

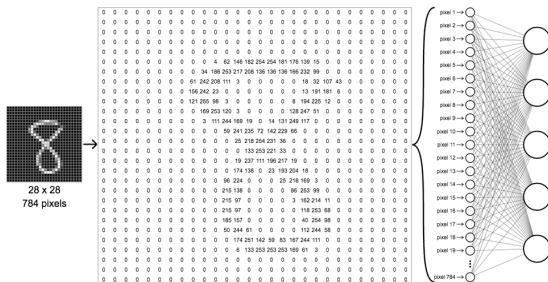
11-695: AI Engineering

Convolution I

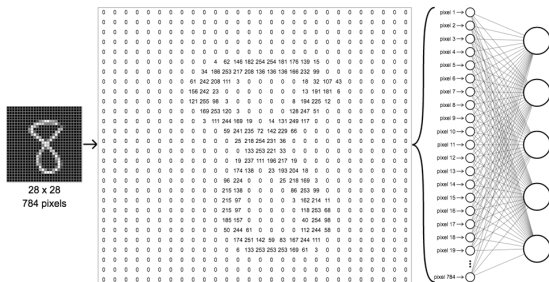
LTI/SCS

Spring 2020

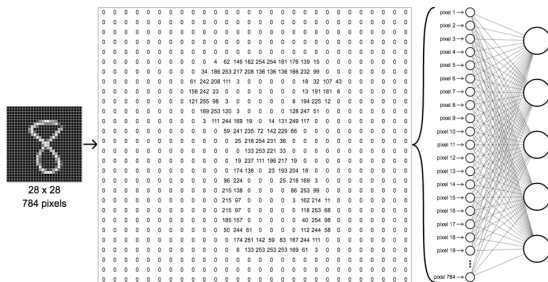
- 1 Motivation from Feedforward NNs
- 2 Convolution
- 3 Convolution in Neural Network
- 4 Pooling
- 5 Case Studies



- The first layer of a NN
- Number of params? $28 * 28 * n$



- The first layer of a NN
- Number of params? $28 * 28 * n$
- Real world: $200 * 200 * 3 * 1000 \approx 1e8$ for the first layer



- The first layer of a NN
- Number of params? $28 * 28 * n$
- Real world: $200 * 200 * 3 * 1000 \approx 1e8$ for the first layer
- What's more? How about spatial correlations?



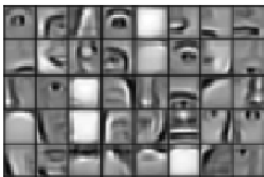
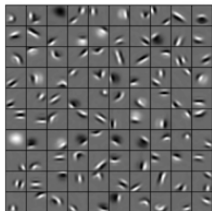
- Local spatial features at different locations might be similar



- Local spatial features at different locations might be similar
- Model parts of the image instead of the whole

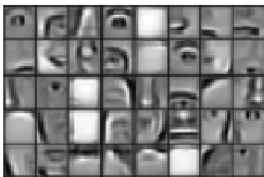
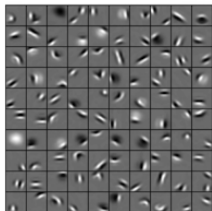


- Local spatial features at different locations might be similar
- Model parts of the image instead of the whole
- Exploit image's redundancy: e.g. edges



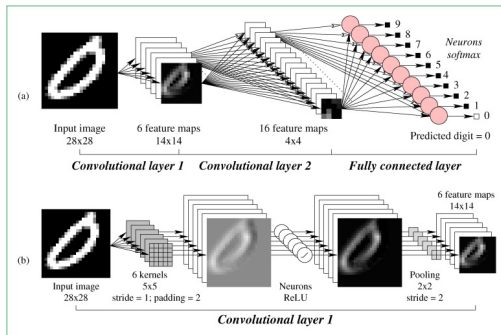
- Human cognition works similarly
 - Eyes detect edges
 - Visual cortex uses **Gabor**-like filters to recognize objects

¹Bengio J. et al. [Representation Learning: A Review and New Perspectives](#) Image credit: Nvidia Developer Blog

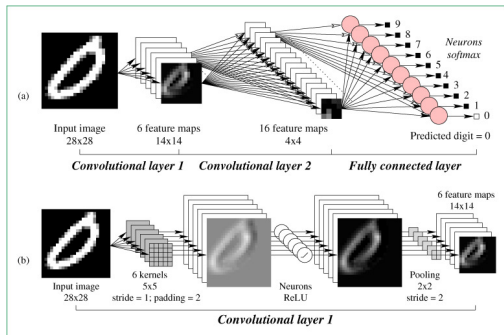


- Human cognition works similarly
 - Eyes detect edges
 - Visual cortex uses **Gabor**-like filters to recognize objects
- Intention to build hierarchical abstract representation¹

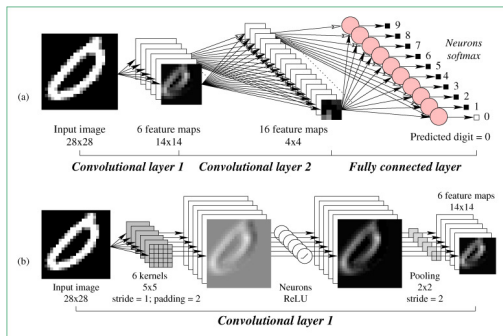
¹Bengio J. et al. [Representation Learning: A Review and New Perspectives](#) Image credit: Nvidia Developer Blog



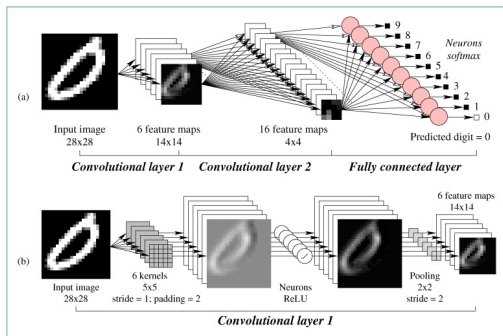
- Similar design to Feedforward NNs



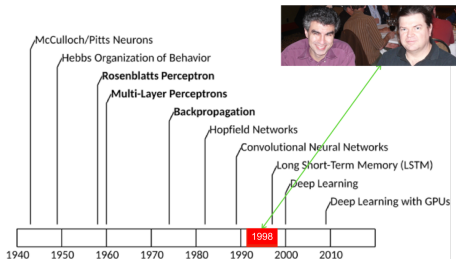
- Similar design to Feedforward NNs
- With main components of not FC (or Dense) layers, but:
 - Convolutional, and



- Similar design to Feedforward NNs
- With main components of not FC (or Dense) layers, but:
 - Convolutional, and
 - Pooling Layers

