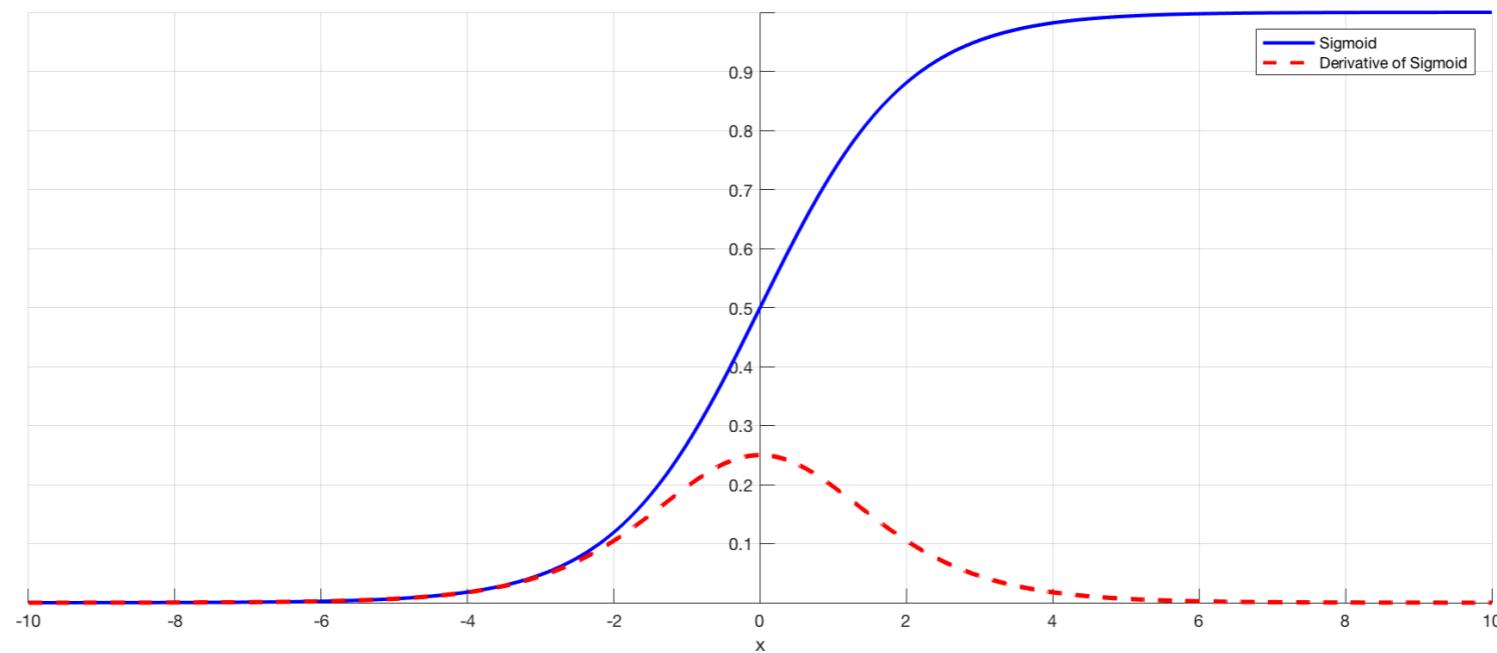
- Often it is helpful to change to learning rate over time when training deep neural networks.
 - If we are far from a minimum we maybe want a high learning rate
 - If we are close to a minimum we want a smaller learning rate.
- If the learning rate is to large the the updated parameters will just bounce around and not reach a minimum
- If it is to small it takes ‘forever’ to reach a minimum

Better optimisation

- A lot have been done to do faster training of deep neural networks
- The gradients in a deep neural networks are unstable
 - Many gradients are multiplied together
 - Vanishing gradient / exploding gradient
 - Better optimisation methods

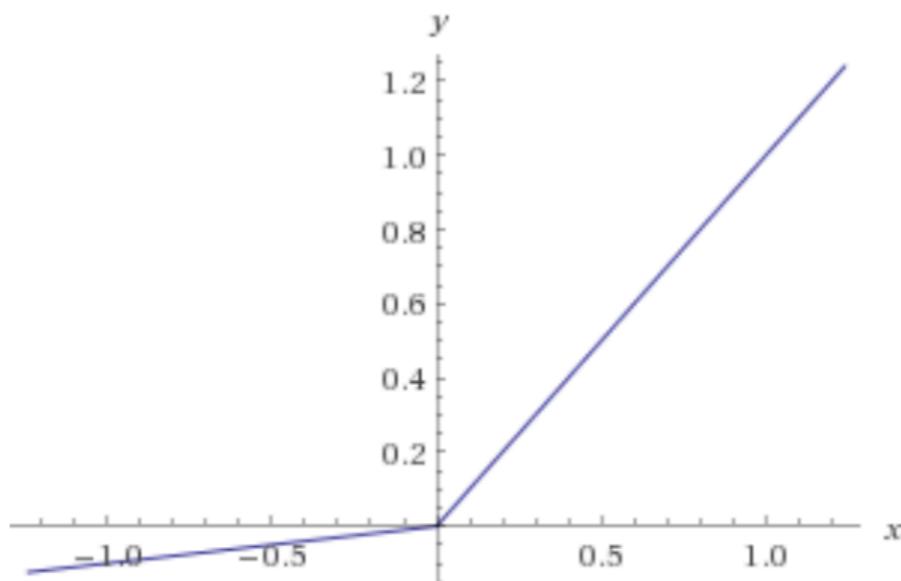
Vanishing gradients

- The gradients in the early layers of a network is the product of the gradients from later layers
 - If we use a sigmoid activation function we may see the ‘vanishing gradient problem’



Vanishing gradient

- Solution?
 - Use ReLU activation function
- Drawbacks
 - No gradient below 0
 - Not differentiable in 0
 - Usually not a problem
- Sometimes people use a ‘Leaky ReLU’ to overcome this



Unstable gradients

- In general, deep neural networks have unstable gradients and layers may not learn at the same rate.

Input normalization

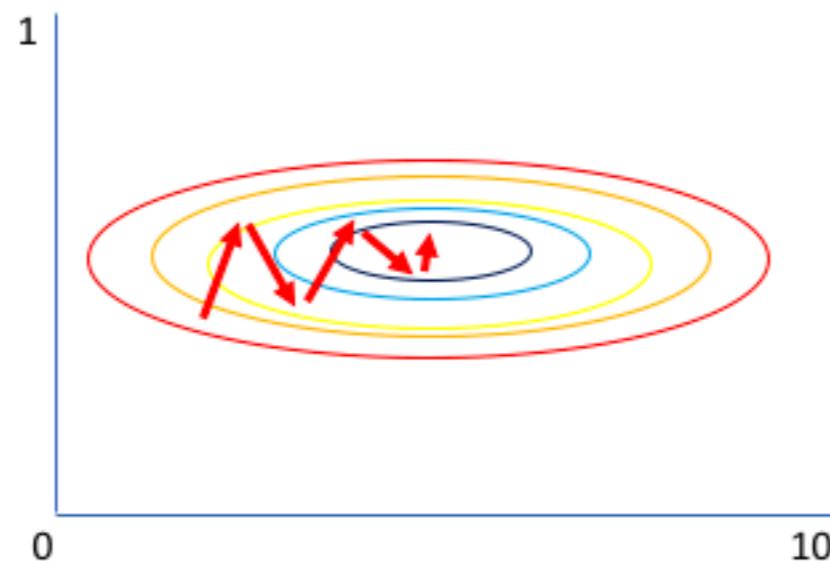
- Consider two images:
 - One of a person in very dark room
 - One of a person in bright sunlight
- We want to train a network on both types of images
- The difference between them is mainly their means (and std. deviations)



