

Session 3 & Assignment

✓ Published

 [Edit](#)

⋮

Advanced Convolutions, Batch Normalization and Receptive Field

Session 3 Video

Session 3 - EIP4

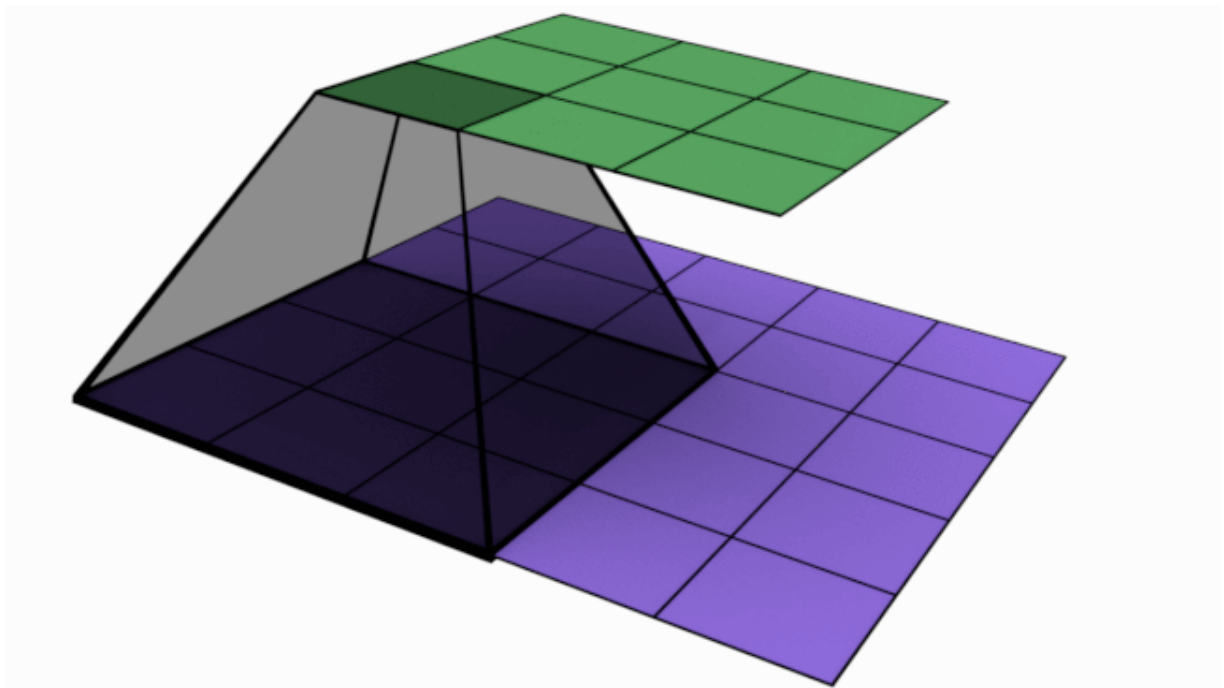


Advanced Convolutions

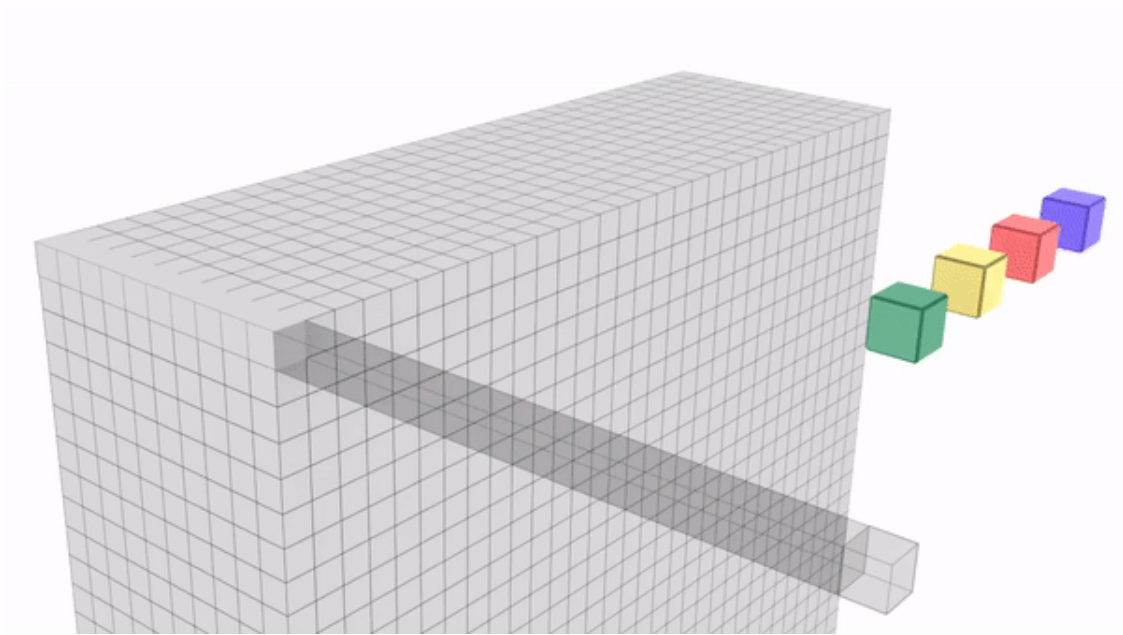
What is Convolution?

The process of extracting features from input data using kernels/filters. Filter moves discretely on top of data-plane/channel, scales the input data equal to the size of it's receptive field, sums this results and creates a new feature map.

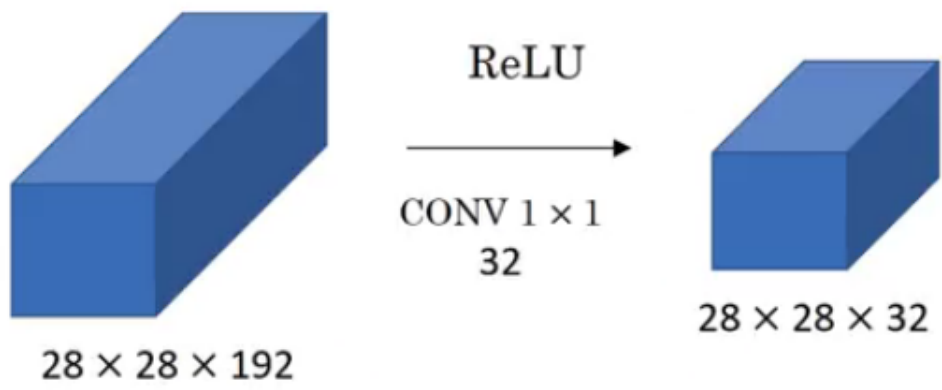
Normal Convolutions



Pointwise Convolution



Basically

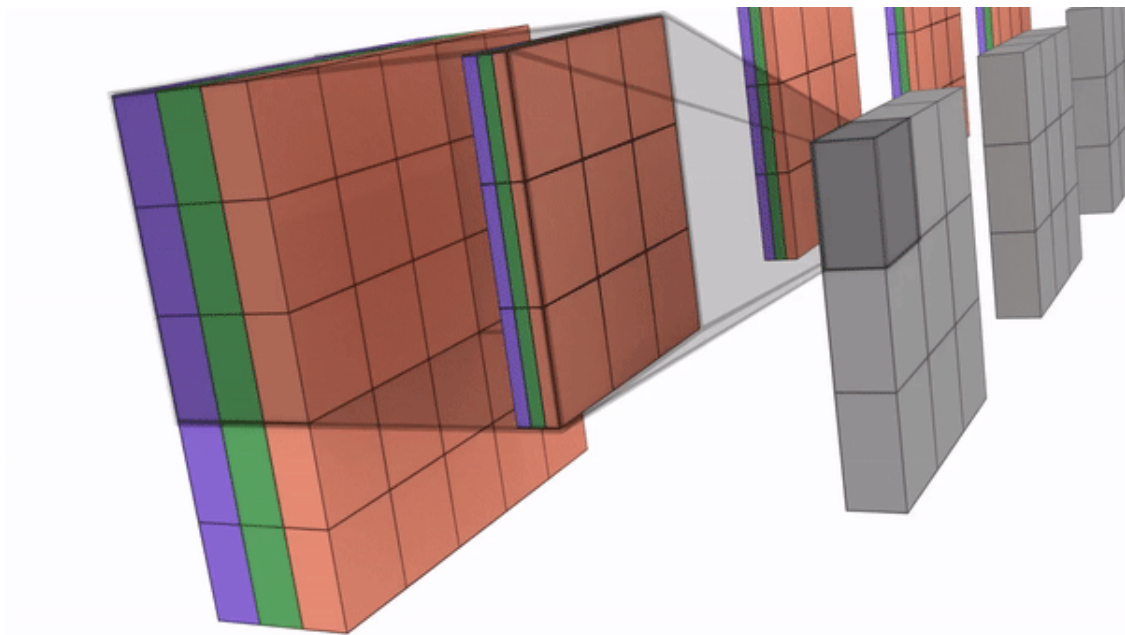


Behind the scenes:

