



Introduction to Machine Learning (CS419M)

Lecture 13:

- Convolutional Neural Networks

Mar 13, 2020

Convolutional Neural Networks (CNNs)

- Fully connected (dense) layers have no awareness of spatial information
- Key concept behind convolutional layers is that of ***kernels*** or ***filters***
- Filters slide across an input space to detect spatial patterns (translation invariance) in local regions (locality)

Convolution Layer

1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

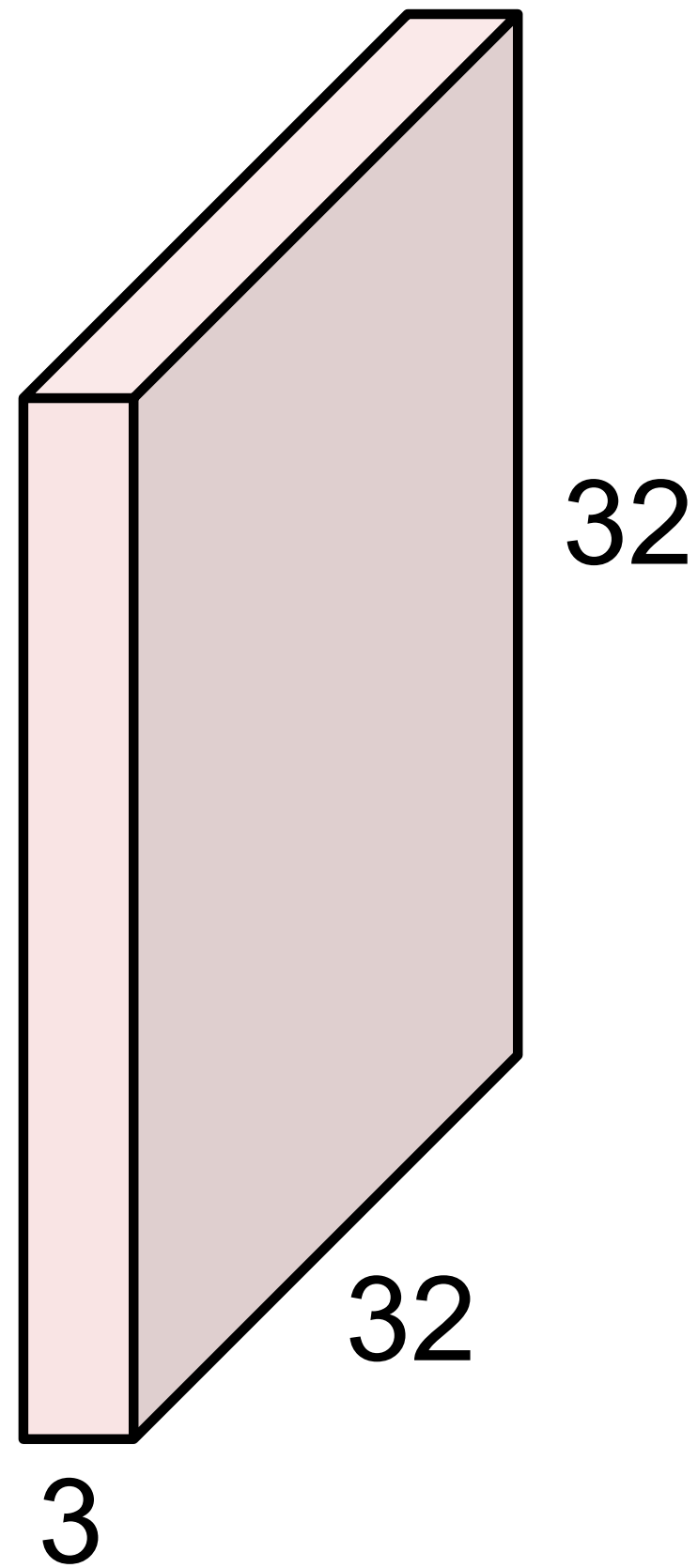
Image

4		

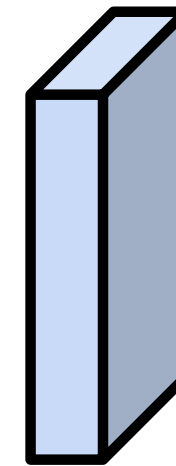
Convolved
Feature

Convolution Layer