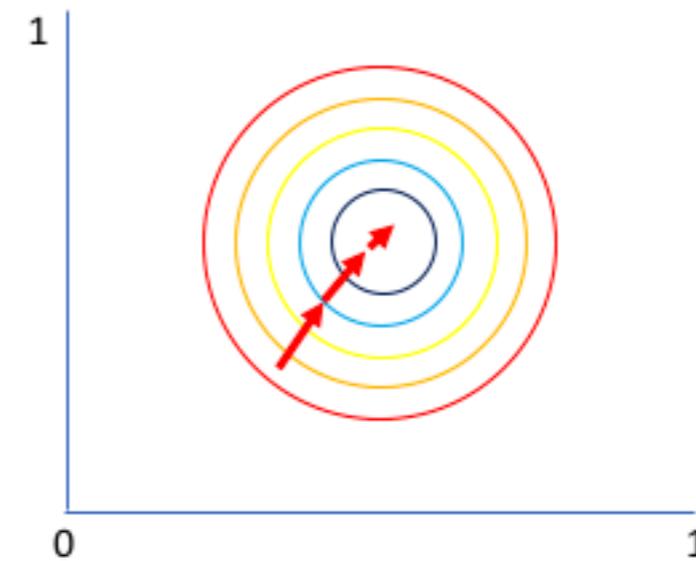
Input normalization



- We (almost) always normalize our input to some specified range ($[0,1]$ or $[-1,1]$) or to zero mean and std. dev. 1
- Intuitively this makes it easier for the network to learn since the output will also be in the same(ish) range



Gradient of larger parameter
dominates the update



Both parameters can be
updated in equal proportions



Input normalization



Input



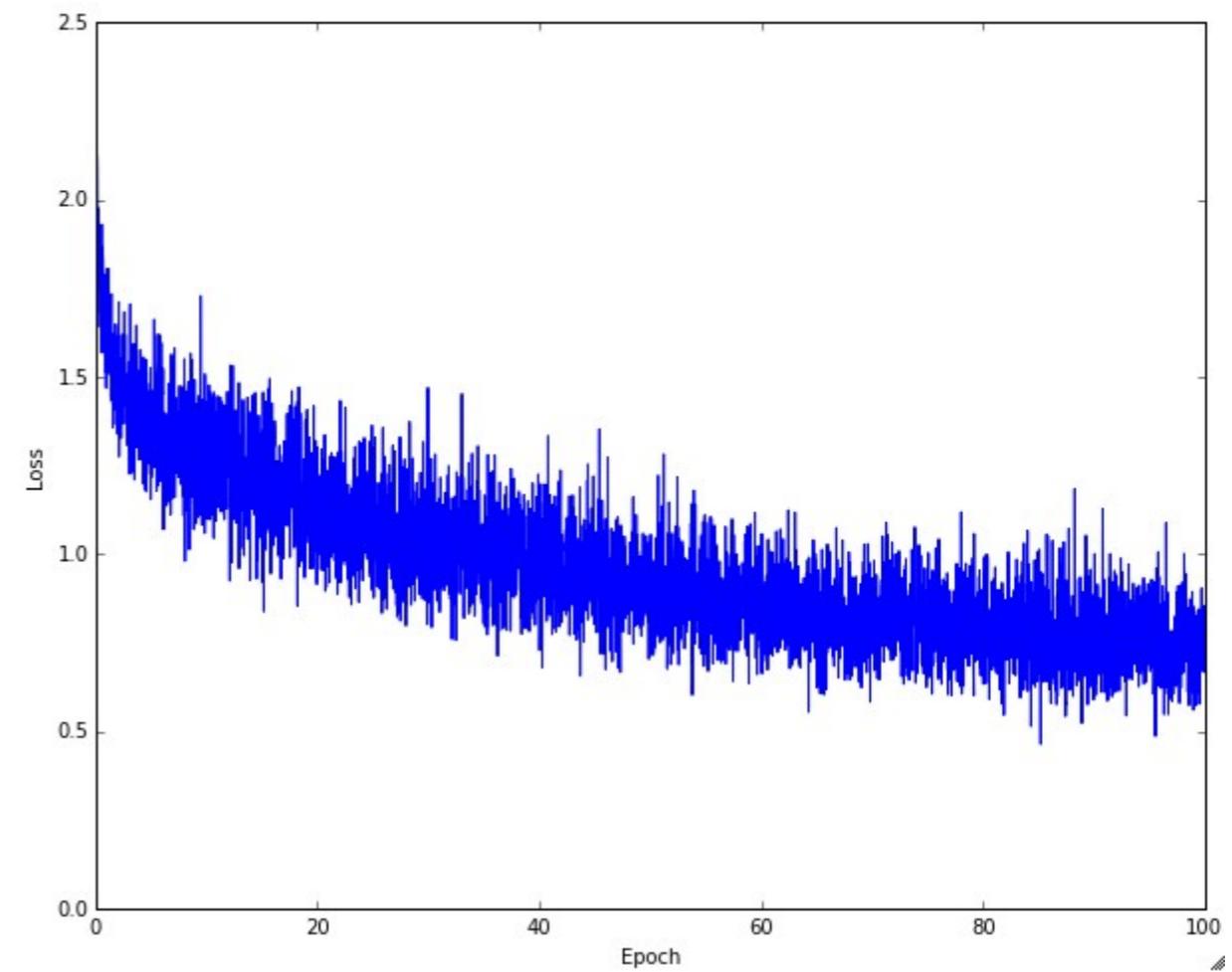
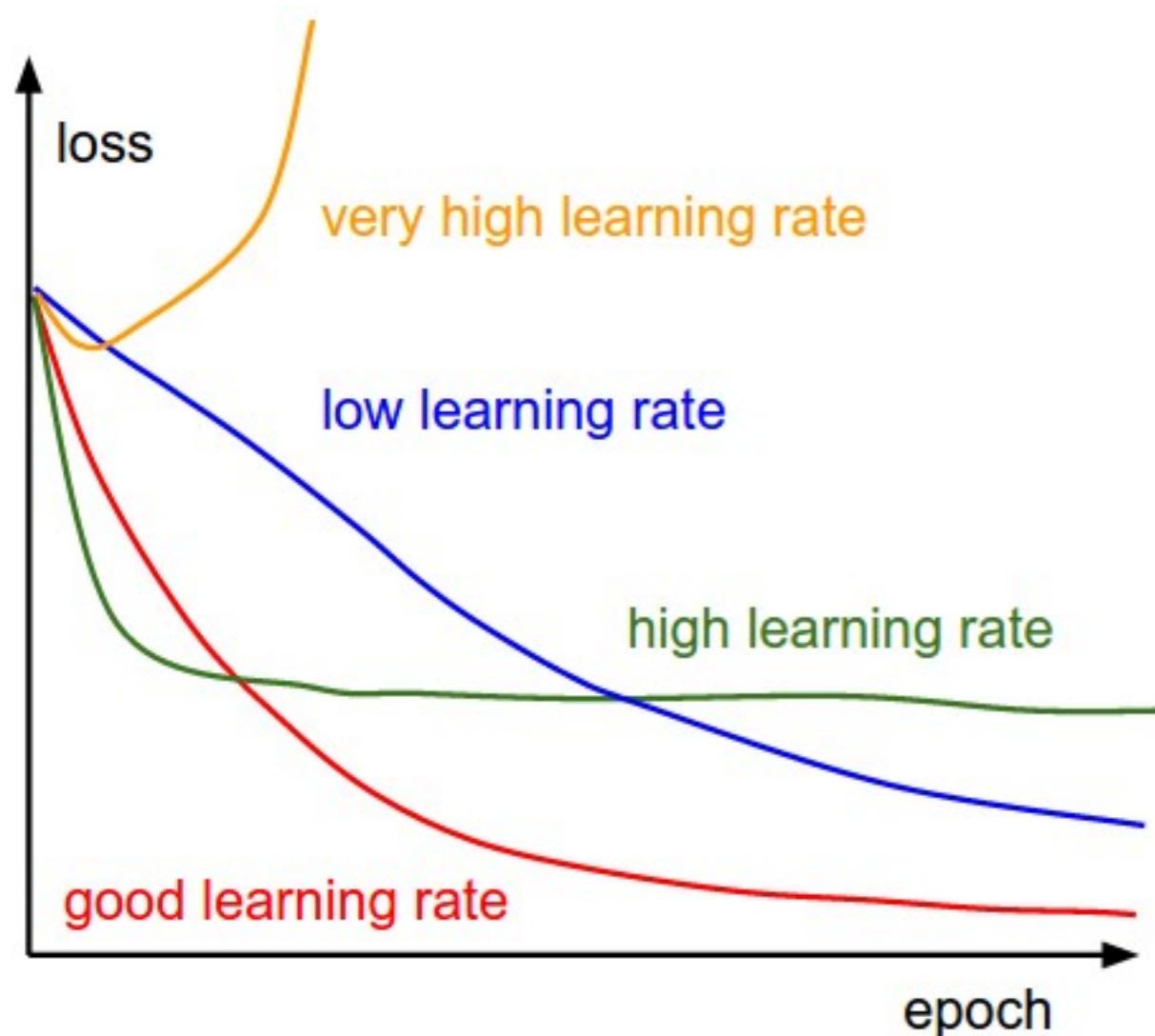
Output