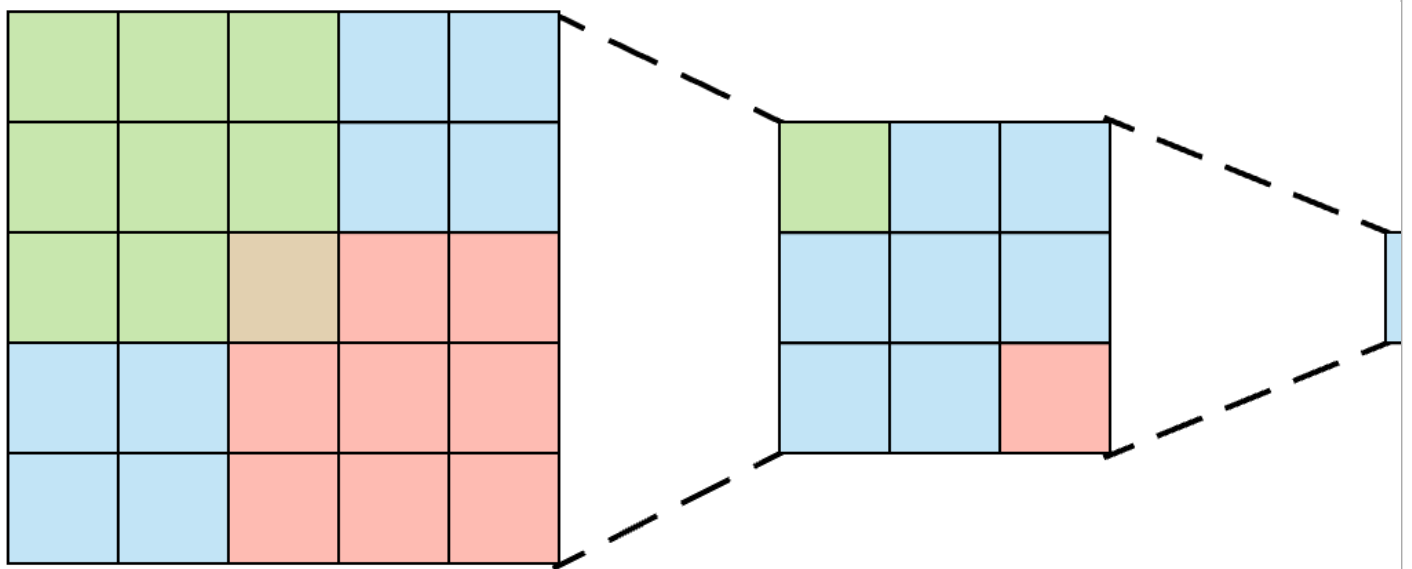
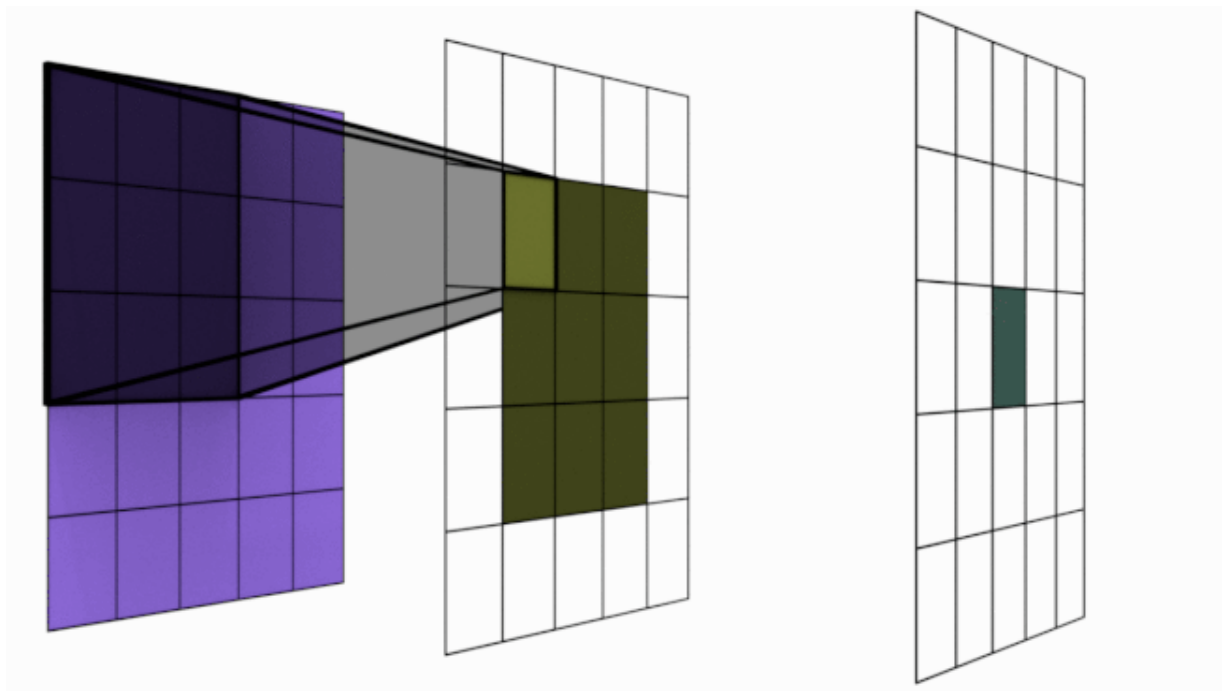
0 ₂	0 ₀	0 ₁	0	0	0	0
0 ₁	2 ₀	2 ₀	3	3	3	0
0 ₀	0 ₁	1 ₁	3	0	3	0
0	2	3	0	1	3	0
0	3	3	2	1	2	0
0	3	3	0	2	3	0
0	0	0	0	0	0	0

1	6	5
7	10	9
7	10	8

Concept of Channels



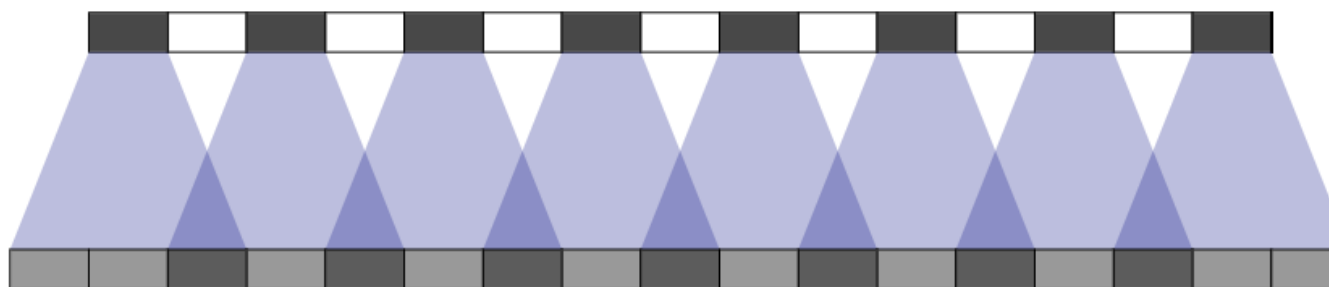
Receptive Field



5x5 receptive field

3x3 receptive field

Checkboard Issue



A Note on The Strides



- why should we not use strides?
- why should we use strides or when do we use them?

We get a receptive field of 3×3 when we convolve by 3×3 . What should we do to get a receptive field of 5×5 when we convolve by 3×3 ?

Atrous Convolutions

Or

Dilated Convolution

Dilated convolution is a way of increasing receptive view (global view) of the network exponentially and linear parameter accretion. With this purpose, it finds usage in applications cares more about integrating the knowledge of the wider context with less cost.

The key application the dilated convolution authors have in mind is dense prediction: vision applications where the predicted object that has similar size and structure to the input image. For example, semantic segmentation with one label per pixel; image super-resolution, denoising, demosaicing, bottom-up saliency, keypoint detection, etc.