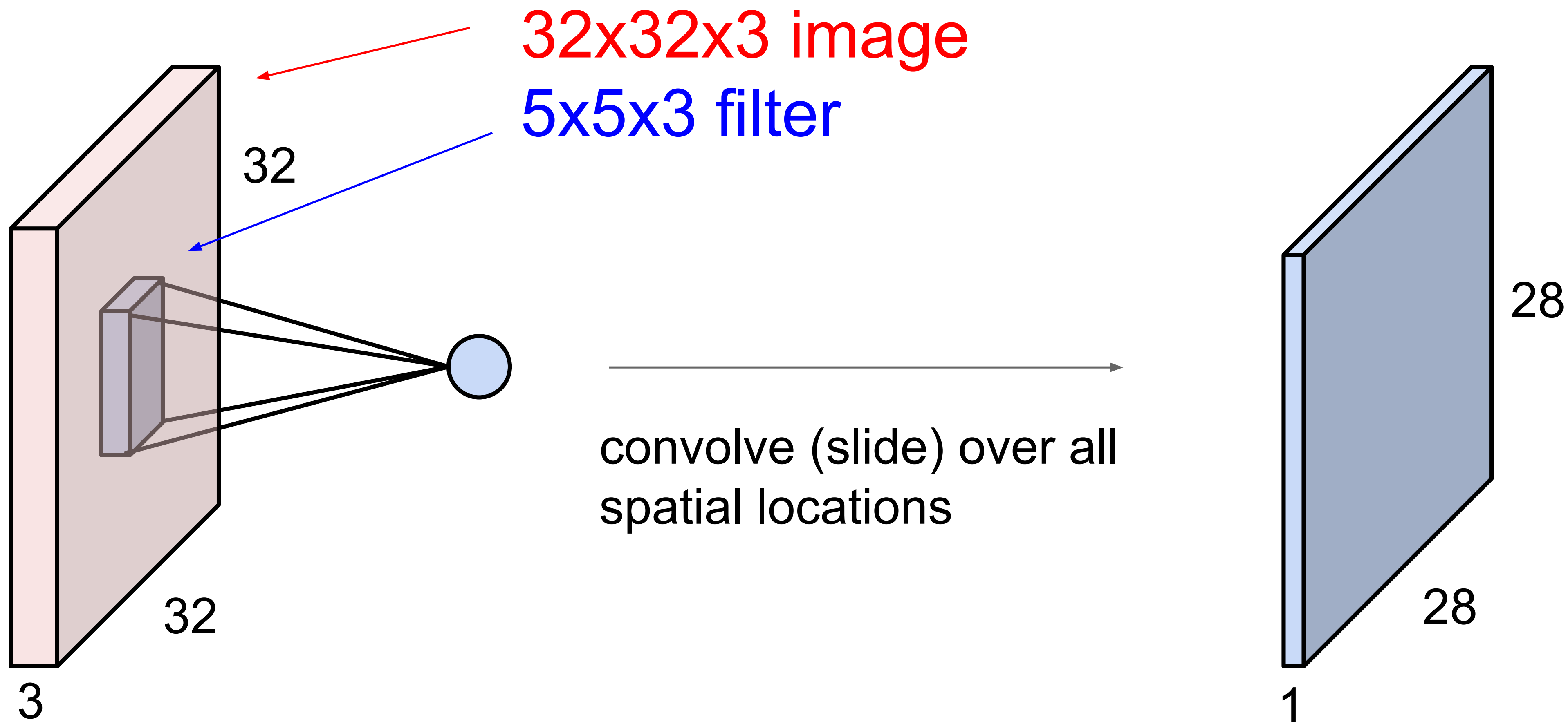
32x32x3 image



5x5x3 filter



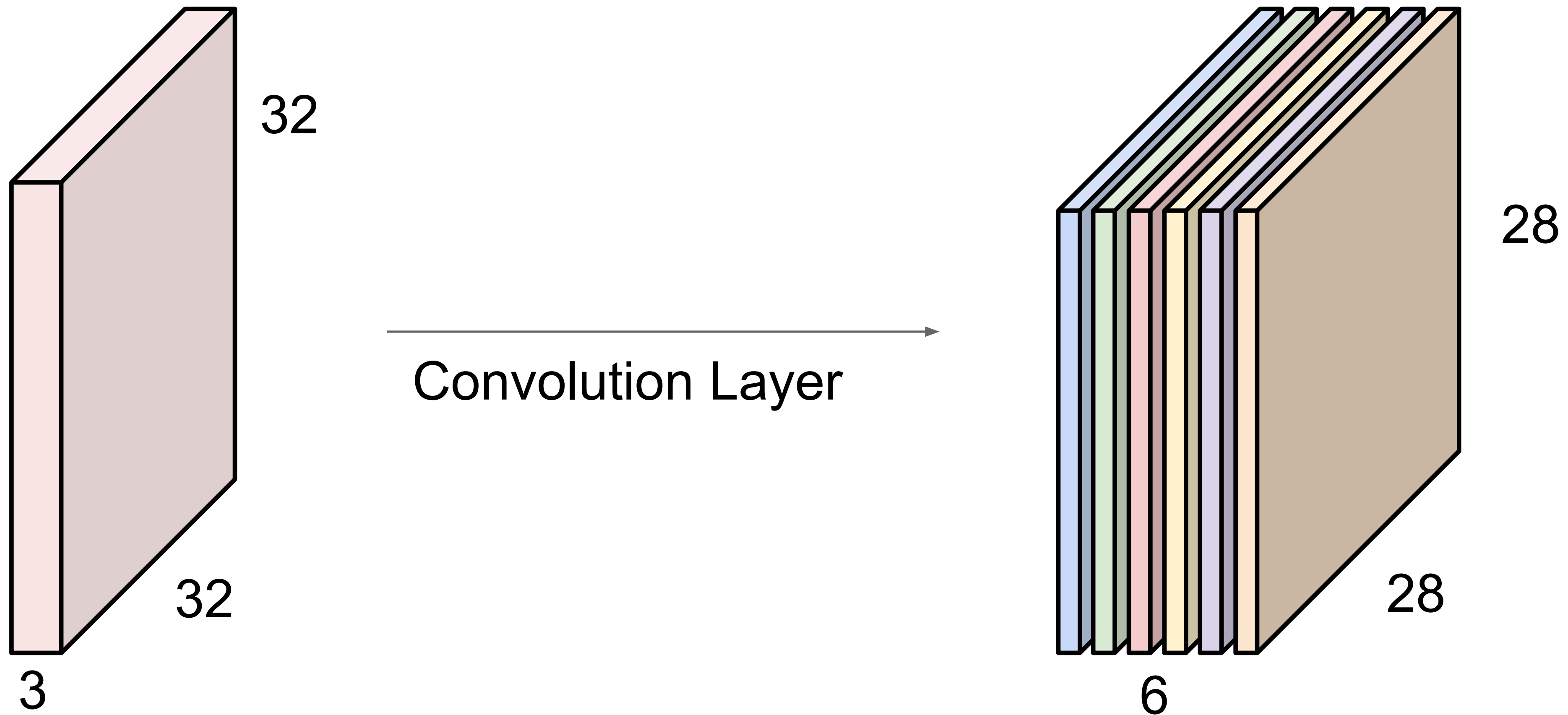
Convolution Layer



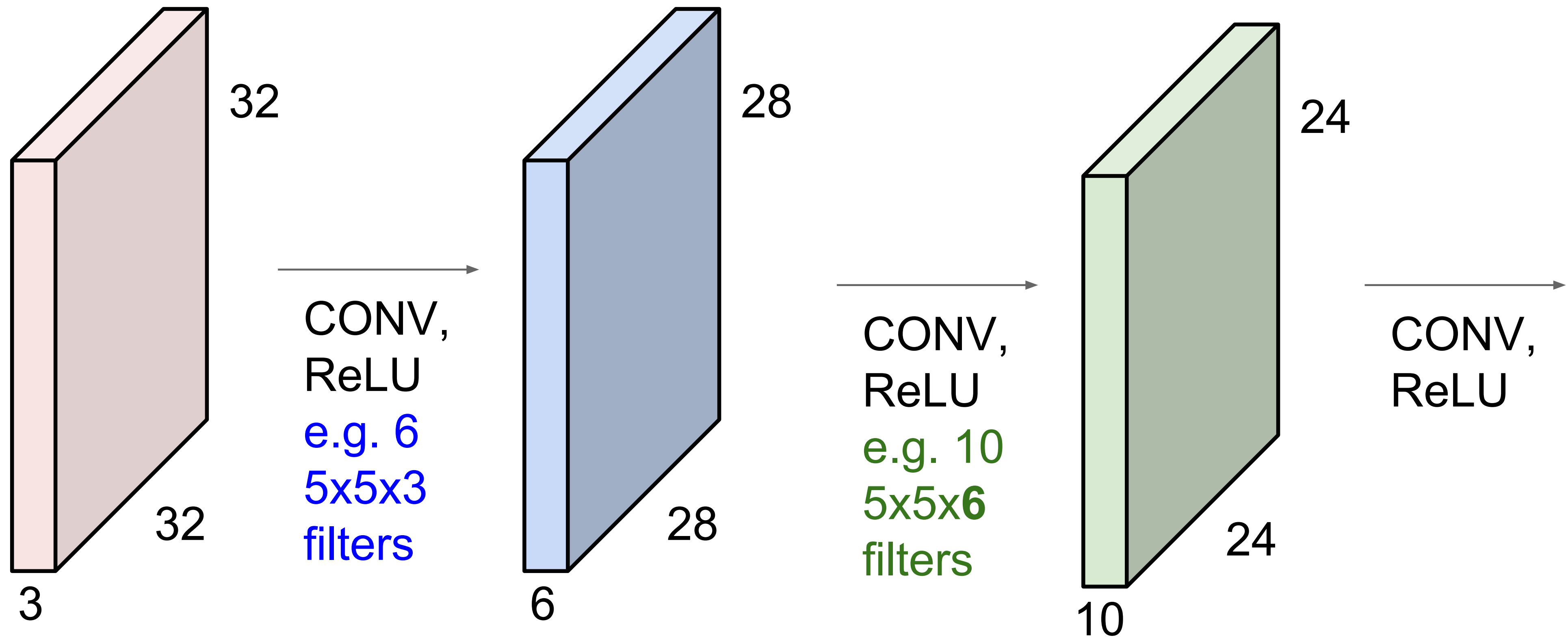
Convolution Layer



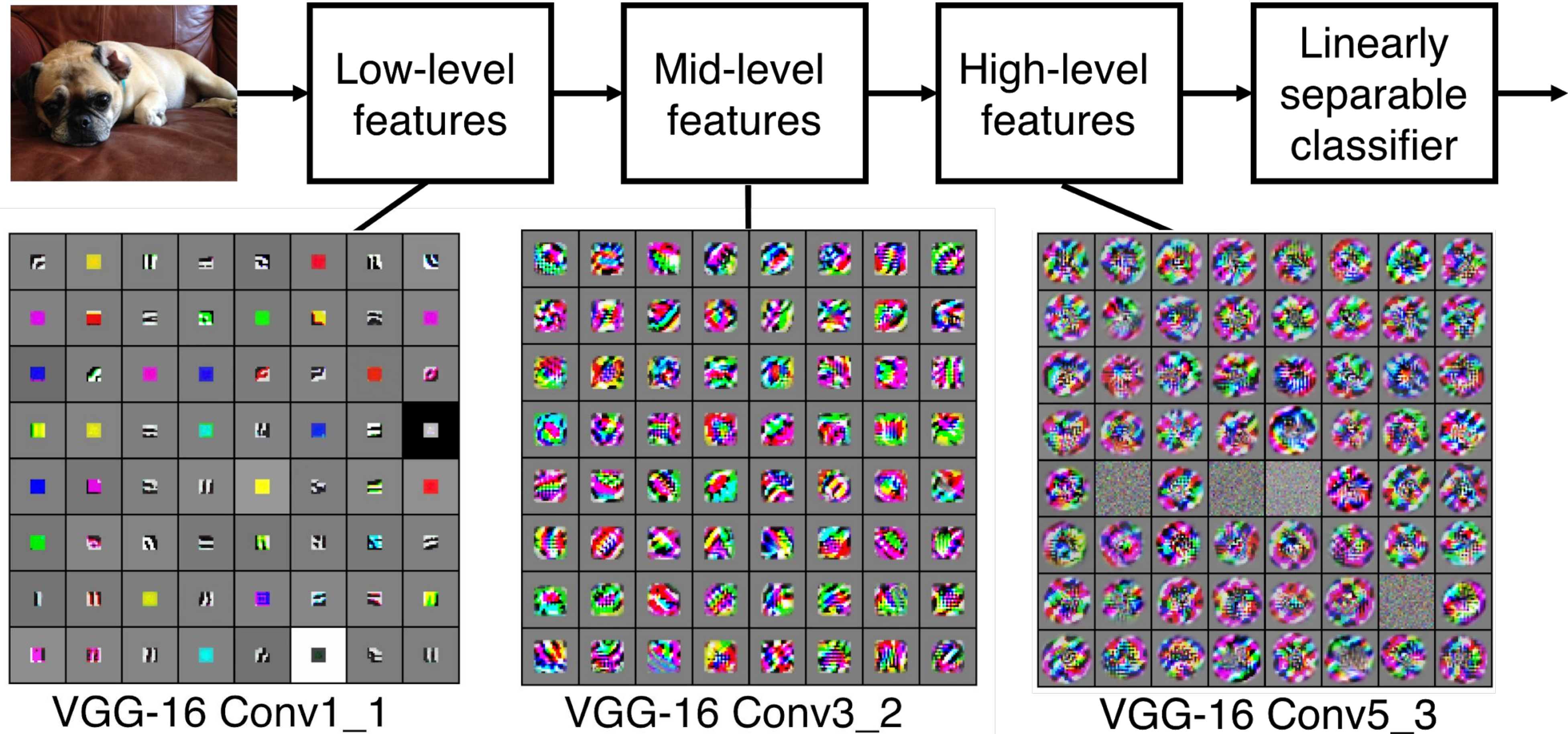