Changes

Learning rate annealing

- Two common types of learning rate decay are
- Step decay
 - Reduce the learning rate by some factor $k < 1$ every few epochs. The factor and how many epochs is very dependent on the problem
- Exponential decay
 - Where a_0, k are hyperparameters and t is the time step (epoch number)
$$\alpha = \alpha_0 e^{-kt}$$

Learning rate annealing

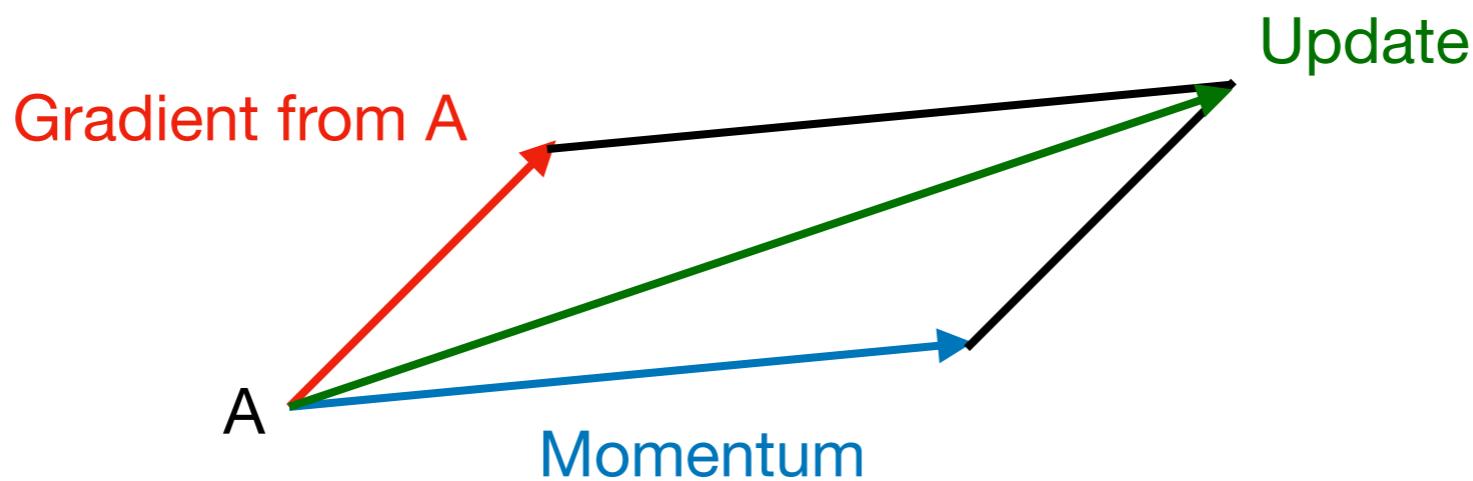


Using other optimisers

- Simply using SGD is sufficient for many tasks
- Sometimes, however, training takes a long time which may be reduced if we can take better steps adapting to the loss surface
- Many SGD alterations have been invented to do just this

SGD with momentum

- Consider a ball rolling down a hillside
 - Just because the ball meets a small obstacle does not mean that it completely changes direction
 - Instead it keeps some of the velocity in its current direction and add some velocity in another direction
 - We call this momentum
 - We can use this idea to take more ideal steps than SGD

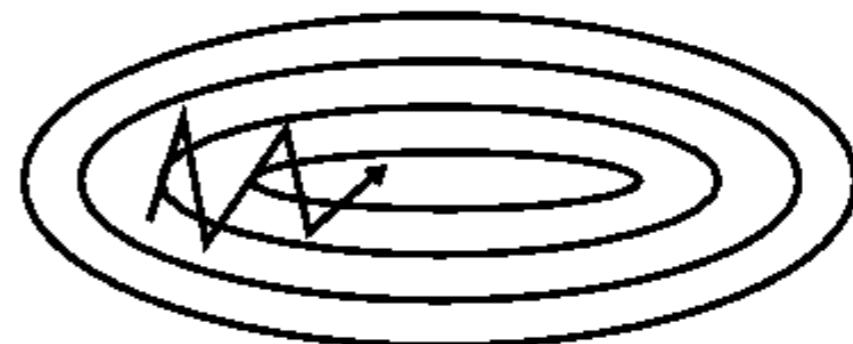


SGD with momentum

- If a surface curves more steeply in one dimension than in another it takes many steps for SGD to move to the minimum
- SGD would oscillate across the slopes with high curvature with only minor progress towards the minima



SGD



SGD with momentum

SGD with momentum

- Standard SGD

$$w_{j,k}^\ell \rightarrow \bar{w}_{j,k}^\ell = w_{j,k}^\ell - \alpha \frac{\partial L}{\partial w_{j,k}^\ell} \quad \mathbf{W} \rightarrow \bar{W} = \mathbf{W} - \alpha \nabla_W \mathcal{L}$$