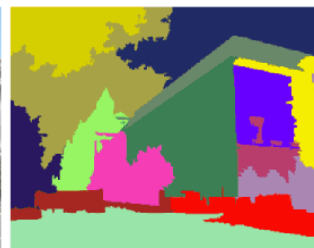


Image
Segmentation



Image
Superresolution

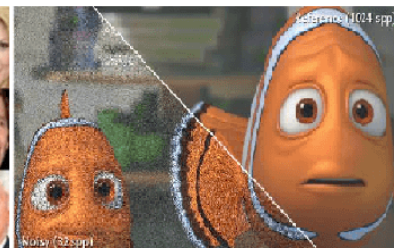


Image
Denoising

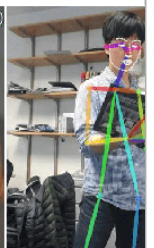
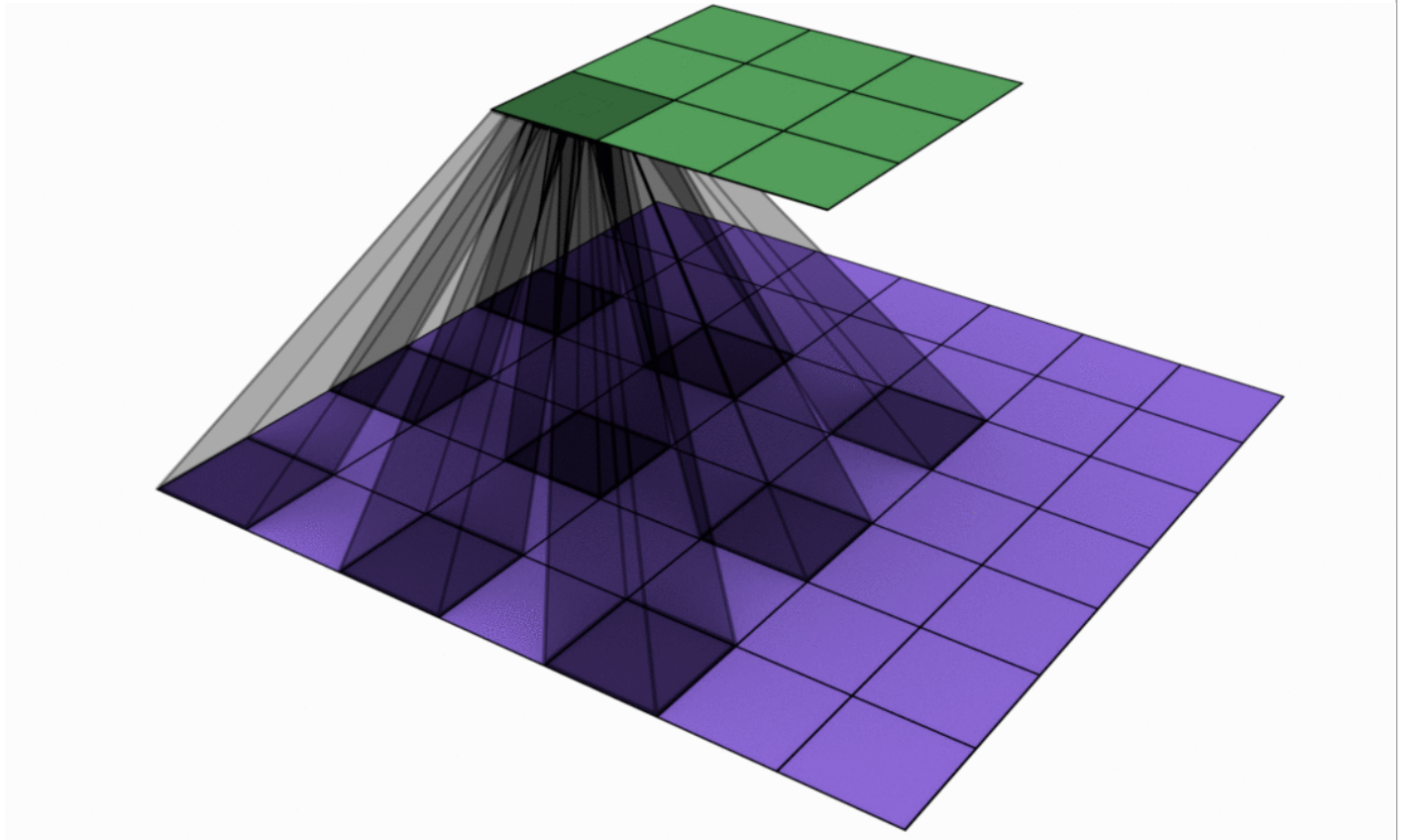


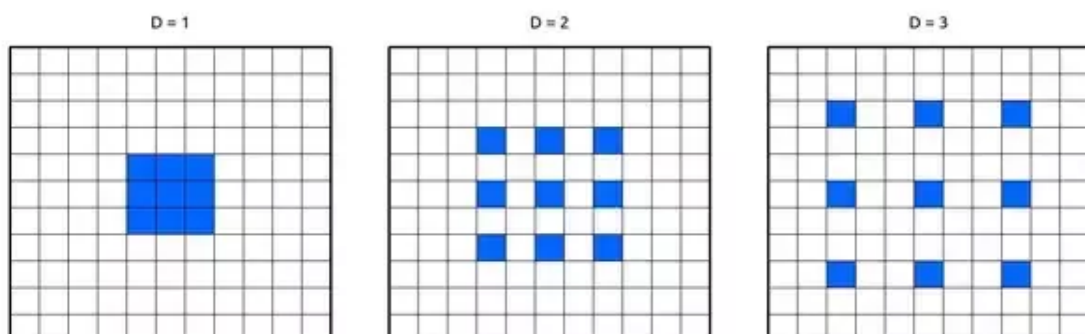
Image Keyp
Detection

In many such applications one wants to integrate information from different spatial scales and balance two properties:

1. local, pixel-level accuracy, such as precise detection of edges, and
2. integrating the knowledge of the wider, global context



Dilation:



The Problem: "convolution is a rather local operation"

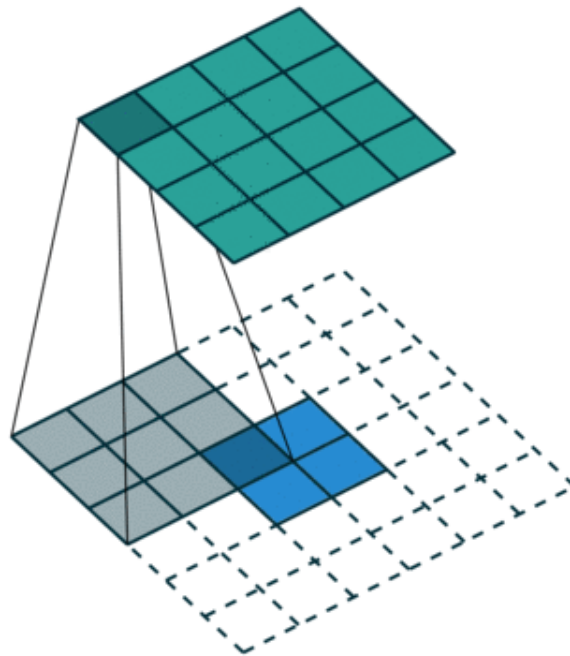
Improving the resolution of segmentation results

Dilated convolutions or atrous convolutions present a potential solution to this problem: they are capable of capturing global context without reducing the resolution of the segmentation map. Normal convolutional neural networks (CNNs) are based on the fact that each convolutional layer collects information from a larger neighborhood around each pixel/voxel. This means that in the upper layers of the network, each pixel of a feature map could potentially hold information about a large region of the image. However, experiments show that during learning this is usually not the case. Rather, the **information in each pixel stays localized** [[source](https://arxiv.org/abs/1412.6856) [\(https://arxiv.org/abs/1412.6856\)](https://arxiv.org/abs/1412.6856)] and the network “does not live up to its full potential”. Dilated convolutions change the rules of the game by using kernels of the same size as normal convolutional layers but spread out over a larger area on the image (see animation).

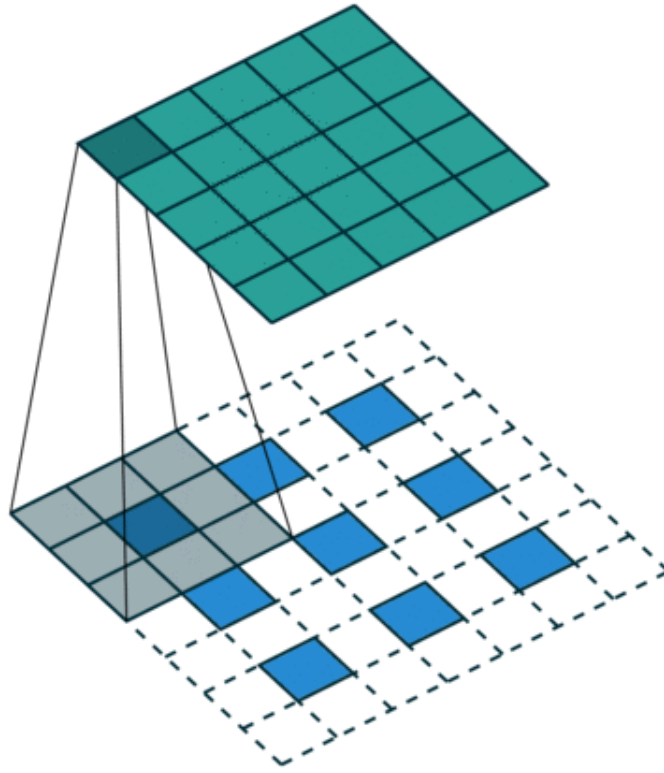
DECONVOLUTION or Fractionally Strided OR Transpose Convolution

We have a kernel of size 3×3 . If we convolve on 5×5 we get 3×3 . What do we do to get 7×7 ? Or get a larger channel size than we started with?

A simple approach, but inefficient.



Deconvolution/Transpose Convolution



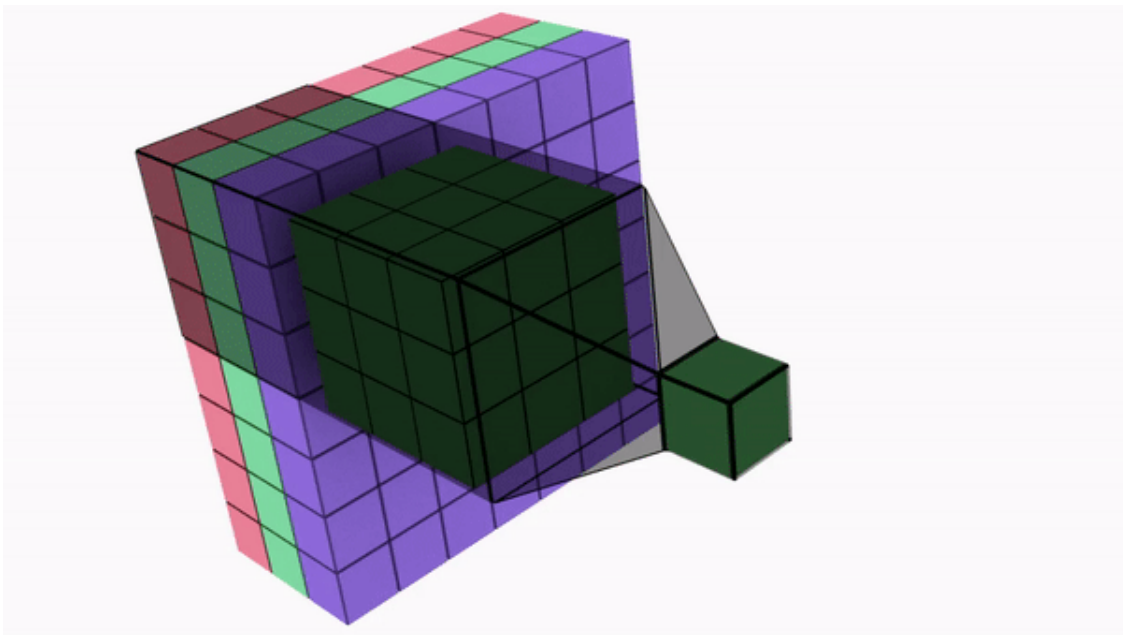
it introduces the checker board issue!