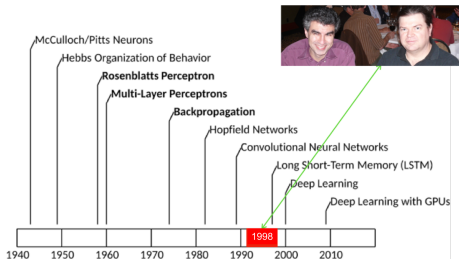


- Similar design to Feedforward NNs
- With main components of not FC (or Dense) layers, but:
 - Convolutional, and
 - Pooling Layers
- FC usually appears at the end of architectures



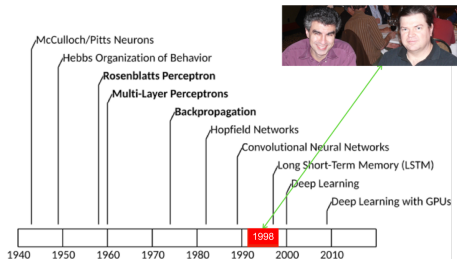
- Be a go-to feature extraction solution for images/videos

²Yin W. et al [Comparative Study of CNN and RNN for NLP](#)



- Be a go-to feature extraction solution for images/videos
- Even outperform RNNs in some NLP tasks!²

²Yin W. et al [Comparative Study of CNN and RNN for NLP](#)



- Be a go-to feature extraction solution for images/videos
- Even outperform RNNs in some NLP tasks!²
- Be an essential part in many SoTA solutions for classifications, detections, recognition, segmentation, OCR, motion, pose, etc.

²Yin W. et al [Comparative Study of CNN and RNN for NLP](#)

- 1 Motivation from Feedforward NNs
- 2 Convolution
- 3 Convolution in Neural Network
- 4 Pooling
- 5 Case Studies

- A mathematics operation [▶ demo](#)

$$(f * g)(t) = \int_{-\infty}^{\infty} f(u)g(t - u)du = \int_{-\infty}^{\infty} f(t - u)g(u)du$$

- A mathematics operation [▶ demo](#)

$$(f * g)(t) = \int_{-\infty}^{\infty} f(u)g(t - u)du = \int_{-\infty}^{\infty} f(t - u)g(u)du$$

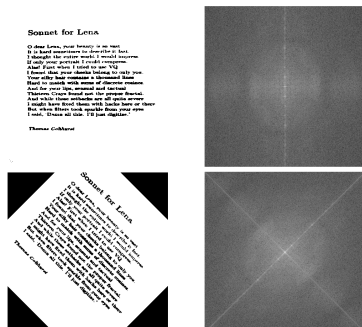
- A linear time invariant (LTI) operation, means that no new frequency components are created, and so output is the pointwise product of input and a transfer function (Wikipedia)

- A mathematics operation [▶ demo](#)

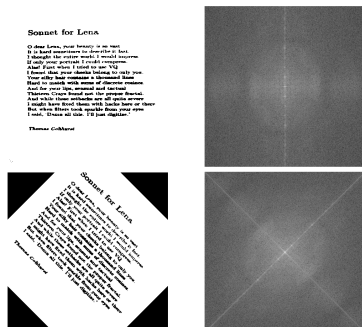
$$(f * g)(t) = \int_{-\infty}^{\infty} f(u)g(t - u)du = \int_{-\infty}^{\infty} f(t - u)g(u)du$$

- A linear time invariant (LTI) operation, means that no new frequency components are created, and so output is the pointwise product of input and a transfer function (Wikipedia)
- Relation to Fourier Transformation [▶ demo](#)

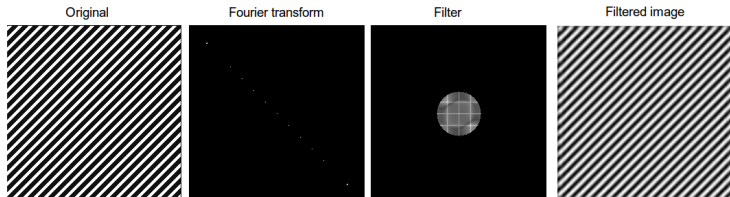
$$(f * g)(t) = \mathcal{F}^{-1}(\sqrt{2\pi} \mathcal{F}|f| \cdot \mathcal{F}|g|)$$



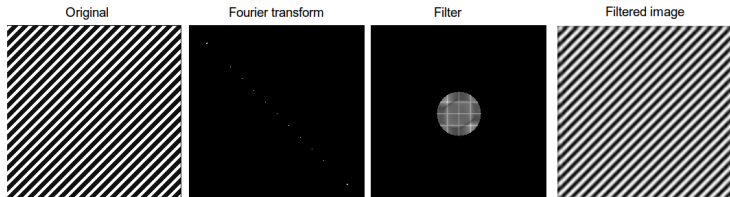
- Convolution is an operation in Fourier domain



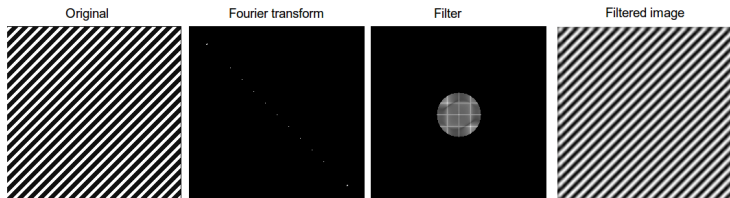
- Convolution is an operation in Fourier domain
- It can capture orientation change



- Convolution is often referred to as filtering



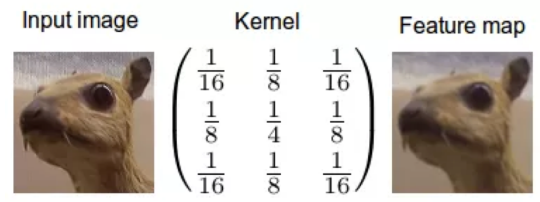
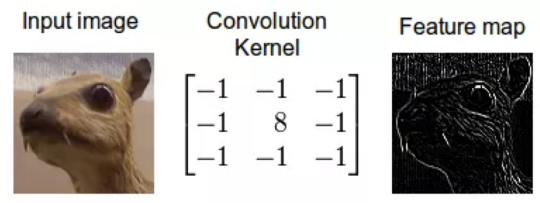
- Convolution is often referred to as filtering
- Other names: kernel, receptive field