- SGD + momentum

$$\mathbf{W} \rightarrow \bar{W} = \mathbf{W} - \mathbf{v}_t$$

$$\mathbf{v}_t = \gamma \mathbf{v}_{t-1} + \alpha \nabla_W \mathcal{L}$$

- Typically $\gamma = 0.9$

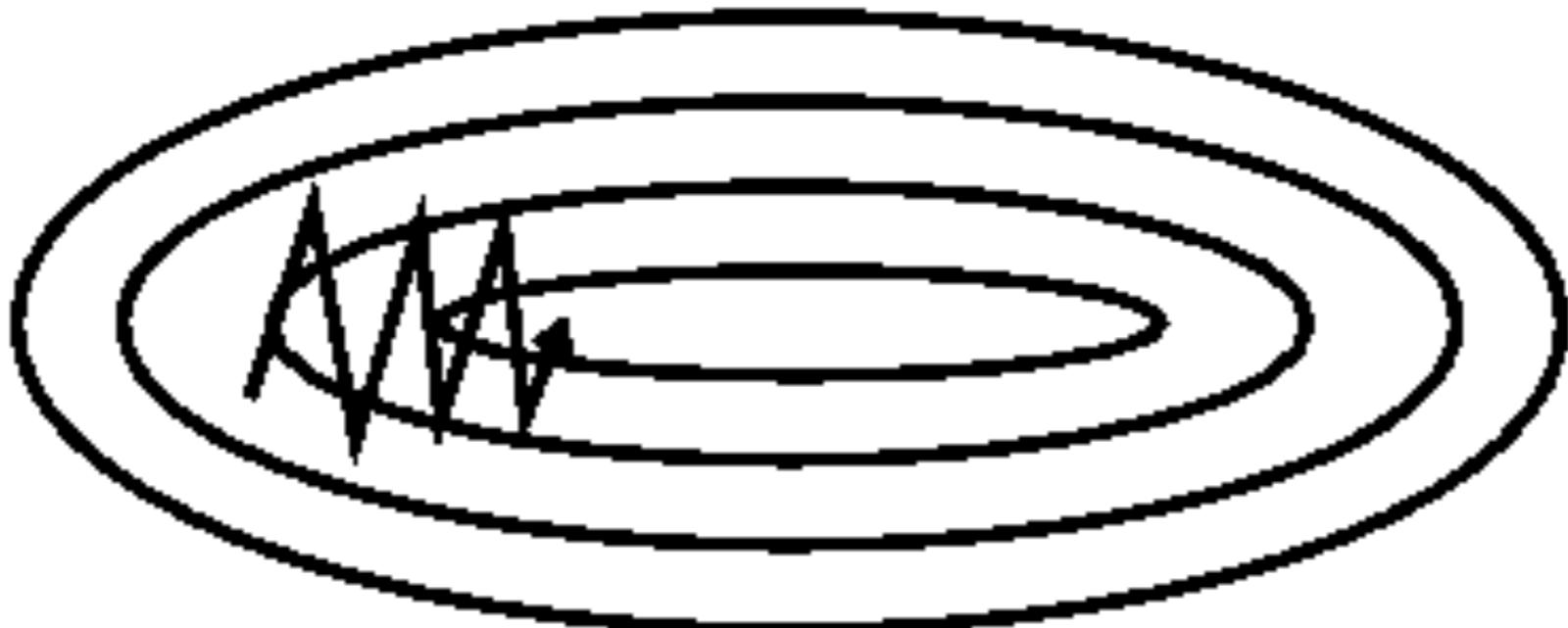
RMSProp

- Root Mean Square Propagation (RMSProp)
- The idea is to adapt the learning rate α for each individual parameter such that we take large steps in “good” dimensions and smaller steps in “bad” dimensions

$$v_t = \gamma v_{t-1} + (1 - \gamma) \left(\frac{\partial \mathcal{L}}{\partial w_{j,k}^\ell} \right)^2$$

$$w_{j,k}^\ell \rightarrow \bar{w}_{j,k}^\ell = w_{j,k}^\ell - \frac{\alpha}{\sqrt{v_t}} \frac{\partial \mathcal{L}}{\partial w_{j,k}^\ell}$$

RMSProp



$$w_{j,k}^\ell \rightarrow \bar{w}_{j,k}^\ell = w_{j,k}^\ell - \frac{\alpha}{\sqrt{v_t}} \frac{\partial \mathcal{L}}{\partial w_{j,k}^\ell}$$

Adam

- Adaptive Moment Estimation (Adam) is a combination of momentum and RMSProp

$$m_w^t = \gamma_1 m_w^{t-1} + (1 - \gamma_1) \frac{\partial \mathcal{L}}{\partial w_{j,k}^\ell}$$
$$v_w^t = \gamma_2 v_w^{t-1} + (1 - \gamma_2) \left(\frac{\partial \mathcal{L}}{\partial w_{j,k}^\ell} \right)^2$$

- Bias corrections

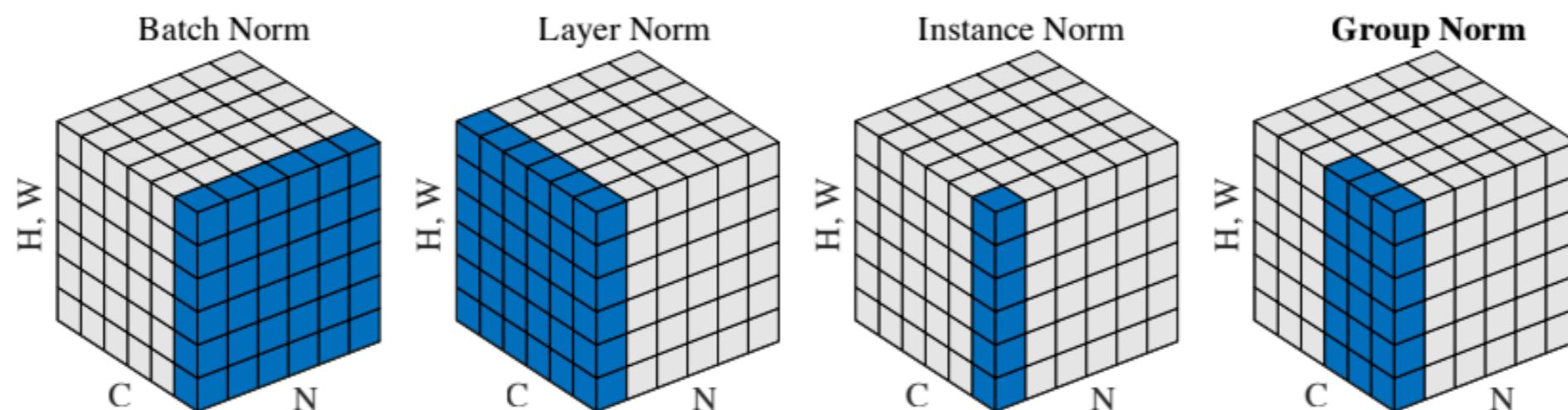
$$\hat{m}_w^t = \frac{m_w^t}{1 - \gamma_1^{t+1}}, \quad \hat{v}_w^t = \frac{v_w^t}{1 - \gamma_2^{t+1}}$$

- Update

$$w_{j,k}^\ell \rightarrow \bar{w}_{j,k}^\ell = w_{j,k}^\ell - \alpha \frac{m_w^t}{\sqrt{v_w^t}}$$

Batch normalisation

- We normalize the input to the network for faster convergence and more stability
 - Why not do the same for the outputs from the hidden layers?
- This is called batch normalisation
- Batch normalisation normalises the output of each layer by subtracting the mean and dividing by the standard deviation of the of mini-batch outputs



Batch normalisation

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;

Parameters to be learned: γ, β

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

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

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

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

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

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

Batch normalisation

- Batch normalisation approximately makes that the output distribution from each layer is always the same
- Batch normalisation adds some noise to the output of each neuron (similar to dropout)
 - This functions as a regulariser



Batch normalisation at test time



- At test time we might only process images one at a time
 - We cannot do mini-batch statistics
- We keep a running average of the minibatch statistics across the minibatches to use at test time

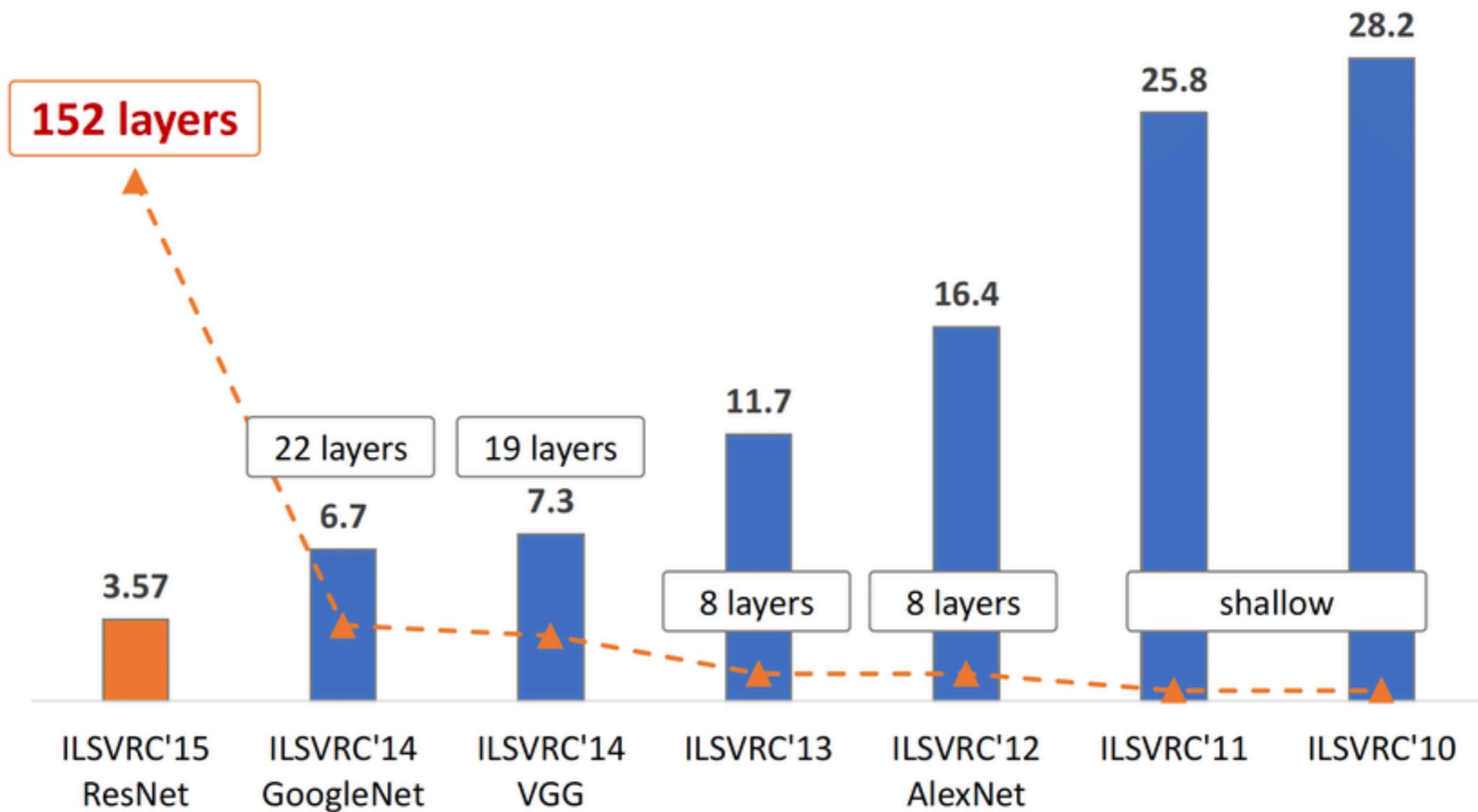
$$\mu_t = \gamma\mu_{t-1} + (1 - \gamma)\mu_B$$