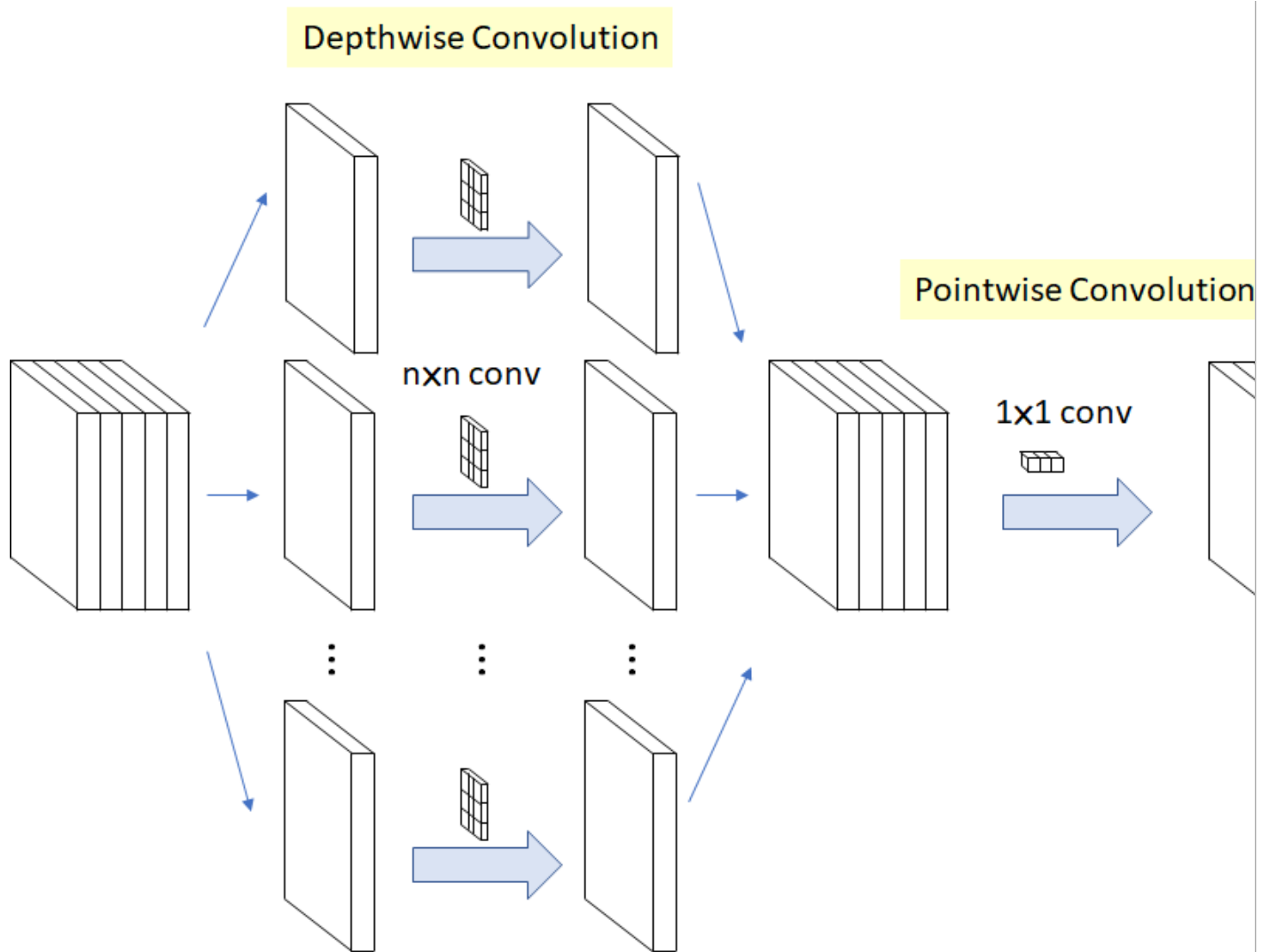
Perceptual Losses for Real-Time Style Transfer and Super-Resolution



Depthwise Separable Convolution

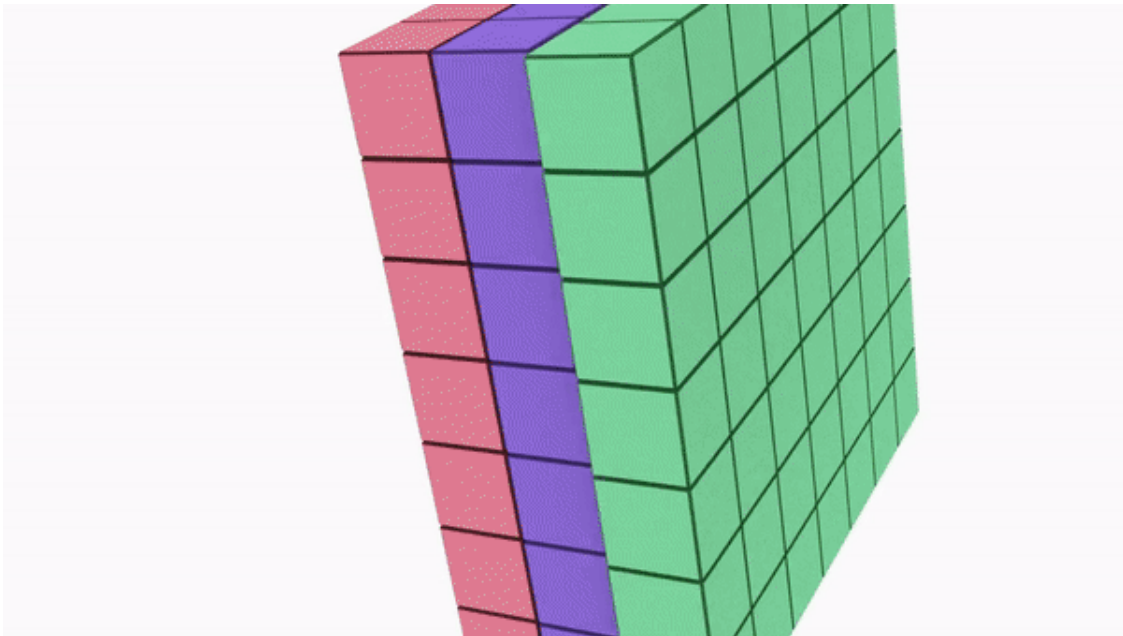
What is wrong in this convolution?





Used by XCEPTION and MOBILENET a lot

Depthwise Convolution



Reference [_\(https://eli.thegreenplace.net/2018/depthwise-separable-convolutions-for-machine-learning/\)_](https://eli.thegreenplace.net/2018/depthwise-separable-convolutions-for-machine-learning/)

For a depthwise separable convolution on the same example, we traverse the 16 channels with 1 3×3 kernel each, giving us 16 feature maps. Now, before merging anything, we traverse these 16 feature maps with 32 1×1 convolutions each and only then start to them add together. This results in 656 $(16 \times 3 \times 3 + 16 \times 32 \times 1 \times 1)$ parameters opposed to the 4608 $(16 \times 32 \times 3 \times 3)$ parameters from above.