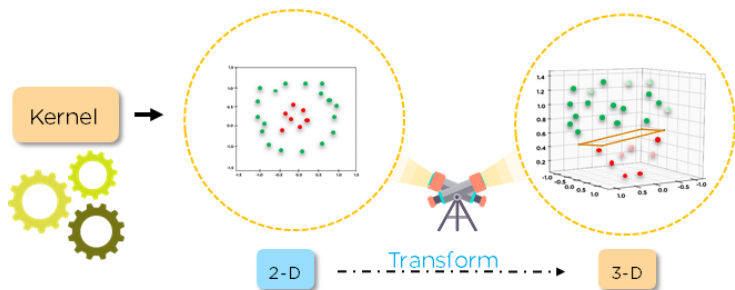
- Convolution is often referred to as filtering
- Other names: kernel, receptive field
- And now:

feature_map = filter * input

$$\begin{aligned} &= \sum_{y=1}^{\text{n_cols}} \left(\sum_{x=1}^{\text{n_rows}} \text{input}(x - a, y - b) * \text{filter}(x, y) \right) \\ &= \mathcal{F}^{-1}(\sqrt{2\pi} \mathcal{F}|\text{input}| \cdot \mathcal{F}|\text{filter}|) \end{aligned}$$

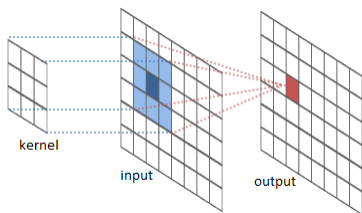


- Similarity: Kernel SVM, PCA?

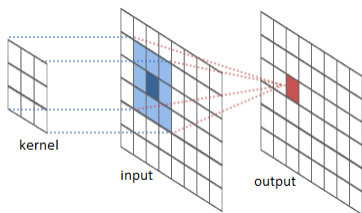


- Key similar intuition: kernel-activated feature space will be helpful.

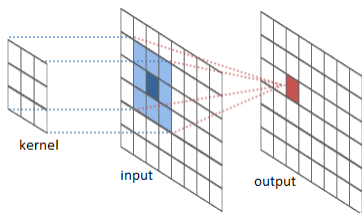
- 1 Motivation from Feedforward NNs
- 2 Convolution
- 3 Convolution in Neural Network**
- 4 Pooling
- 5 Case Studies



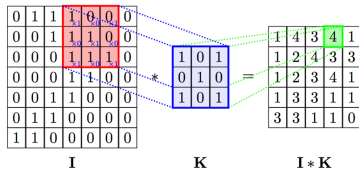
- Given a $200 * 200$ image
- A fully-connected layer with n hidden neurons: $200 * 200 * n$ params



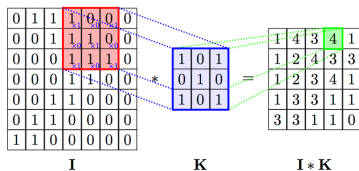
- Given a $200 * 200$ image
- A fully-connected layer with n hidden neurons: $200 * 200 * n$ params
- With convolution with filter size $3 * 3$, then how many params?



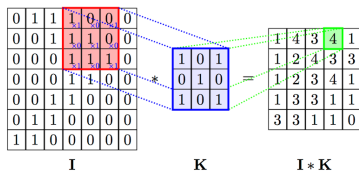
- Given a $200 * 200$ image
- A fully-connected layer with n hidden neurons: $200 * 200 * n$ params
- With convolution with filter size $3 * 3$, then how many params?
 $3 * 3 * n$



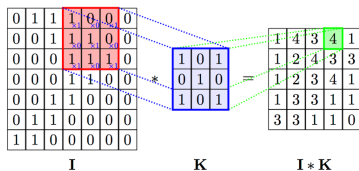
- How it works?



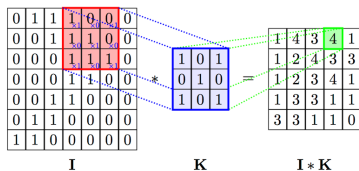
- How it works? Simply, it's a dot product



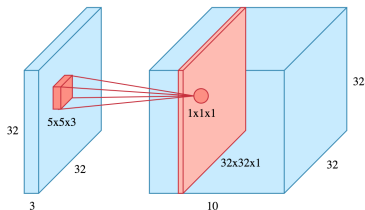
- How it works? Simply, it's a dot product
- A filter scan through the whole image which is convolved by it
- *Important:* a filter has a fixed weight \rightarrow output is similar if the region is similar



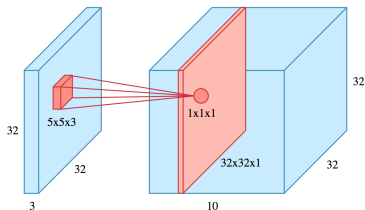
- How it works? Simply, it's a dot product
- A filter scan through the whole image which is convolved by it
- *Important:* a filter has a fixed weight \rightarrow output is similar if the region is similar
- Practice: $n * n$ image, $k * k$ filter, output?



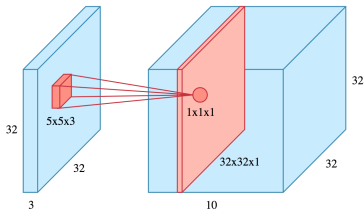
- How it works? Simply, it's a dot product
- A filter scan through the whole image which is convolved by it
- *Important:* a filter has a fixed weight \rightarrow output is similar if the region is similar
- Practice: $n * n$ image, $k * k$ filter, output? $(n - k + 1) * (n - k + 1)$



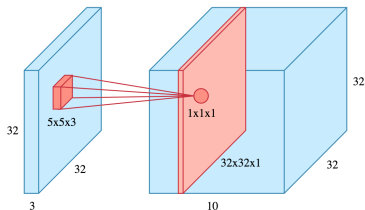
- Filter shares the *same depth* with input
- If input image has 3 channels: $200 * 200 * 3 \rightarrow$ the filter should be $n * n * 3$



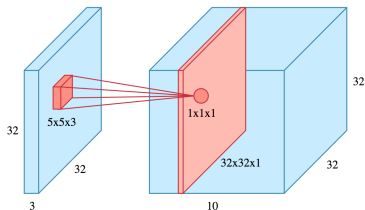
- Filter shares the *same depth* with input
- If input image has 3 channels: $200 * 200 * 3 \rightarrow$ the filter should be $n * n * 3$
- Dimension of each convolution operation?



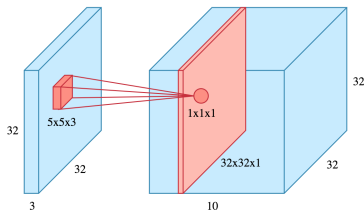
- Filter shares the *same depth* with input
- If input image has 3 channels: $200 * 200 * 3 \rightarrow$ the filter should be $n * n * 3$
- Dimension of each convolution operation? Simply, a scalar.