Common architectures

ImageNet Large Scale Visual Recognition Challenge

- The ILSVRC is a large scale challenge for object detection and image classification
- ImageNet has more than 14 million images in more than 20.000 categories
- The yearly challenge has spurred some of the most widely used architectures for visual recognition today

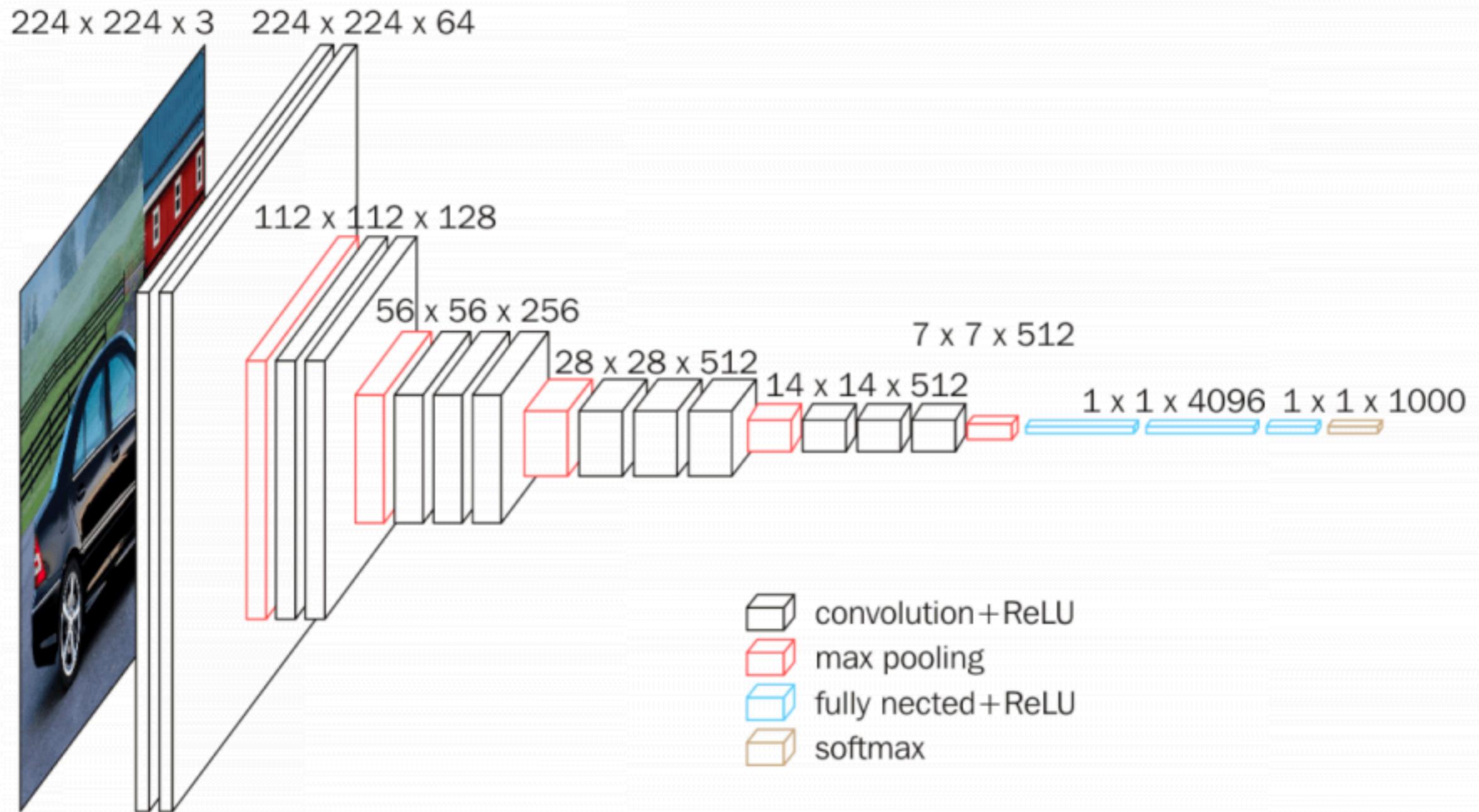
ILSVRC



VGGNet

- Simonyan & Zisserman pushed the standard convolution neural network with 3x3 convolutions to 19 layers in 2014
- Runner-up of ILSVRC 2014 (Their best model had 16 layers)
- VGG16 has around 138 million parameters
 - Using 4 GPUs it took 2-3 weeks to train