Other convolutions you should learn about:

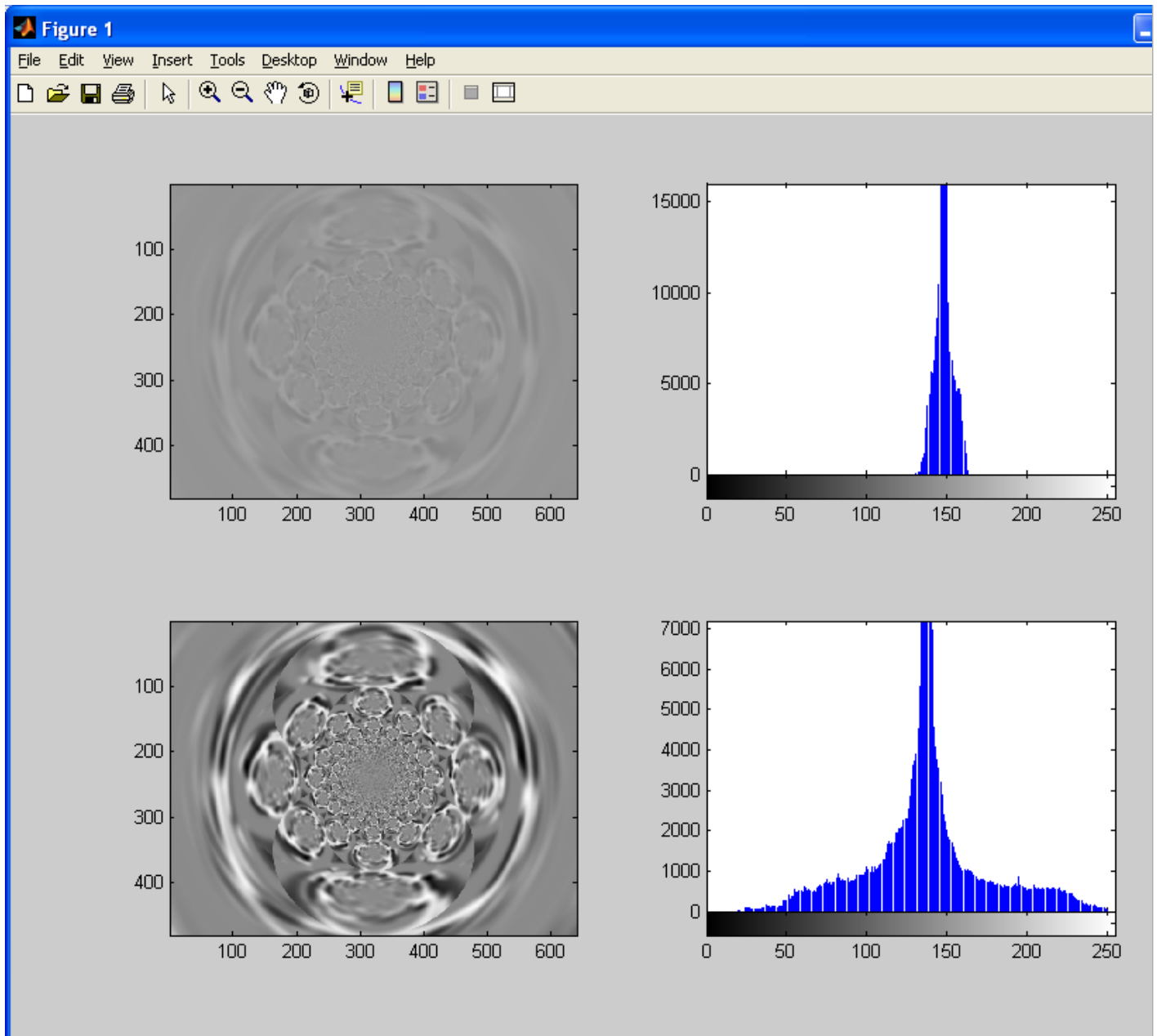
- spacially separable convolution
- grouped convolutions
- roots concept in convolutions

Batch Normalization & Regularization

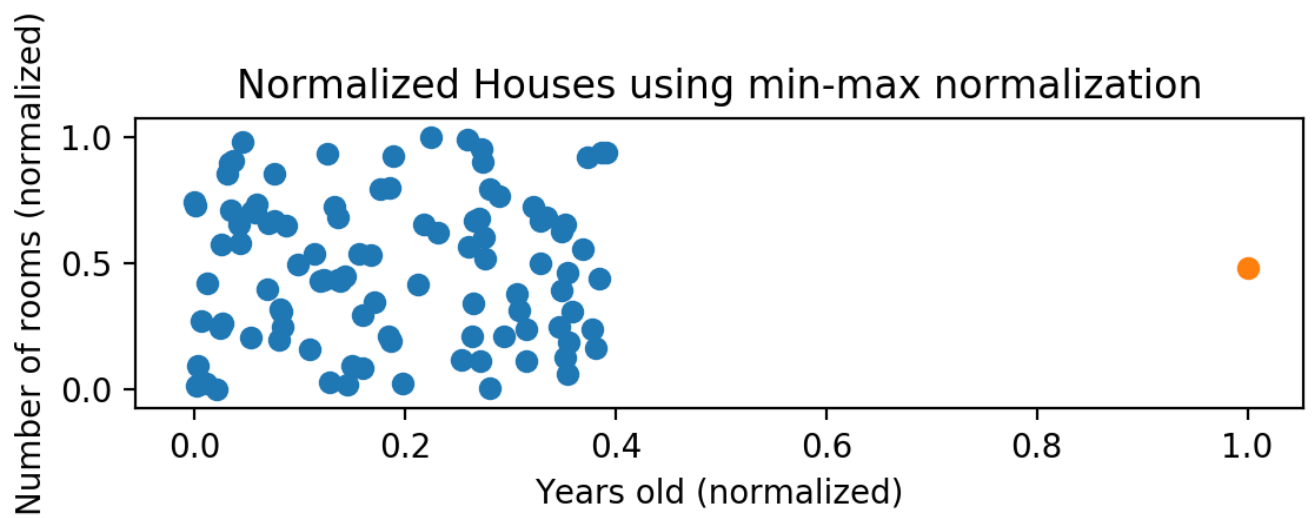
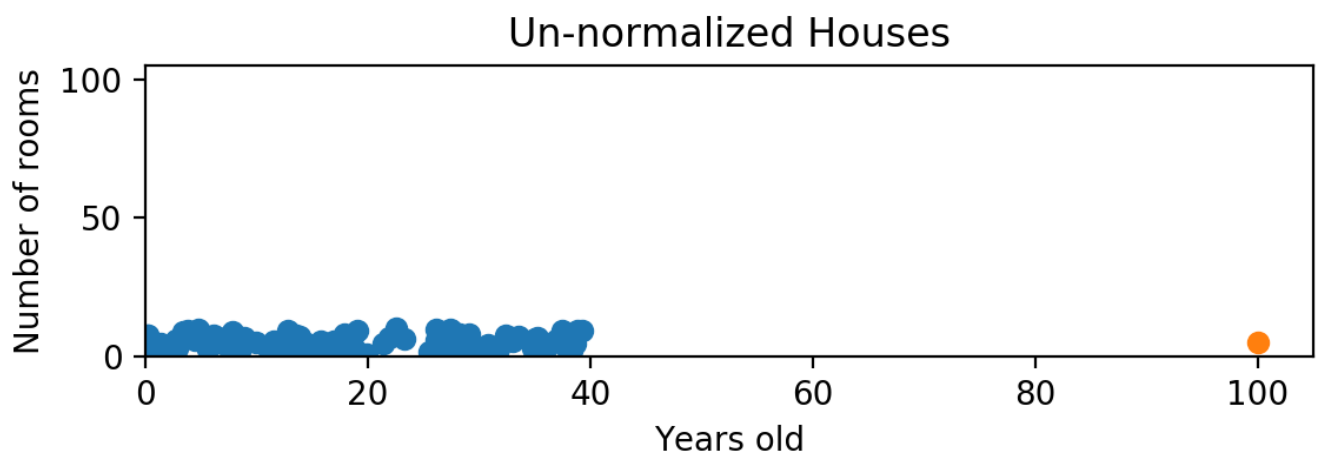
Image Normalization

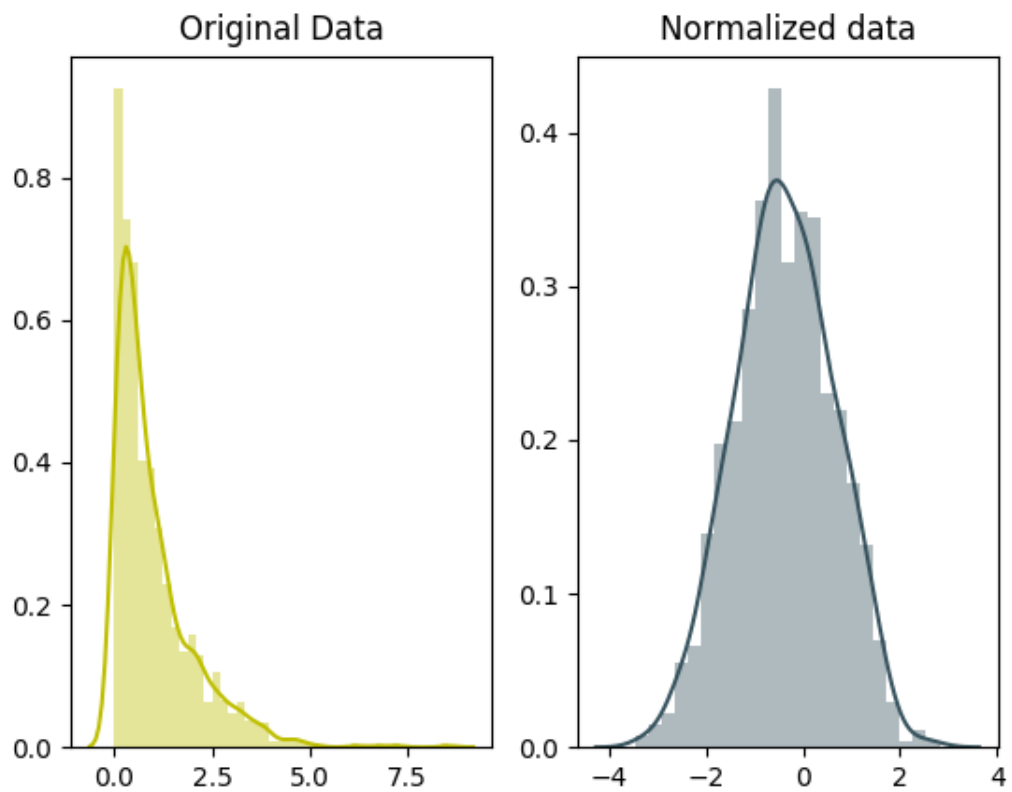
Normalizing input (LeCun et al 1998 “Efficient Backprop”)