Convolution Layer



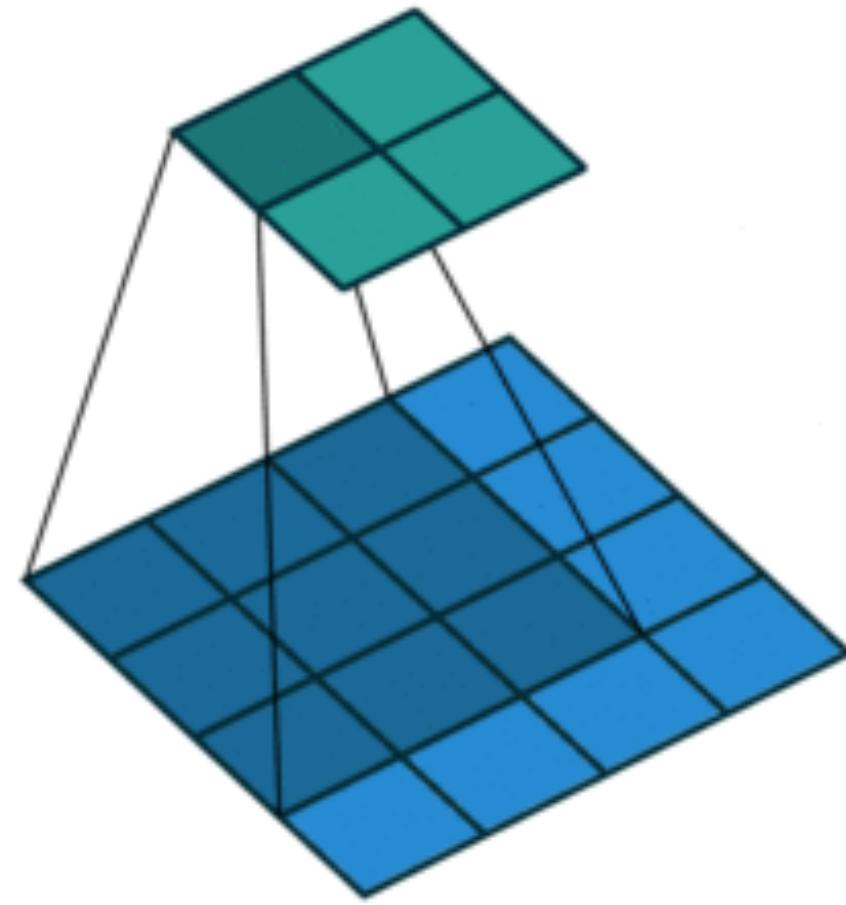
Convolutional Neural Network



What do these layers learn?

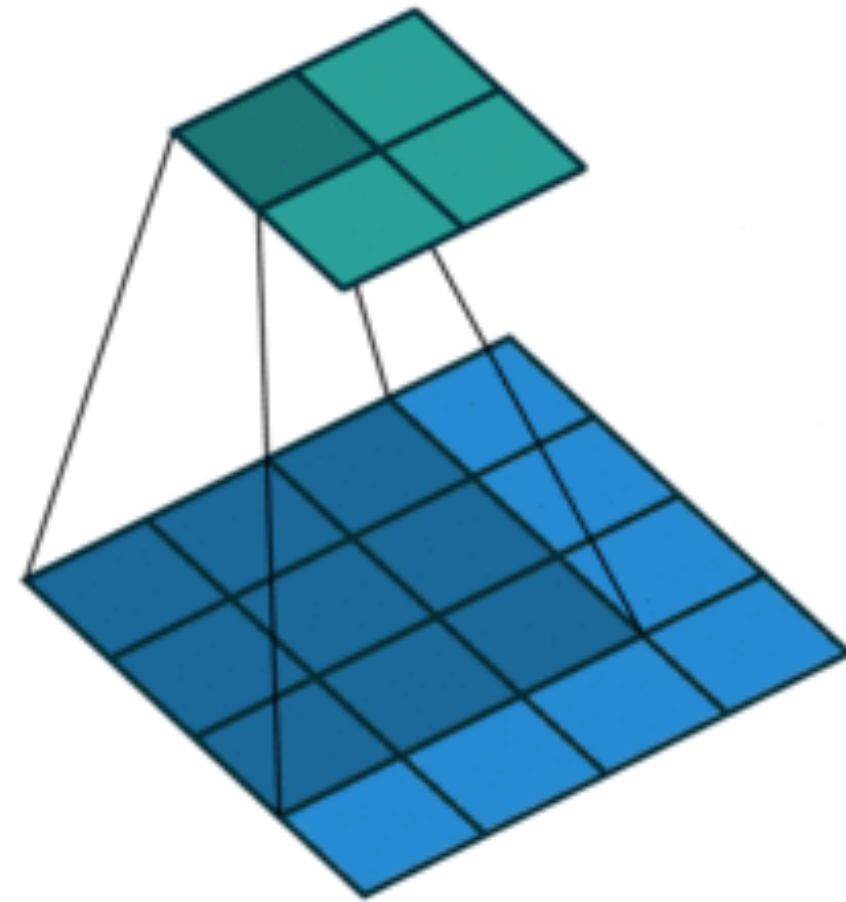


Convolutional Neural Networks (CNNs)