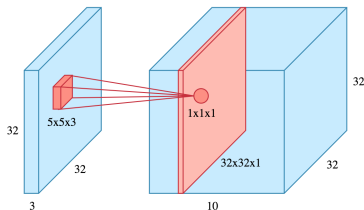
- Filter shares the *same depth* with input
- If input image has 3 channels: $200 * 200 * 3 \rightarrow$ the filter should be $n * n * 3$
- Dimension of each convolution operation? Simply, a scalar.
- Recall, what if we want many feature maps?



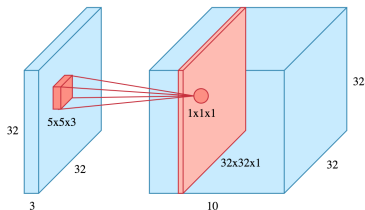
- Filter shares the *same depth* with input
- If input image has 3 channels: $200 * 200 * 3 \rightarrow$ the filter should be $n * n * 3$
- Dimension of each convolution operation? Simply, a scalar.
- Recall, what if we want many feature maps? Use many filters.
- Practice: $n * n * 3$ image, $k * k * 3$ filter, output?



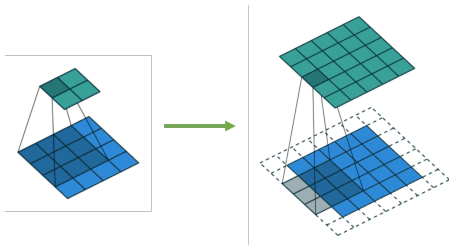
- Filter shares the *same depth* with input
- If input image has 3 channels: $200 * 200 * 3 \rightarrow$ the filter should be $n * n * 3$
- Dimension of each convolution operation? Simply, a scalar.
- Recall, what if we want many feature maps? Use many filters.
- Practice: $n * n * 3$ image, $k * k * 3$ filter, output?
 $(n - k + 1) * (n - k + 1) * 3$



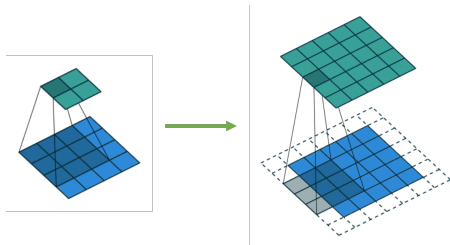
- Filter shares the *same depth* with input
- If input image has 3 channels: $200 * 200 * 3 \rightarrow$ the filter should be $n * n * 3$
- Dimension of each convolution operation? Simply, a scalar.
- Recall, what if we want many feature maps? Use many filters.
- Practice: $n * n * 3$ image, $k * k * 3$ filter, output?
 $(n - k + 1) * (n - k + 1) * 1$



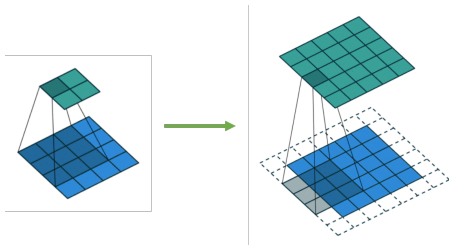
- Filter shares the *same depth* with input
- If input image has 3 channels: $200 * 200 * 3 \rightarrow$ the filter should be $n * n * 3$
- Dimension of each convolution operation? Simply, a scalar.
- Recall, what if we want many feature maps? Use many filters.
- Practice: $n * n * 3$ image, $k * k * 3$ filter, output?
 $(n - k + 1) * (n - k + 1) * 1$
- Practice: If there are 10 filters?



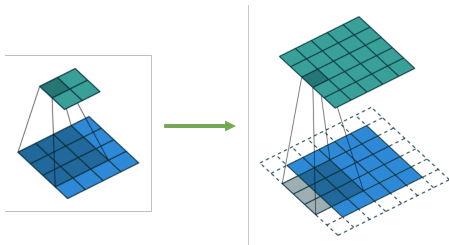
- Change the size of output, normally with zero [▶ Demo](#)



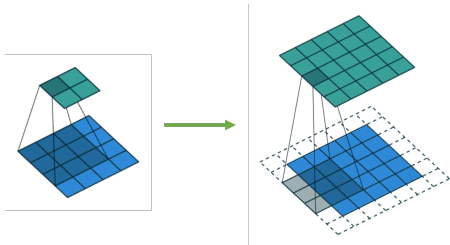
- Change the size of output, normally with zero [▶ Demo](#)
- There are 2 types of padding
 - *Valid*: no padding



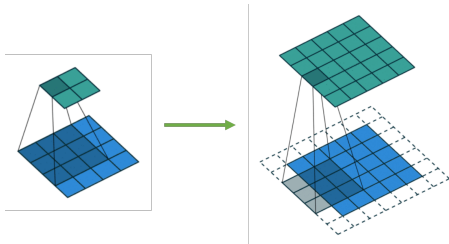
- Change the size of output, normally with zero [▶ Demo](#)
- There are 2 types of padding
 - *Valid*: no padding
 - *Same*: pad so that output size is the same as input's



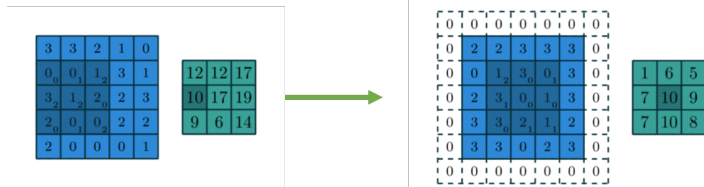
- Change the size of output, normally with zero [▶ Demo](#)
- There are 2 types of padding
 - *Valid*: no padding
 - *Same*: pad so that output size is the same as input's
- Practice: $n * n$ convolves $k * k$, padding p , output?



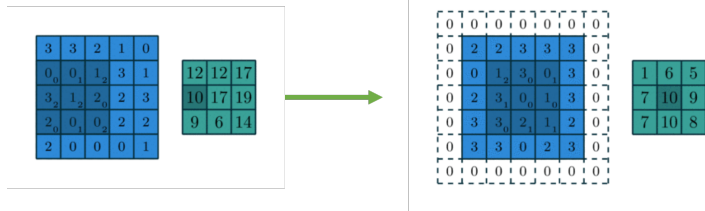
- Change the size of output, normally with zero [▶ Demo](#)
- There are 2 types of padding
 - *Valid*: no padding
 - *Same*: pad so that output size is the same as input's
- Practice: $n * n$ convolves $k * k$, padding p , output?
 $(n - k + 1 + 2p) * (n - k + 1 + 2p)$



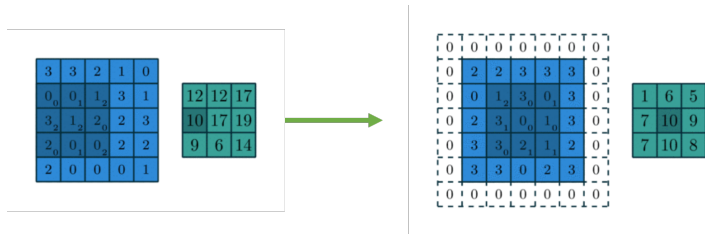
- Change the size of output, normally with zero [▶ Demo](#)
- There are 2 types of padding
 - *Valid*: no padding
 - *Same*: pad so that output size is the same as input's
- Practice: $n * n$ convolves $k * k$, padding p , output?
 $(n - k + 1 + 2p) * (n - k + 1 + 2p)$
- Assumption of this formula (and the ones above)?