Let's understand through some examples:

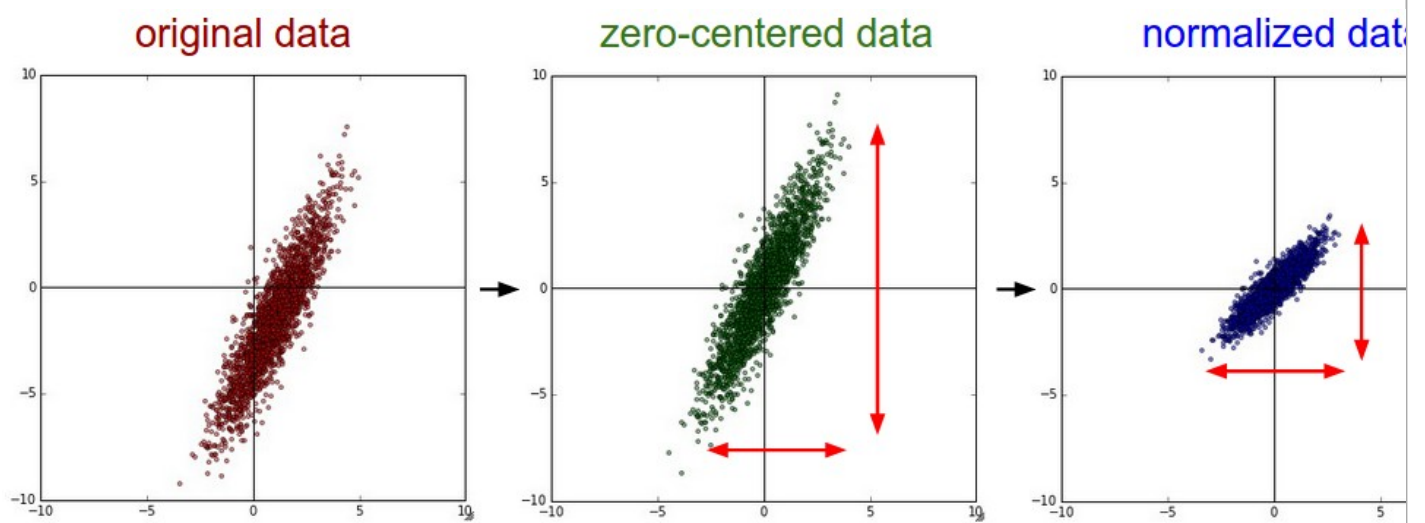


Let's look at normalization working on other kinds of data.





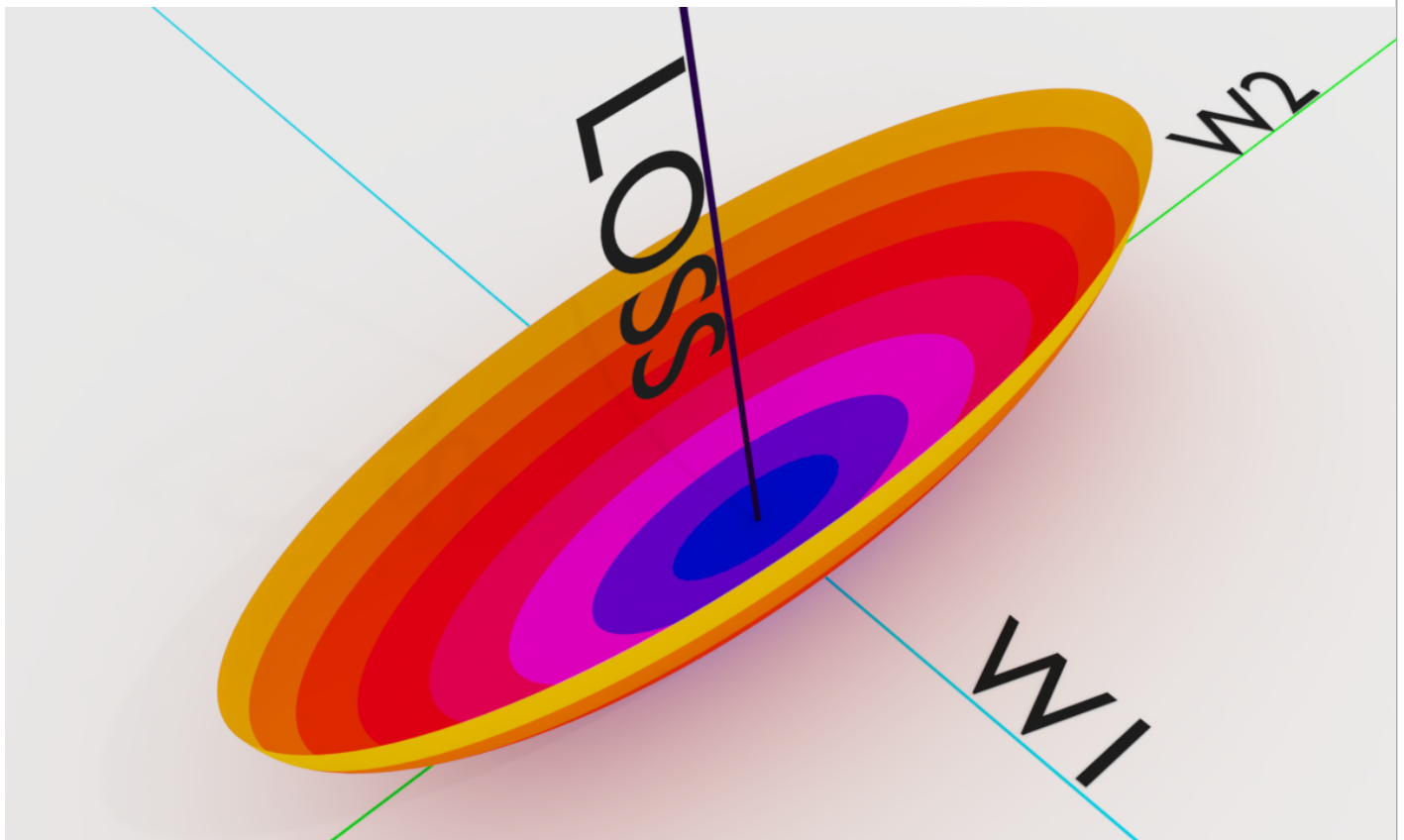
How do we achieve normalization?



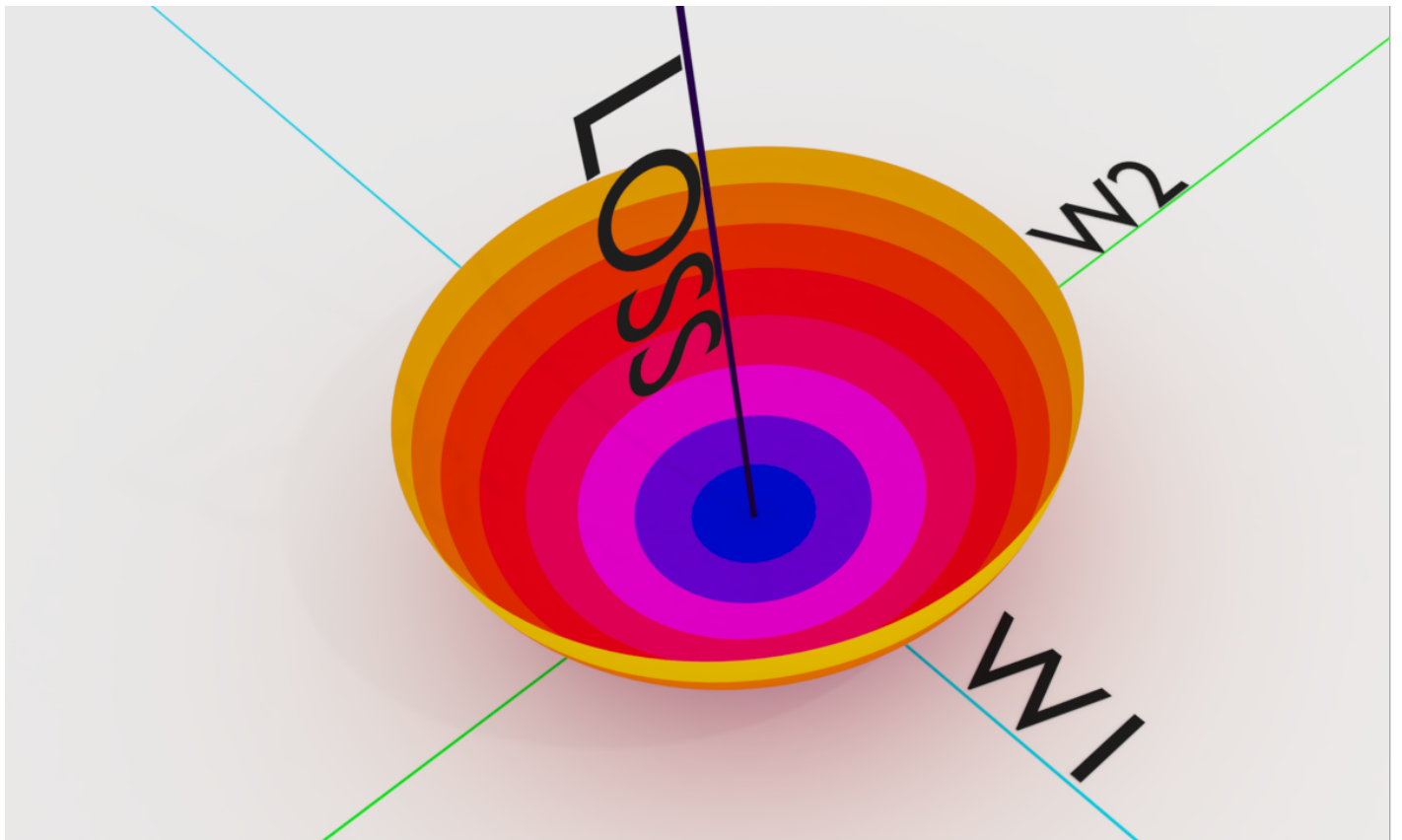
Why limit normalization to images only then?