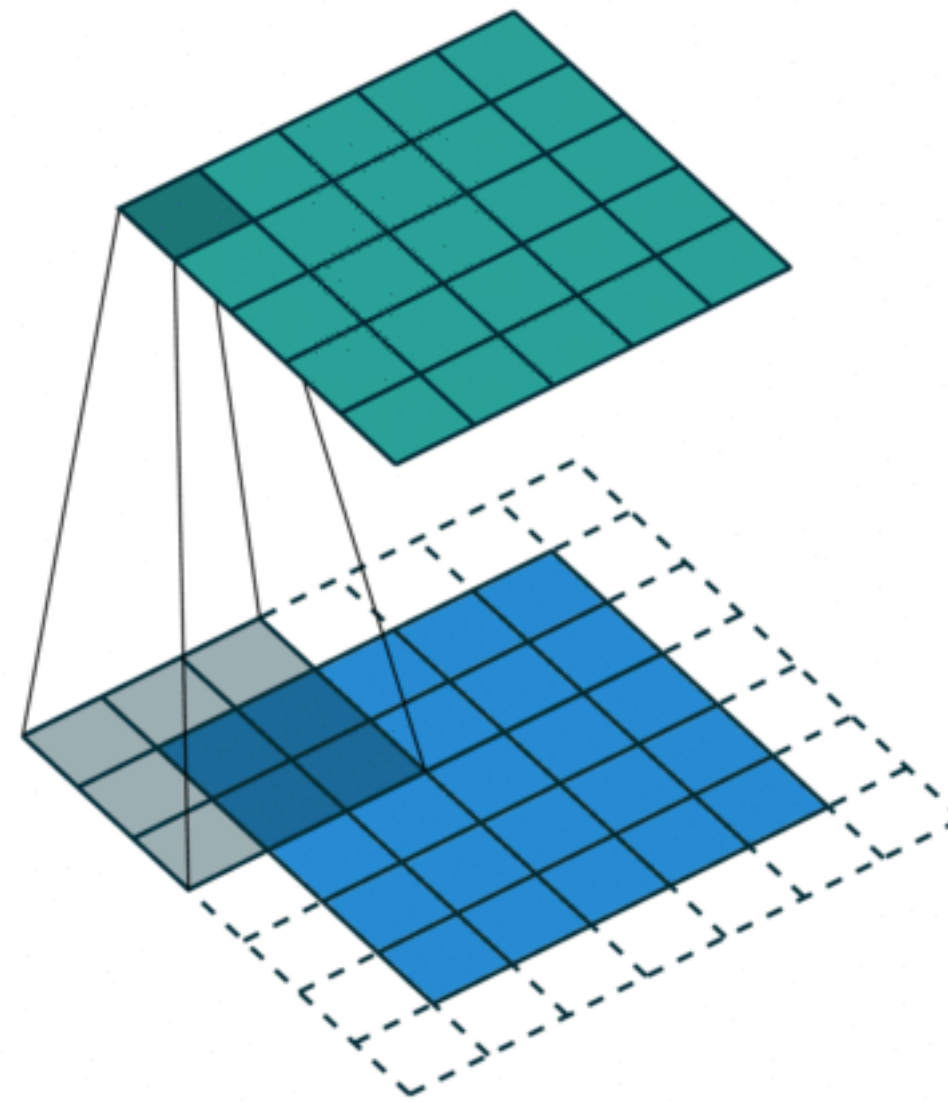


Stride=1, No padding

Convolutional Neural Networks (CNNs)