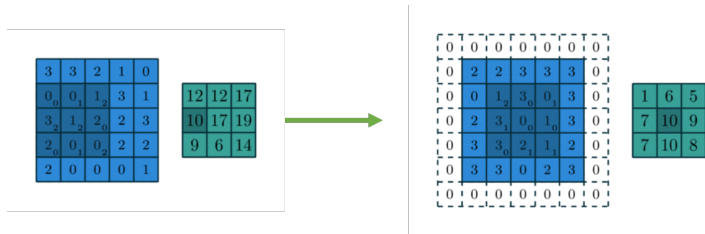
- Velocity of convolution [▶ Demo](#)
- Practice: 5x5 convolves 3x3
 - Valid padding, stride 1, output?



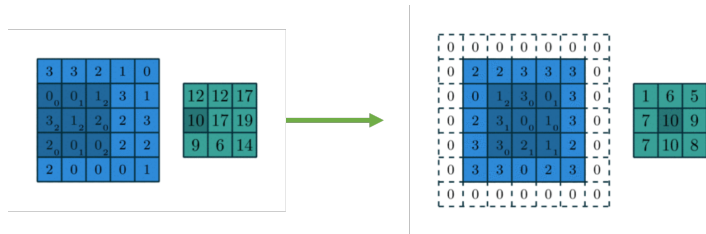
- Velocity of convolution [▶ Demo](#)
- Practice: 5x5 convolves 3x3
 - Valid padding, stride 1, output? $(5 - 3 + 1)(5 - 3 + 1)$
 - Pad 1, stride 1, output?



- Velocity of convolution [▶ Demo](#)
- Practice: 5x5 convolves 3x3
 - Valid padding, stride 1, output? $(5 - 3 + 1)(5 - 3 + 1)$
 - Pad 1, stride 1, output? $(5 - 3 + 1 + 2)(5 - 3 + 1 + 2)$
 - Pad 1, stride 2, output?

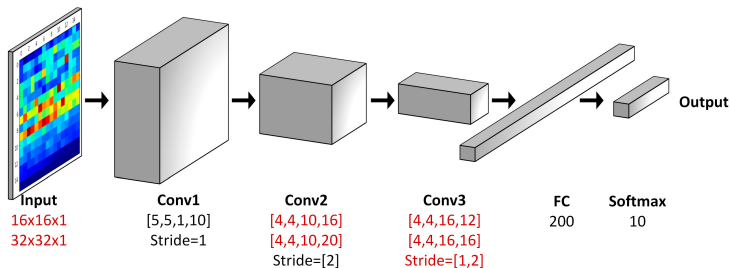


- Velocity of convolution [▶ Demo](#)
- Practice: 5x5 convolves 3x3
 - Valid padding, stride 1, output? $(5 - 3 + 1)(5 - 3 + 1)$
 - Pad 1, stride 1, output? $(5 - 3 + 1 + 2)(5 - 3 + 1 + 2)$
 - Pad 1, stride 2, output? $((5 - 3 + 2)/2 + 1)((5 - 3 + 2)/2 + 1)$

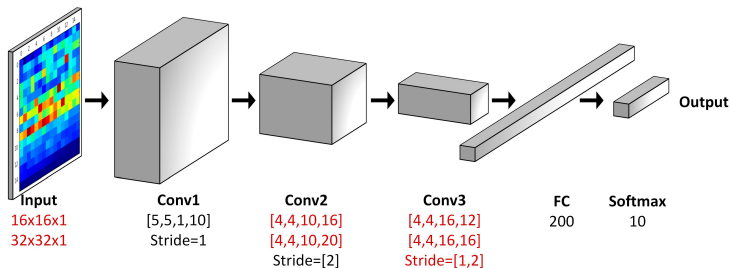


- Velocity of convolution [▶ Demo](#)
- Practice: 5x5 convolves 3x3
 - Valid padding, stride 1, output? $(5 - 3 + 1)(5 - 3 + 1)$
 - Pad 1, stride 1, output? $(5 - 3 + 1 + 2)(5 - 3 + 1 + 2)$
 - Pad 1, stride 2, output? $((5 - 3 + 2)/2 + 1)((5 - 3 + 2)/2 + 1)$
 - Which one is Same padding?
- So the formula is: $(\frac{n-k+2p}{2} + 1)$, Note: n can vary by dimensions
- Practice: *Same* convolution, how p relates to k ?

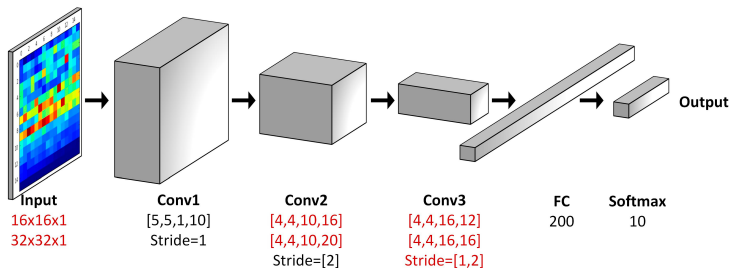
- 1 Motivation from Feedforward NNs
- 2 Convolution
- 3 Convolution in Neural Network
- 4 Pooling
- 5 Case Studies



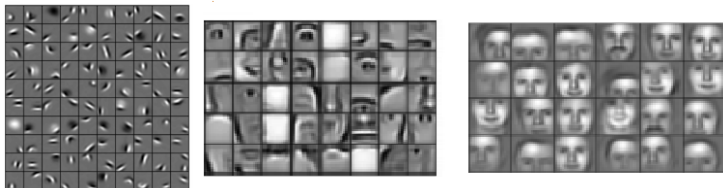
- Practice: given valid conv and zero padding, identify filter size at each step? Then calculate number of parameters at each step?