Channel Normalization -aka Batch Normalization

Let's think about our un-normalized kernels and how loss function would look like:



If we had normalized our kernels (indirectly channels), this is how it would look like:



If features are found at a similar scale, then weights would be in a similar scale, and then backprop would actually make sense, ponder!

Batch Normalization (<http://jmlr.org/proceedings/papers/v37/ioffe15.pdf>) (BN), introduced in 2015– and is now the defacto standard for all CNNs and RNNs.

Batch Normalization solves a problem called **Internal Covariate shift**. To understand BN we need to understand what is CS.

Covariate means input features. Covariate shift means that the distribution of the features is different in different parts of the training/test data.

Internal Covariate shift refers to changes within the neural network, between layers. *A kernel always giving out higher activation makes next layer kernels always expect this higher activation and so on.*

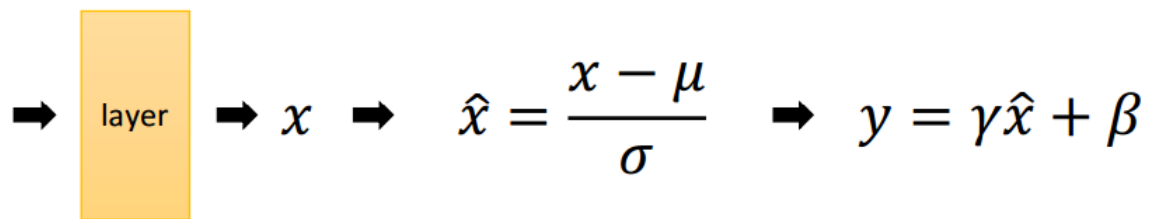
Imagine what would happen if One channel ranges between -1 to 1 and another between -10 to 10

Very Deep nets can be trained faster and generalize better when the distribution of activations is kept normalized during BackProp.

We regularly see Ultra-Deep ConvNets like Inception, Highway Networks, and ResNet. And giant RNNs for speech recognition, [machine translation](https://research.googleblog.com/2016/09/a-neural-network-for-machine.html) (<https://research.googleblog.com/2016/09/a-neural-network-for-machine.html>), etc.

Let's look at some math:

Batch Normalization (BN)



- μ : mean of x in mini-batch
- σ : std of x in mini-batch
- γ : scale
- β : shift
- μ, σ : functions of x , analogous to responses
- γ, β : parameters to be learned, analogous to weights

This is how we implement "batch" normalization:

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;

Parameters to be learned: γ, β

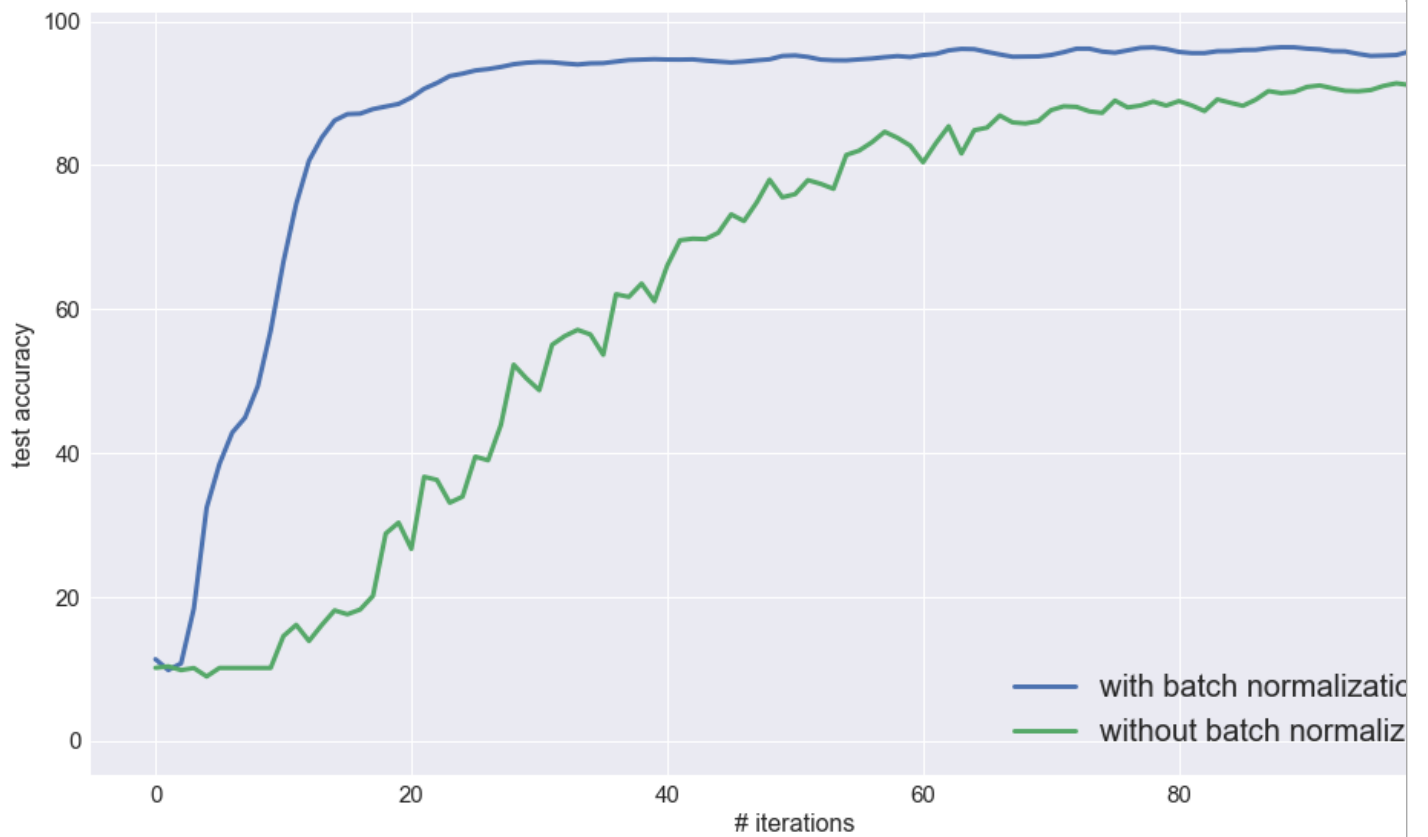
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

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

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

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

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



ReLU before BN or BN before ReLU?

Reference Videos:

Normalizing Activations in a Network (C2W3L04)



Fitting Batch Norm Into Neural Networks (C2W3L05)



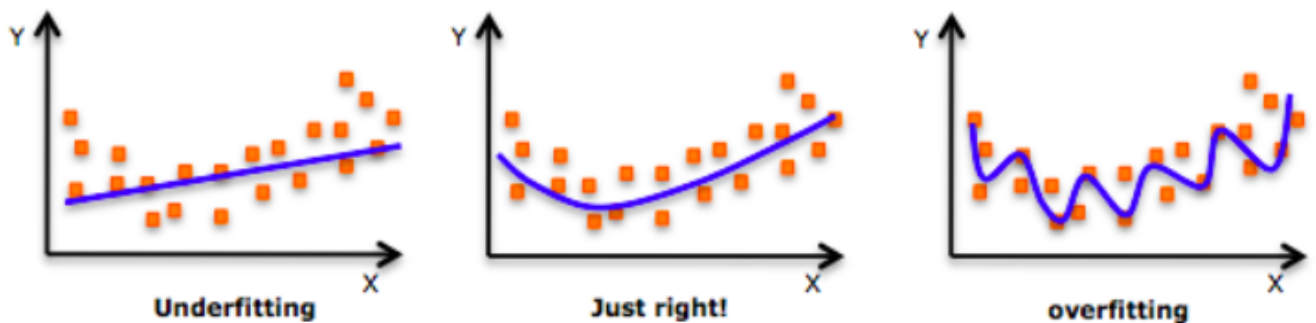
Why Does Batch Norm Work? (C2W3L06)



Batch Norm At Test Time (C2W3L07)