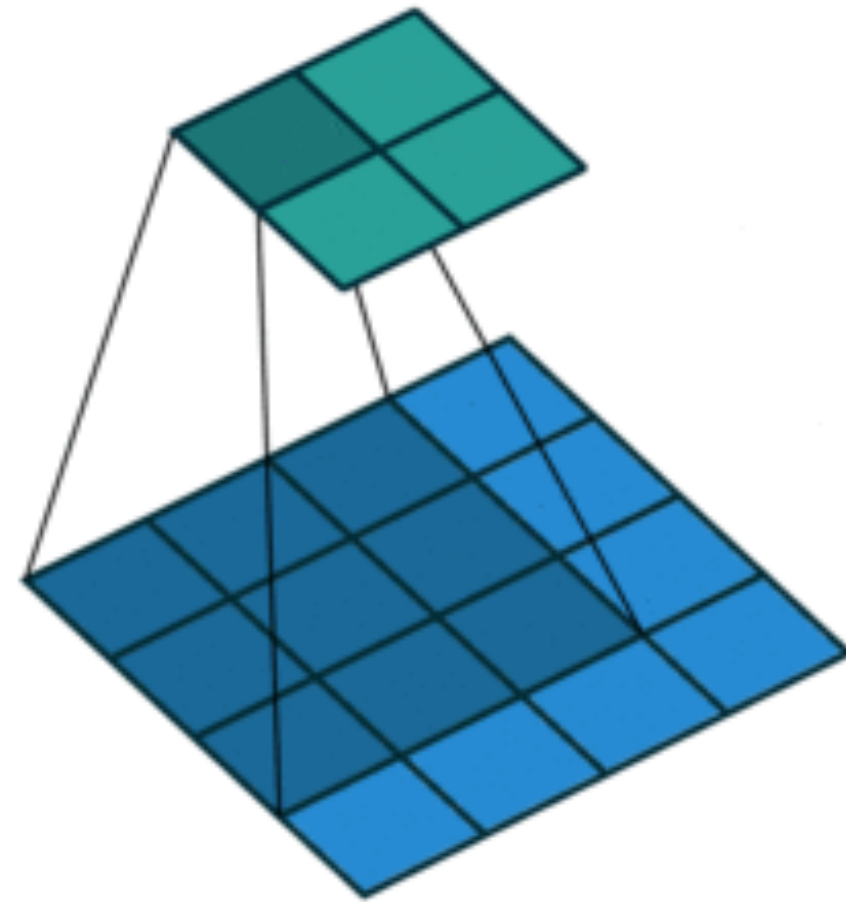


Stride=1, No padding

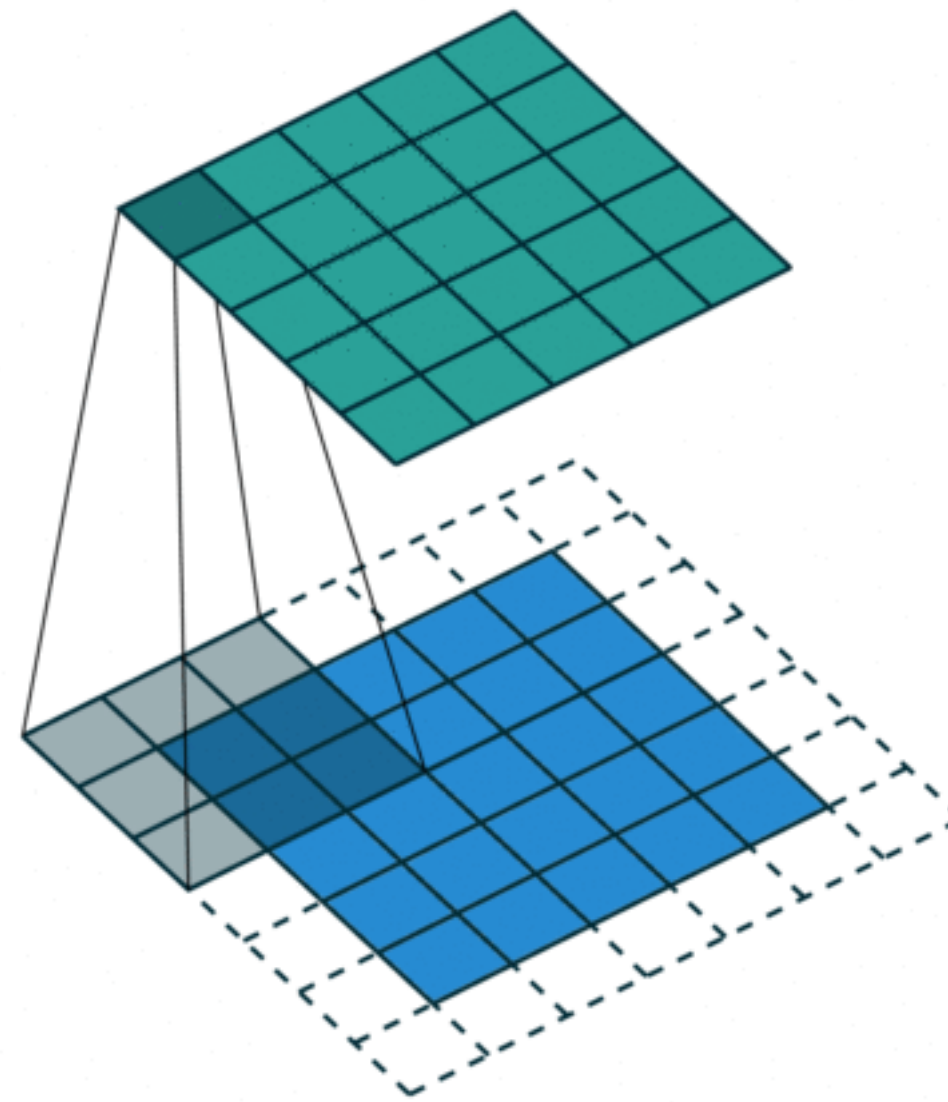