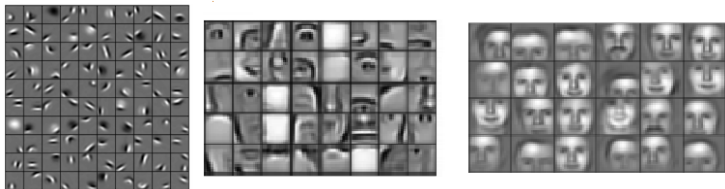
- Practice: given valid conv and zero padding, identify filter size at each step? Then calculate number of parameters at each step?
- Can we go deeper with this? Imagine the real world cases of 200x200, or even 1024x1024.



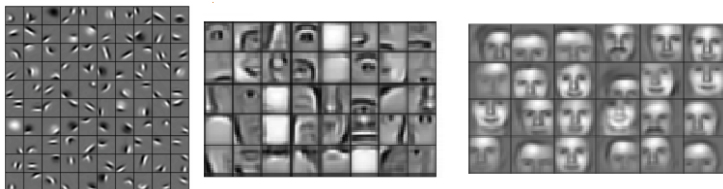
- Practice: given valid conv and zero padding, identify filter size at each step? Then calculate number of parameters at each step?
- Can we go deeper with this? Imagine the real world cases of 200x200, or even 1024x1024.
- We need to reduce dimensionality



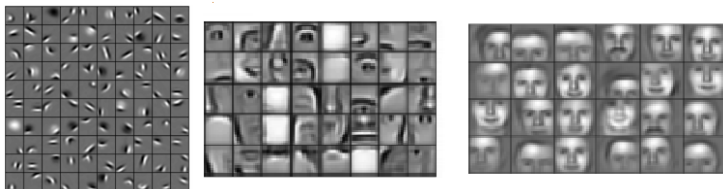
- We want to reduce size of feature maps to easily control



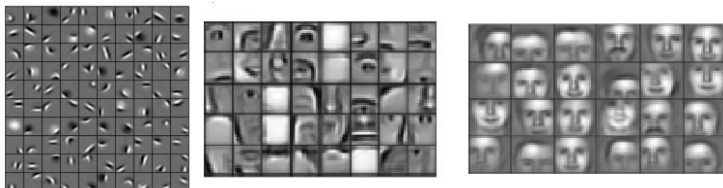
- We want to reduce size of feature maps to easily control
- To prevent loss, we *summarize* features:
 - Each node (scalar) in a feature map represents features of a kernel-size region



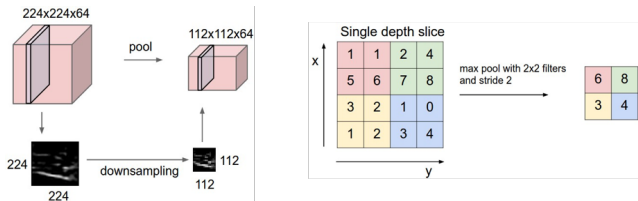
- We want to reduce size of feature maps to easily control
- To prevent loss, we *summarize* features:
 - Each node (scalar) in a feature map represents features of a kernel-size region
 - Adjacent nodes represent a local region's features



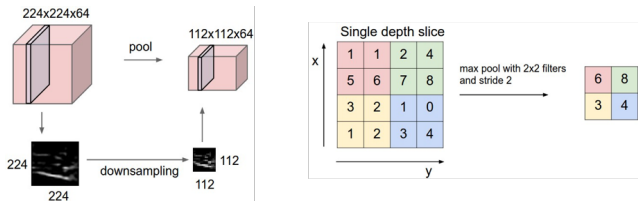
- We want to reduce size of feature maps to easily control
- To prevent loss, we *summarize* features:
 - Each node (scalar) in a feature map represents features of a kernel-size region
 - Adjacent nodes represent a local region's features
 - Summarize adjacent nodes, we have features of a 'big' local region



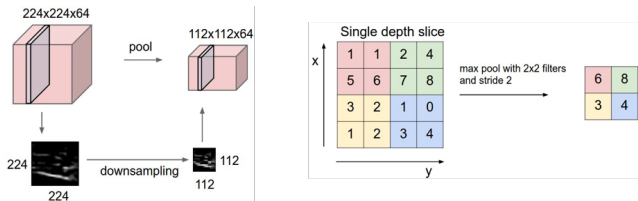
- We want to reduce size of feature maps to easily control
- To prevent loss, we *summarize* features:
 - Each node (scalar) in a feature map represents features of a kernel-size region
 - Adjacent nodes represent a local region's features
 - Summarize adjacent nodes, we have features of a 'big' local region
→ *hierarchy*
- To get *translation invariance* for higher object levels



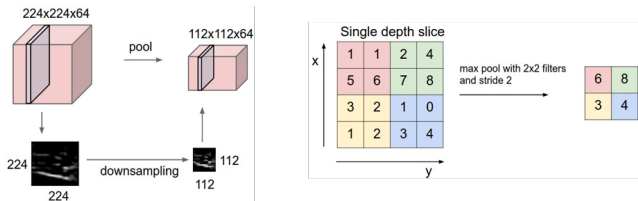
- Process each feature map *independently*
- *a.k.a* subsampling or downsampling process.



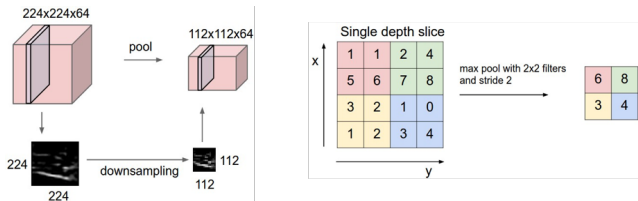
- Process each feature map *independently*
- *a.k.a* subsampling or downsampling process.
- The window is also called a kernel/filter
- Popular types: MAX, AVG
- Practice: 5x5 with filter 3x3, stride 1, output?



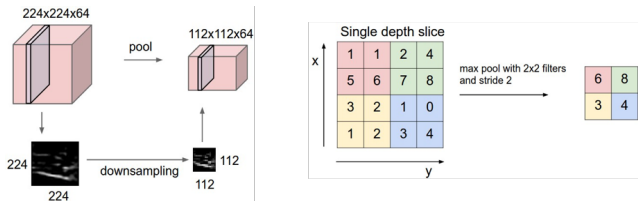
- Process each feature map *independently*
- *a.k.a* subsampling or downsampling process.
- The window is also called a kernel/filter
- Popular types: MAX, AVG
- Practice: 5×5 with filter 3×3 , stride 1, output? 3×3
- Practice: And how many params?