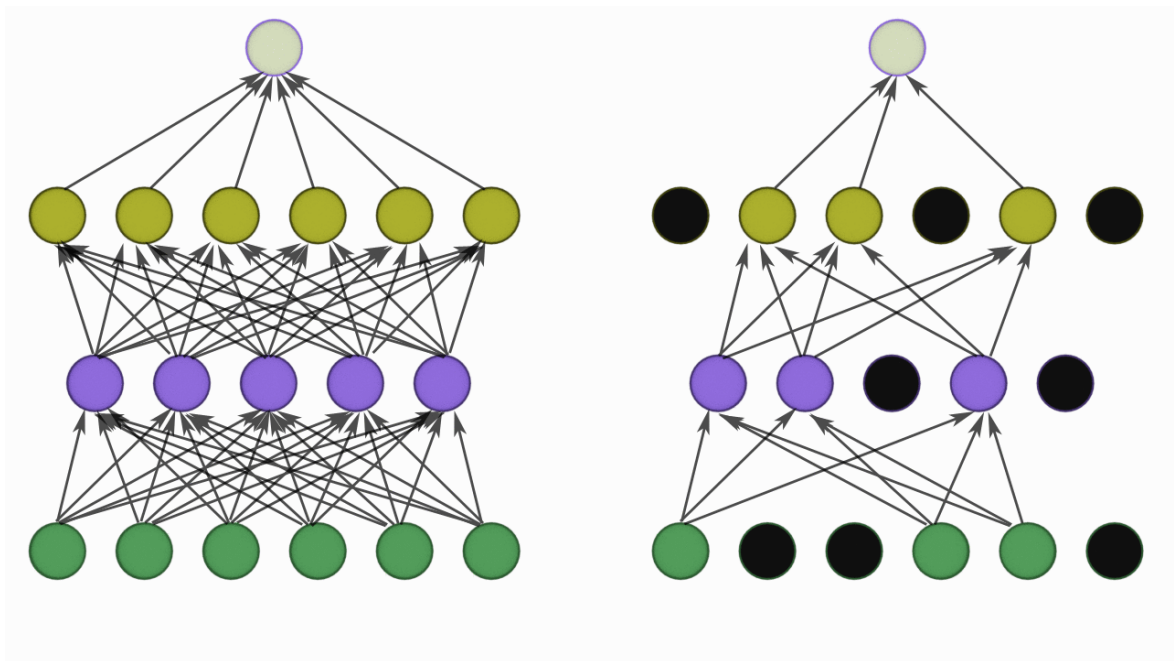


Regularization



Regularization is a key component in preventing overfitting. Also, some techniques of regularization can be used to reduce model parameters while maintaining accuracy, for example, to drive some of the parameters to zero. This might be desirable for reducing the model size or driving down the cost of evaluation in a mobile environment where processor power is constrained.

Drop Out



Most common techniques of regularization used nowadays in the industry are:

- dataset augmentation
- dropout
- weight decay

Recommended:

Why Regularization Reduces Overfitting (C2W1L05)



Dropout Regularization (C2W1L06)



Other Regularization Methods (C2W1L08)



Receptive Fields

$$n_{out} = \left\lfloor \frac{n_{in} + 2p - k}{s} \right\rfloor + 1$$

n_{in} : number of input features

n_{out} : number of output features

k : convolution kernel size

p : convolution padding size

s : convolution stride size

$$1 = (3 + 2*0 - 3)/1 + 1$$

$$3 = (3 + 2*1 - 3)/1 + 1$$

$$3 = (5 + 2*1 - 3)/2 + 1$$

ONLINE CALCULATOR [_ \(https://fomoro.com/research/article/receptive-field-calculator\)](https://fomoro.com/research/article/receptive-field-calculator)

Receptive Field

$$n_{out} = \left\lfloor \frac{n_{in} + 2p - k}{s} \right\rfloor + 1$$

$$j_{out} = j_{in} * s$$

$$r_{out} = r_{in} + (k - 1) * j_{in}$$

$$start_{out} = start_{in} + \left(\frac{k - 1}{2} - p \right) * j_{in}$$

j is the jump, r is the receptive field (ignore start for now)

$$RF1 = 1 + (3-1)*1 = 3$$

$$RF2 = 1 + (3-1)*(1*1) = 3 \text{ (padding has no effect)}$$

$$RF3 = 1 + (3-1)*(1*2) = 4 \text{ (stride has effect!)}$$

The network is 10x10 | 8x8 | 6x6. Our RF currently is 5x5. IF we apply MaxPooling then

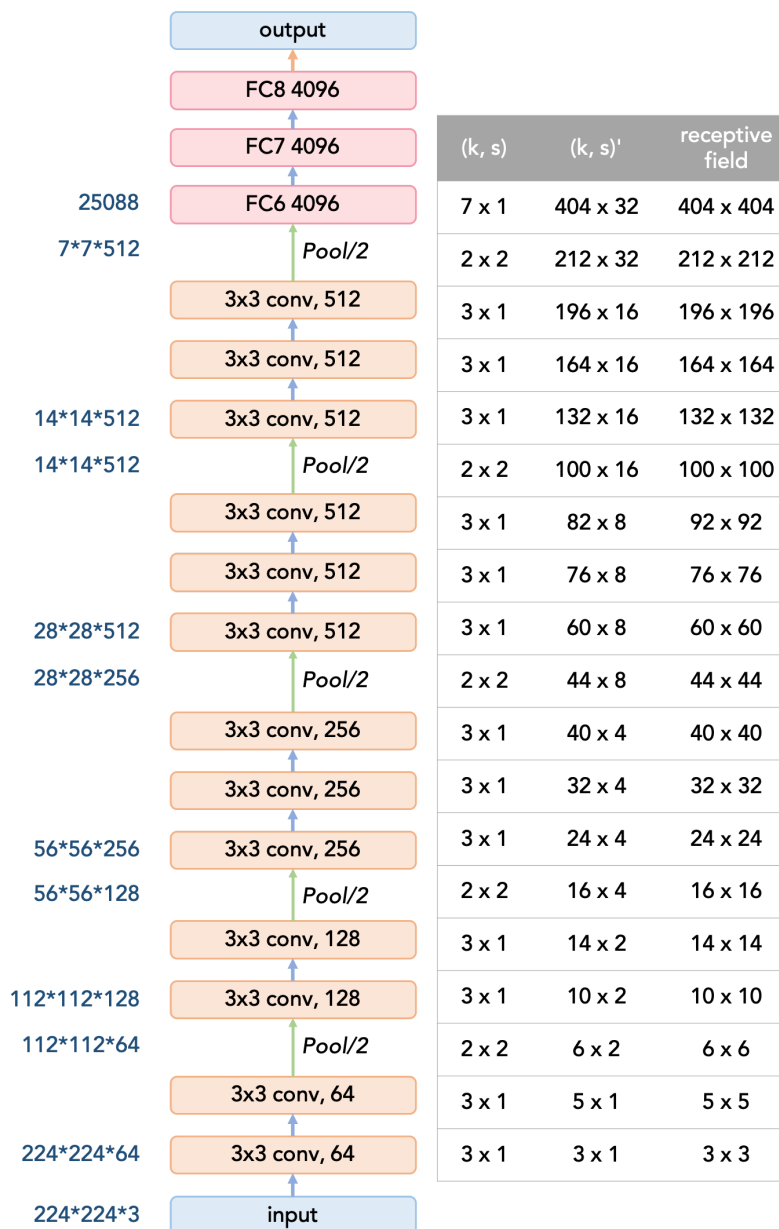
$$\text{First Formula: } n_{out} = (6 + 2*0 - 2)/2 = 3$$

$$\text{Second Formula: } RF = 5 + (2-1)*1 = 6$$

Add 3x3 again

$$\text{First Formula: } (3 + 2*0 - 3) + 1 = 1$$

$$\text{Second Formula: } RF = 6 + (3-1)*2 = 10$$



Source [_ \(http://zike.io/posts/calculate-receptive-field-for-vgg-16/\)](http://zike.io/posts/calculate-receptive-field-for-vgg-16/)

Assignment 3:

1. Run [this](https://colab.research.google.com/drive/1STOg33u7haqSptyjUL40FZlxNW4XdBQK) [_ \(https://colab.research.google.com/drive/1STOg33u7haqSptyjUL40FZlxNW4XdBQK\)](https://colab.research.google.com/drive/1STOg33u7haqSptyjUL40FZlxNW4XdBQK) network (base network) for 50 epochs, report Validation Accuracy **after** 50 epochs.
2. Add new cells at the bottom of the code, and write your own network such that:
 1. it uses depthwise separable convolution **ONLY** (no Conv2D)
 2. it uses BatchNormalization
 3. has less than 100,000 parameters
 4. it uses proper dropout values
 5. you've mentioned the output size for each layer
 6. you've mentioned the receptive field for each layer
 7. runs for 50 epochs
 8. beats the validation score within 50 epochs (at any epoch run, doesn't need to be final one)

3. Submit the github link
4. Your github links must have:
 1. Assignment 3.ipynb file with your code as well as logs for both the models
 2. ReadMe.md which should have:
 1. Final Validation accuracy for Base Network
 2. Your model definition (model.add...) with output channel size and receptive field
 3. Your 50 epoch logs

Points 2,000

Submitting a website url

Due	For	Available from	Until
Nov 29 at 5:30am	F6	Nov 22 at 6:30am	Nov 29 at 5:30am
Nov 30 at 11:30am	S12	Nov 23 at 12:30pm	Nov 30 at 11:30am
Dec 2 at 5:30am	M6	Nov 25 at 6:30am	Dec 2 at 5:30am
Dec 4 at 5:30am	W6	Nov 27 at 6:30am	Dec 4 at 5:30am

+ [Rubric](#)