VGNet

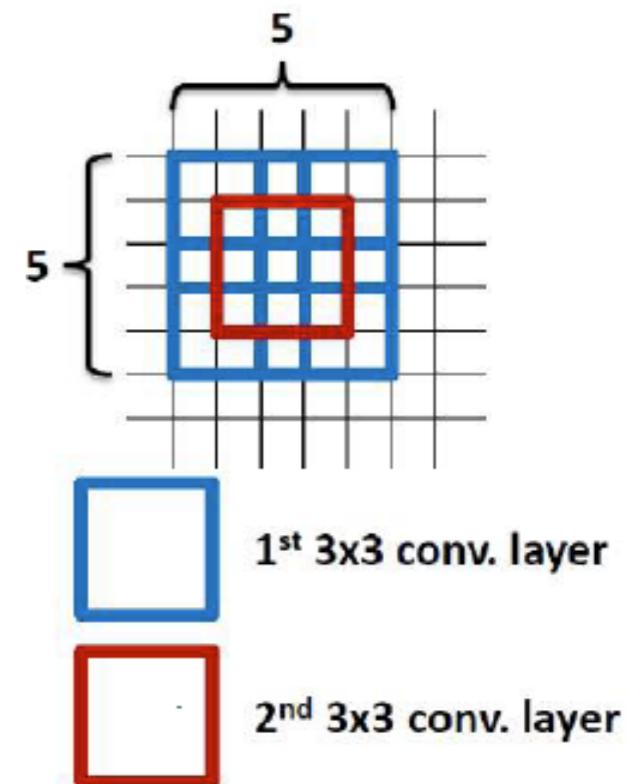


VGGNet

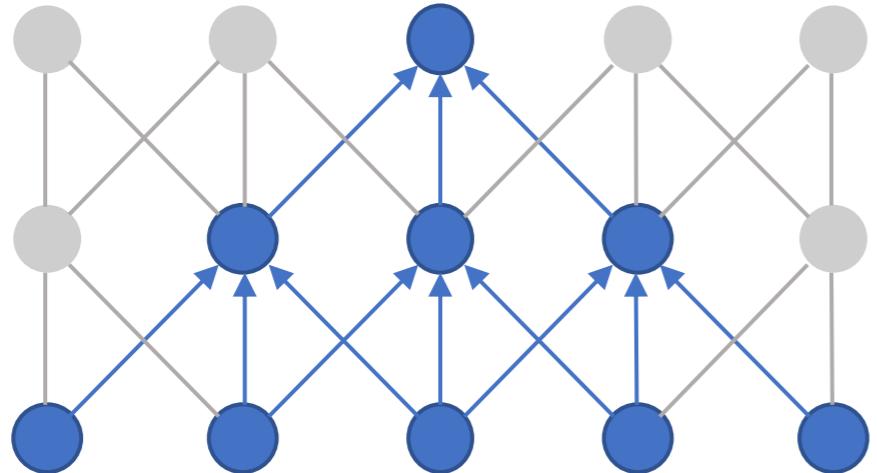
- Why only use 3x3 layers?

VGGNet

- Why only use 3x3 layers?
 - Stacking layers have a large receptive field
 - Two 3x3 layers - 5x5 receptive field
 - Three 3x3 layers 7x7 receptive field
 - ...
 - Less parameters to learn



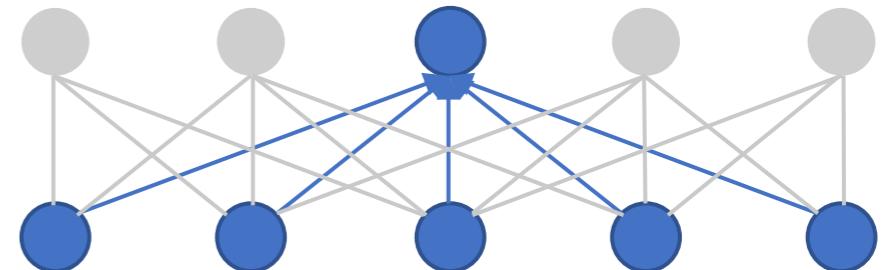
VGGNet



Two 3x3 convolutions

$$\text{Parameters: } 2 \cdot (3^2 + 1) = 20$$

$$\text{Operations: } 2 \cdot (2 \cdot 3^2) = 36$$



One 5x5 convolution

$$\text{Parameters: } 5^2 + 1 = 26$$

$$\text{Operations: } 2 \cdot 5^2 = 50$$

Three 3x3 convolutions

$$\text{Parameters: } 3 \cdot (3^2 + 1) = 30$$

$$\text{Operations: } 3 \cdot (2 \cdot 3^2) = 54$$

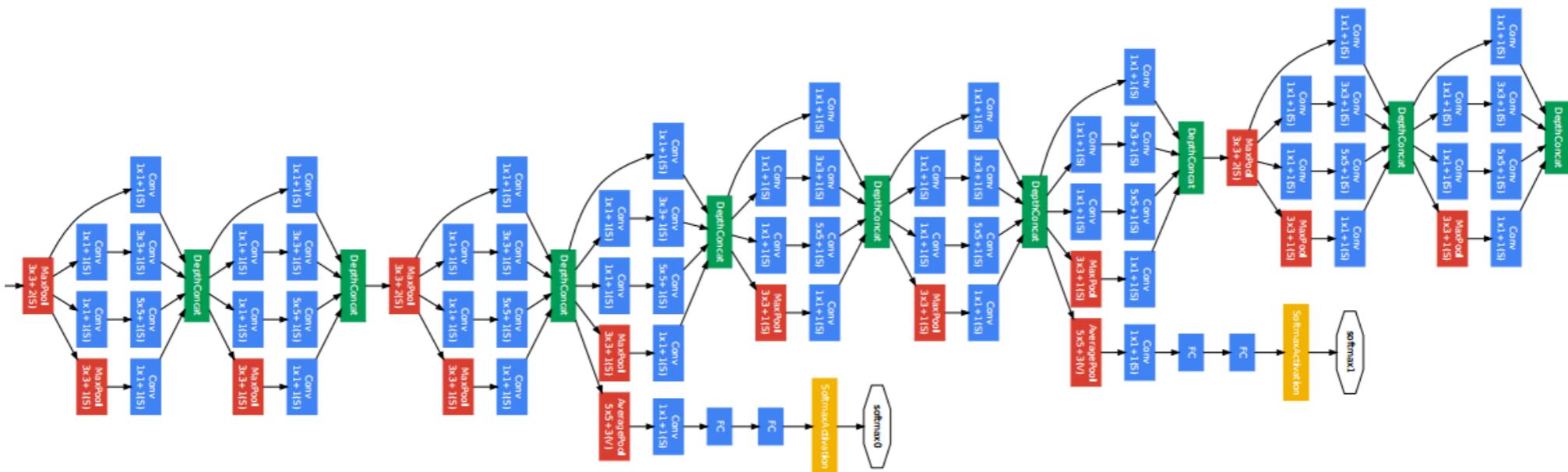
One 7x7 convolution

$$\text{Parameters: } 7^2 + 1 = 50$$

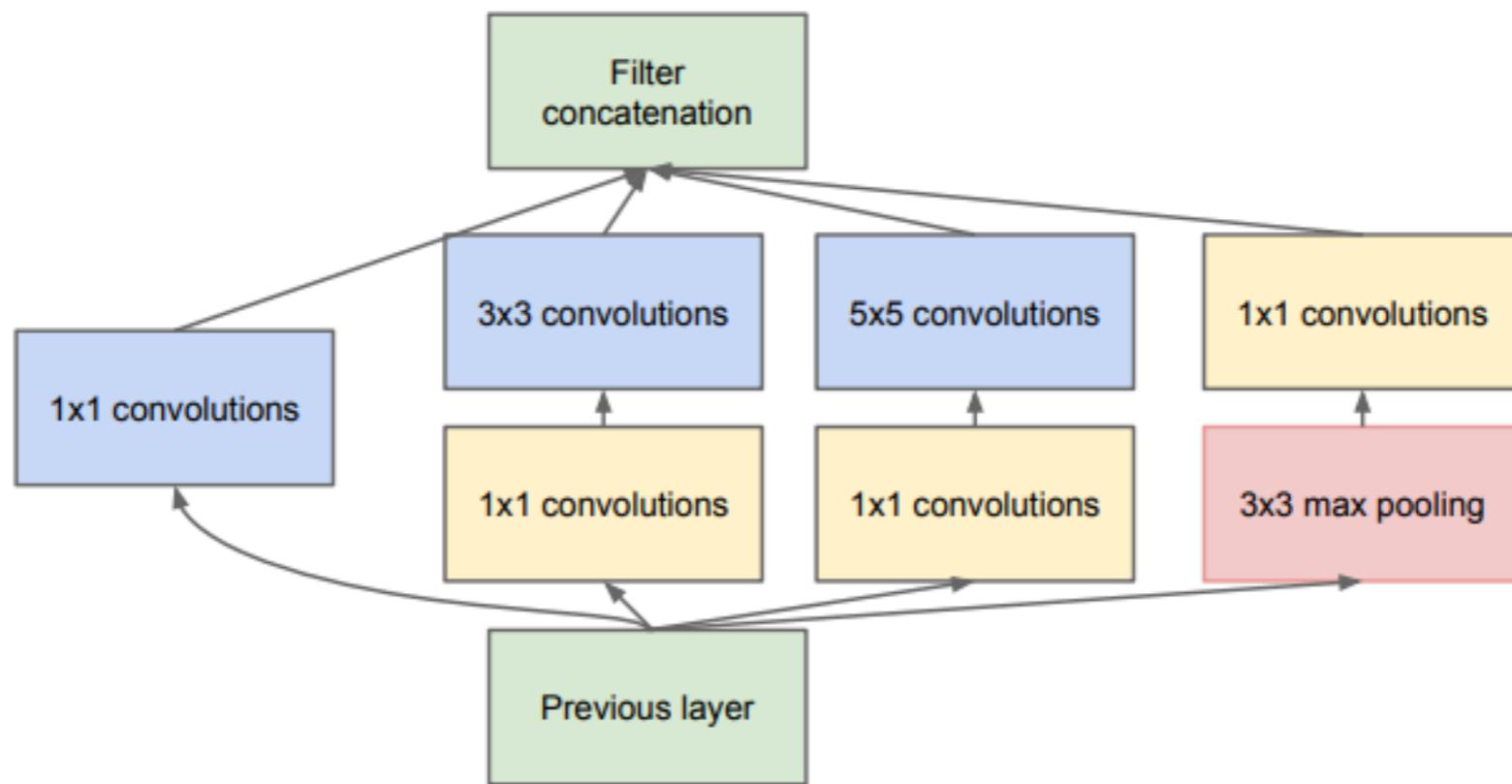
$$\text{Operations: } 2 \cdot 5^2 = 98$$