- Process each feature map *independently*
- *a.k.a* subsampling or downsampling process.
- The window is also called a kernel/filter
- Popular types: MAX, AVG
- Practice: 5x5 with filter 3x3, stride 1, output? 3x3
- Practice: And how many params? Zero

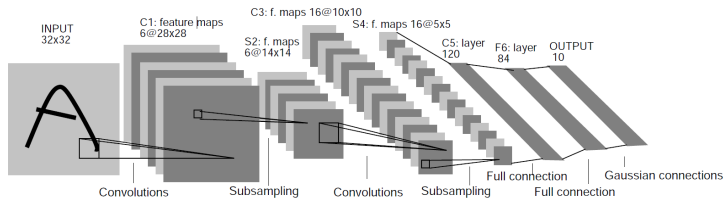


- Process each feature map *independently*
- *a.k.a* subsampling or downsampling process.
- The window is also called a kernel/filter
- Popular types: MAX, AVG
- Practice: 5×5 with filter 3×3 , stride 1, output? 3×3
- Practice: And how many params? Zero, and None padding.
- Output formula?

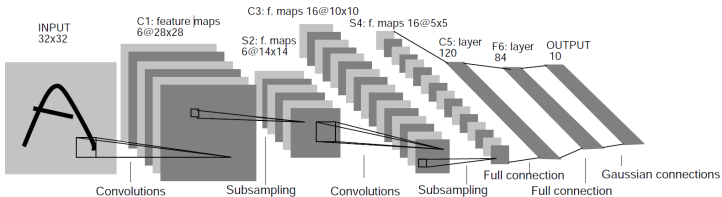


- Process each feature map *independently*
- *a.k.a* subsampling or downsampling process.
- The window is also called a kernel/filter
- Popular types: MAX, AVG
- Practice: 5x5 with filter 3x3, stride 1, output? 3x3
- Practice: And how many params? Zero, and None padding.
- Output formula? $(n - k)/s + 1$

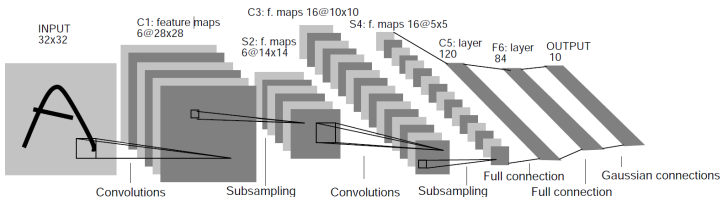
- 1 Motivation from Feedforward NNs
- 2 Convolution
- 3 Convolution in Neural Network
- 4 Pooling
- 5 Case Studies



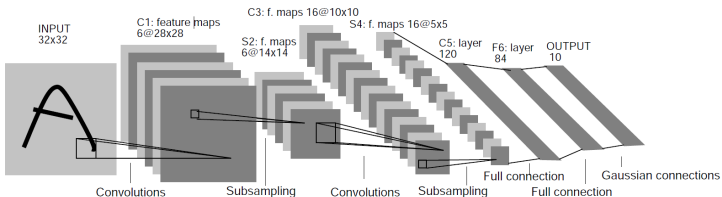
- Practice: what are the sizes of kernels and strides for Conv, and for Pooling?



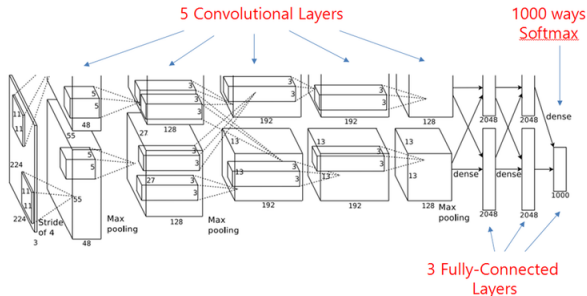
- Practice: what are the sizes of kernels and strides for Conv, and for Pooling? conv: 5×5 and 1, pool: 2×2 and 2
- Practice: how many params at each step?



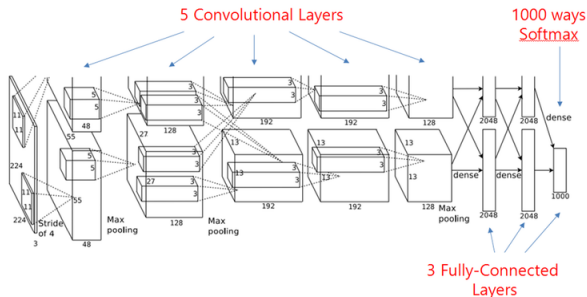
- Arguably the first CNN that *really* works [▶ Demo](#)
- Practice: what are the sizes of kernels and strides for Conv, and for Pooling?



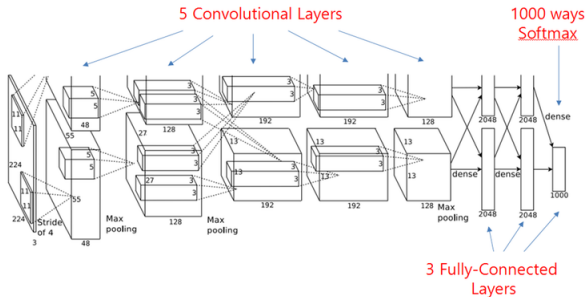
- Arguably the first CNN that *really* works [▶ Demo](#)
- Practice: what are the sizes of kernels and strides for Conv, and for Pooling? conv: $5 * 5$ and 1, pool: $2 * 2$ and 2
- Practice: how many params at each step?



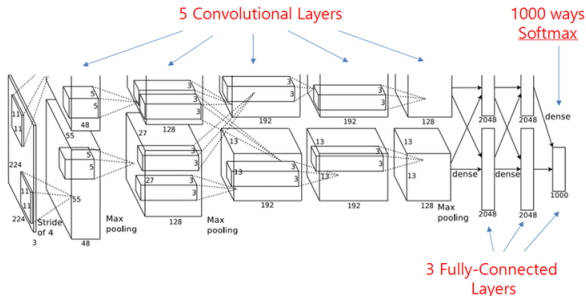
- First use of ReLU, Norm layers, heavy augmentation, dropout, momentum
- Input: $227 * 227 * 3$
- Conv1: 96 of $11 * 11$ filters, stride 4, output?



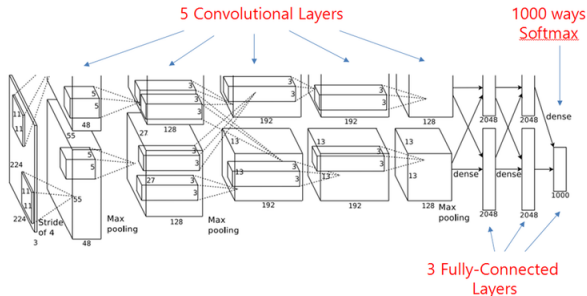
- First use of ReLU, Norm layers, heavy augmentation, dropout, momentum
- Input: $227 * 227 * 3$
- Conv1: 96 of $11 * 11$ filters, stride 4, output? $55 * 55 * 96$,



- First use of ReLU, Norm layers, heavy augmentation, dropout, momentum
- Input: $227 * 227 * 3$
- Conv1: 96 of $11 * 11$ filters, stride 4, output? $55 * 55 * 96$, params?