GoogLeNet/Inception

- Winner of the ILSVRC 2014 Classification challenge
 - Able to train network with 23 layers



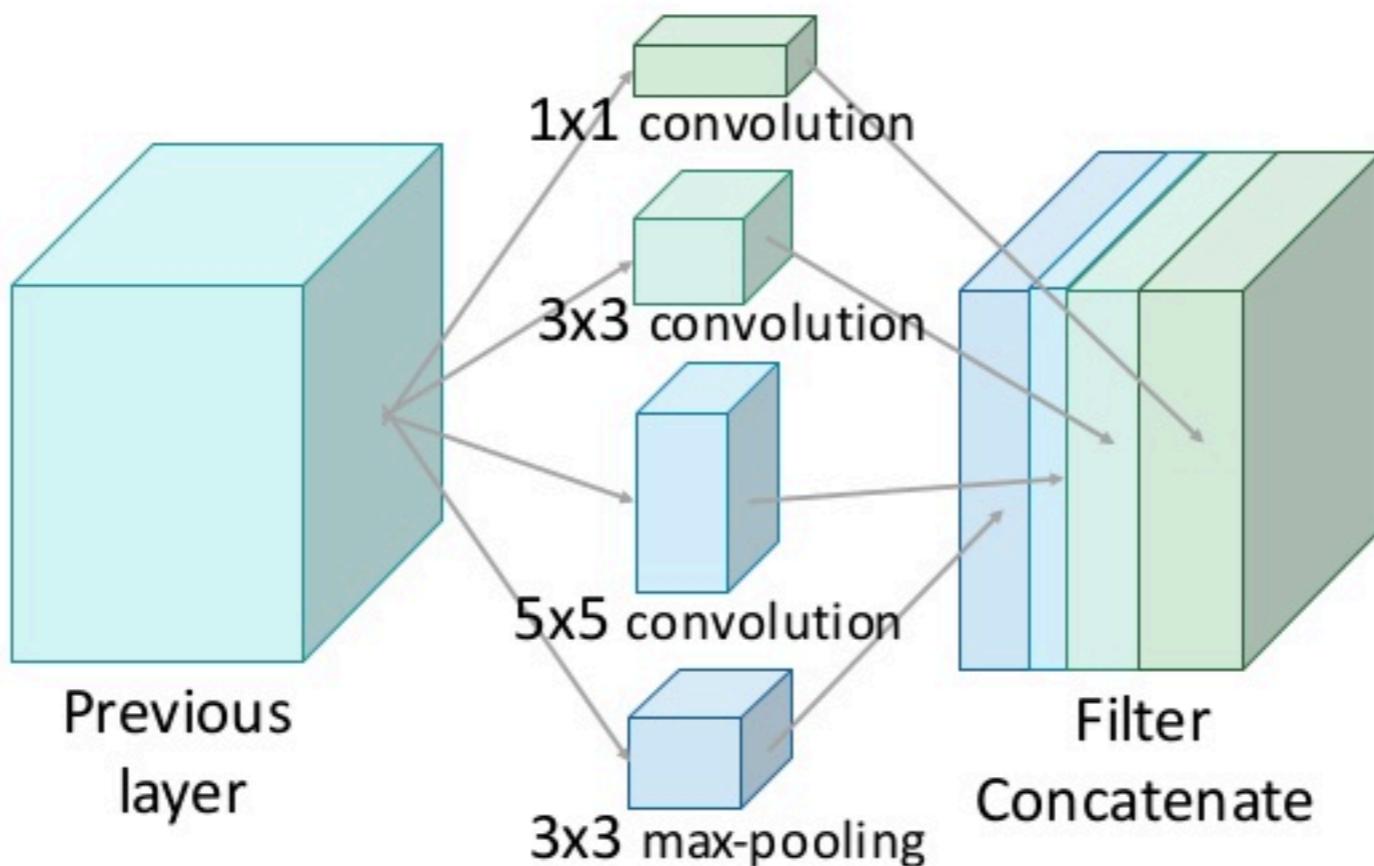
InceptionNet



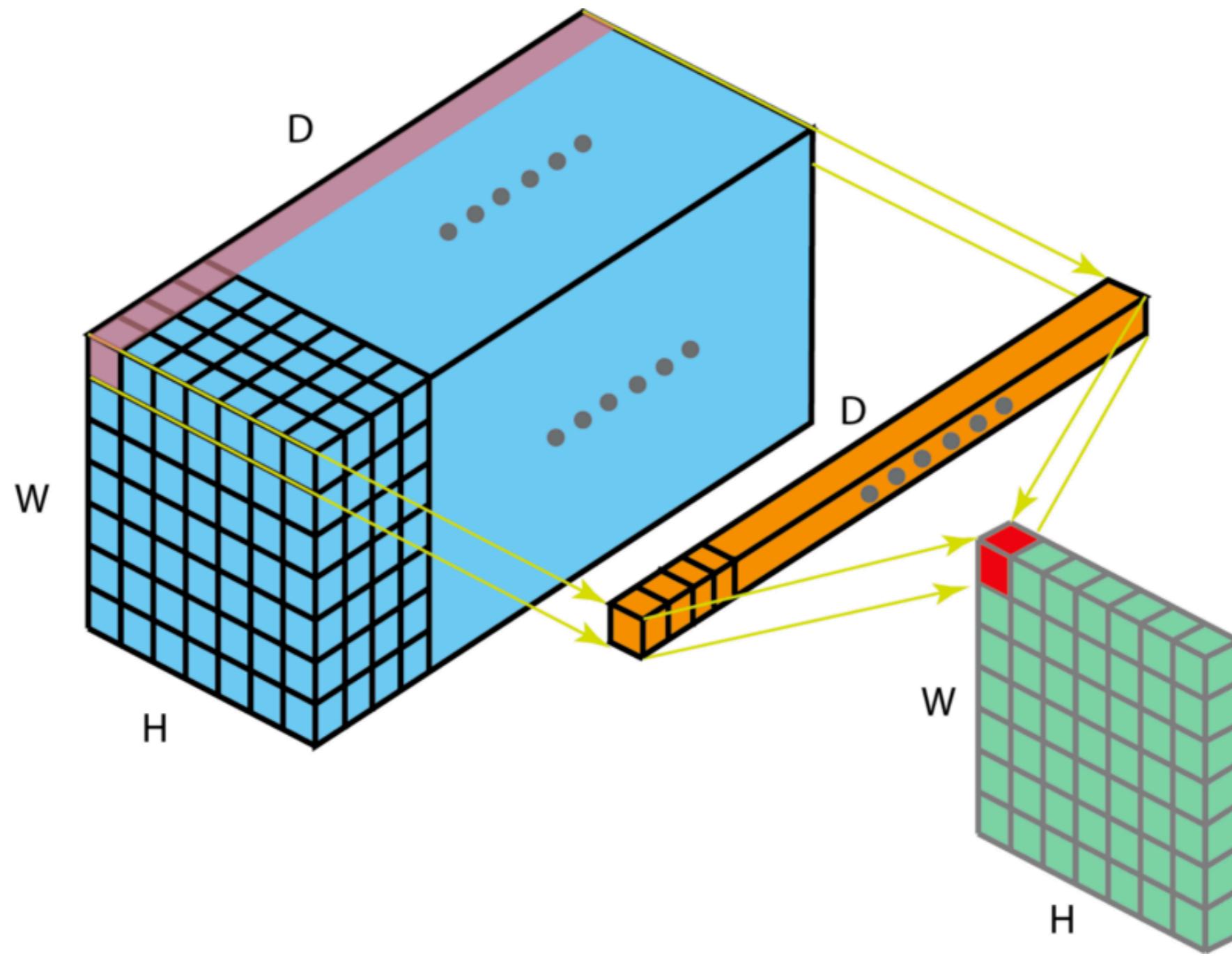
(b) Inception module with dimension reductions