- First use of ReLU, Norm layers, heavy augmentation, dropout, momentum
- Input: $227 * 227 * 3$
- Conv1: 96 of $11 * 11$ filters, stride 4, output? $55 * 55 * 96$, params? $(11 * 11 * 3) * 96$
- Pool1: $3 * 3$ filter, stride 2, output?



- First use of ReLU, Norm layers, heavy augmentation, dropout, momentum
- Input: $227 * 227 * 3$
- Conv1: 96 of $11 * 11$ filters, stride 4, output? $55 * 55 * 96$, params? $(11 * 11 * 3) * 96$
- Pool1: $3 * 3$ filter, stride 2, output? $27 * 27 * 96$

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding),

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output?
 $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding),

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding),

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 384) * 384$
- Conv5: 256 of $3 * 3$ filters, stride 1, padding 1, output?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 384) * 384$
- Conv5: 256 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 256$ (same padding),

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 384) * 384$
- Conv5: 256 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 256$ (same padding), params?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 384) * 384$
- Conv5: 256 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 256$ (same padding), params? $(3 * 3 * 384) * 256$
- Pool3: $3 * 3$ filter, stride 2, output?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 384) * 384$
- Conv5: 256 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 256$ (same padding), params? $(3 * 3 * 384) * 256$
- Pool3: $3 * 3$ filter, stride 2, output? $6 * 6 * 256$
- FC6: 4096 neurons, output?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 384) * 384$
- Conv5: 256 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 256$ (same padding), params? $(3 * 3 * 384) * 256$
- Pool3: $3 * 3$ filter, stride 2, output? $6 * 6 * 256$
- FC6: 4096 neurons, output? 4096,

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 384) * 384$
- Conv5: 256 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 256$ (same padding), params? $(3 * 3 * 384) * 256$
- Pool3: $3 * 3$ filter, stride 2, output? $6 * 6 * 256$
- FC6: 4096 neurons, output? 4096, params?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 384) * 384$
- Conv5: 256 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 256$ (same padding), params? $(3 * 3 * 384) * 256$
- Pool3: $3 * 3$ filter, stride 2, output? $6 * 6 * 256$
- FC6: 4096 neurons, output? 4096, params? $6 * 6 * 256 * 4096$
- FC7: 4096 neurons, output?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 384) * 384$
- Conv5: 256 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 256$ (same padding), params? $(3 * 3 * 384) * 256$
- Pool3: $3 * 3$ filter, stride 2, output? $6 * 6 * 256$
- FC6: 4096 neurons, output? 4096, params? $6 * 6 * 256 * 4096$
- FC7: 4096 neurons, output? 4096,

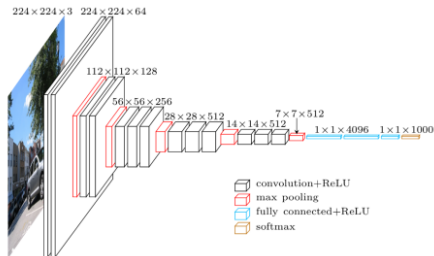
- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 384) * 384$
- Conv5: 256 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 256$ (same padding), params? $(3 * 3 * 384) * 256$
- Pool3: $3 * 3$ filter, stride 2, output? $6 * 6 * 256$
- FC6: 4096 neurons, output? 4096, params? $6 * 6 * 256 * 4096$
- FC7: 4096 neurons, output? 4096, params?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 384) * 384$
- Conv5: 256 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 256$ (same padding), params? $(3 * 3 * 384) * 256$
- Pool3: $3 * 3$ filter, stride 2, output? $6 * 6 * 256$
- FC6: 4096 neurons, output? 4096, params? $6 * 6 * 256 * 4096$
- FC7: 4096 neurons, output? 4096, params? $4096 * 4096$
- FC8: 1000 neurons, output?

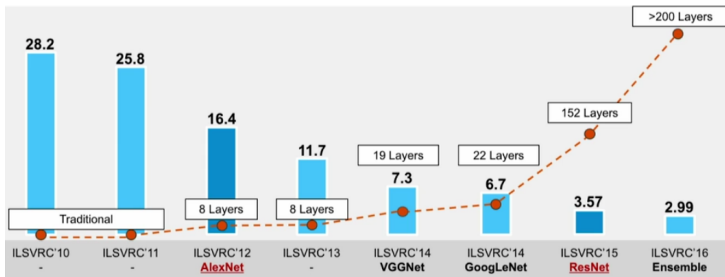
- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 384) * 384$
- Conv5: 256 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 256$ (same padding), params? $(3 * 3 * 384) * 256$
- Pool3: $3 * 3$ filter, stride 2, output? $6 * 6 * 256$
- FC6: 4096 neurons, output? 4096, params? $6 * 6 * 256 * 4096$
- FC7: 4096 neurons, output? 4096, params? $4096 * 4096$
- FC8: 1000 neurons, output? 1000,

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 384) * 384$
- Conv5: 256 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 256$ (same padding), params? $(3 * 3 * 384) * 256$
- Pool3: $3 * 3$ filter, stride 2, output? $6 * 6 * 256$
- FC6: 4096 neurons, output? 4096, params? $6 * 6 * 256 * 4096$
- FC7: 4096 neurons, output? 4096, params? $4096 * 4096$
- FC8: 1000 neurons, output? 1000, params?

- Conv2: 256 of $5 * 5$ filters, stride 1, padding 2, output? $27 * 27 * 256$ (same padding), params? $(5 * 5 * 96) * 256$
- Pool2: $3 * 3$ filter, stride 2, output? $13 * 13 * 256$
- Conv3: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 256) * 384$
- Conv4: 384 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 384$ (same padding), params? $(3 * 3 * 384) * 384$
- Conv5: 256 of $3 * 3$ filters, stride 1, padding 1, output? $13 * 13 * 256$ (same padding), params? $(3 * 3 * 384) * 256$
- Pool3: $3 * 3$ filter, stride 2, output? $6 * 6 * 256$
- FC6: 4096 neurons, output? 4096, params? $6 * 6 * 256 * 4096$
- FC7: 4096 neurons, output? 4096, params? $4096 * 4096$
- FC8: 1000 neurons, output? 1000, params? $4096 * 1000$



- Got 11.2% top 5 error at ILSVRC 2013
- Very widely used for its lightweight yet efficient architecture.
- Popular with 16 and 19 layers
- Reduce params: all fixed 3 * 3 filters, and stride 2 poolings: 138M params (VGG16)



End of Module I.