InceptionNet

Inception Module

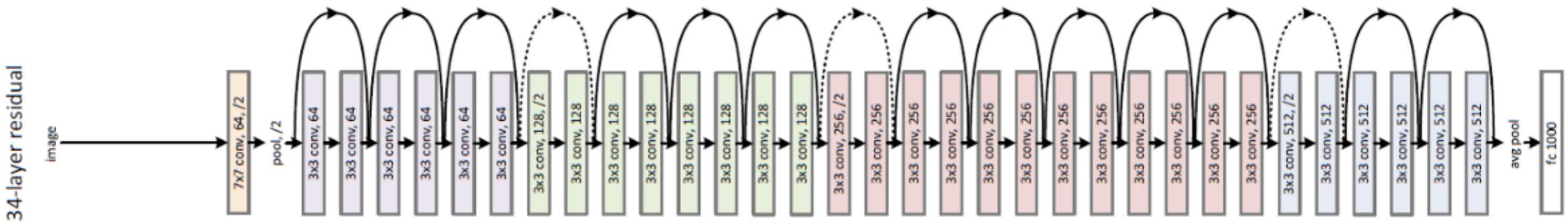


1×1 convolution

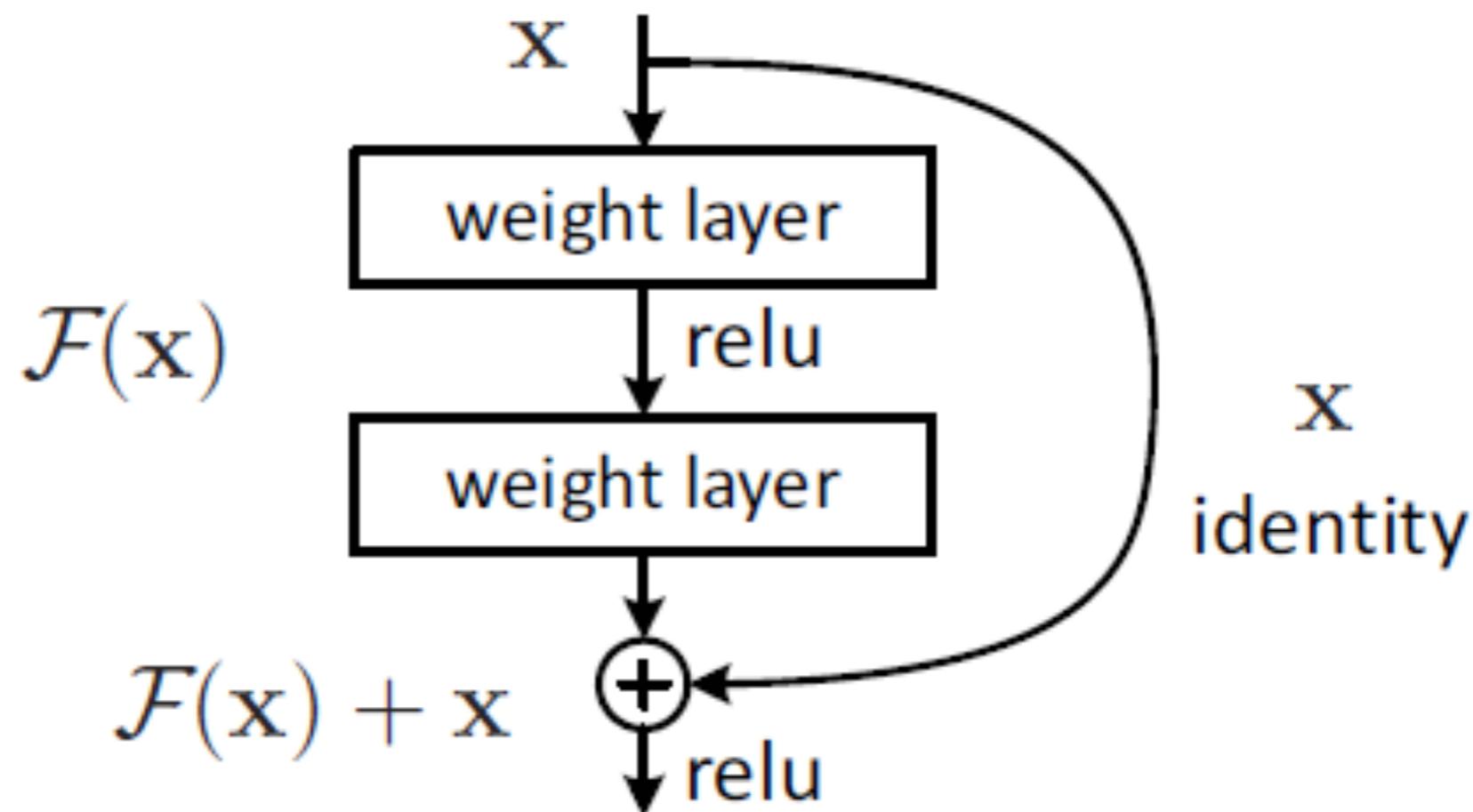


ResNet

- Won the ILSVRC 2015 Challenge
 - Could train 152 layers on ImageNet



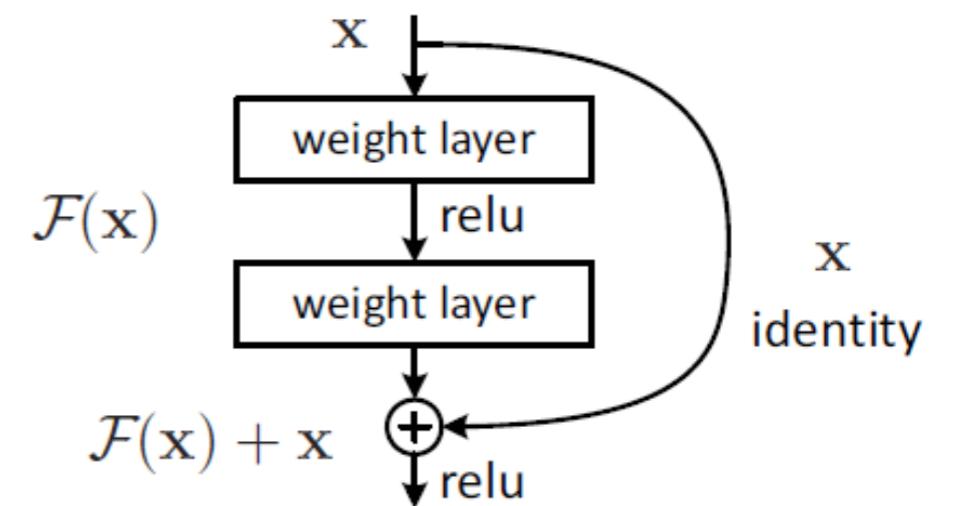
The building block of ResNets



ResNets



- Why can we train very large networks using residual blocks?
- If our weights are zero in a residual block we just feed the values through - so it can function as a more shallow network
- We can now train ResNets with more than 1000 layers
- In back propagation the gradients can use the identity mappings to move faster back to the early layers



conda install matplotlib

**Remember to use
gpustat**