Stride=1, Padding, P=1

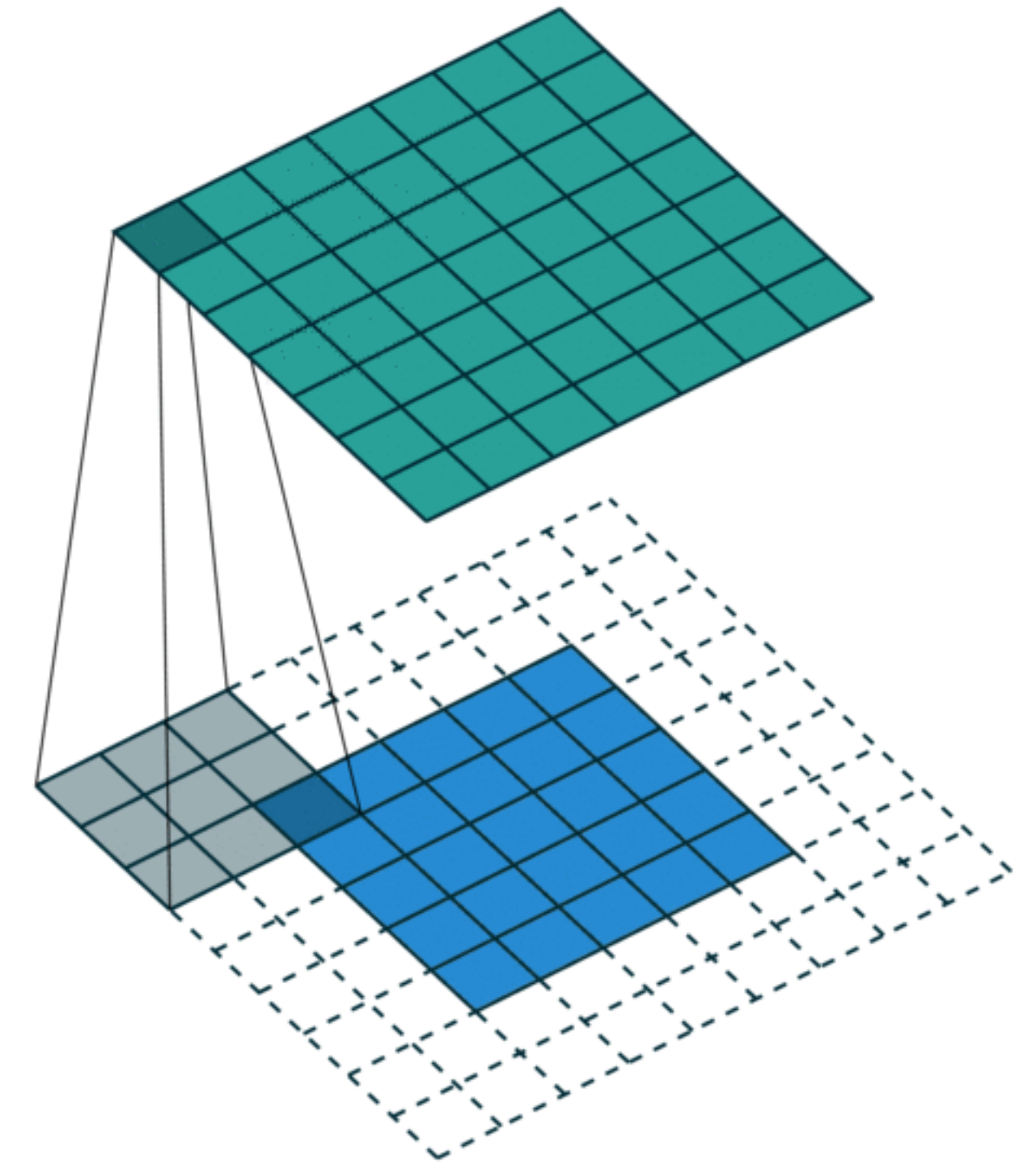
Convolutional Neural Networks (CNNs)



Stride=1, No padding

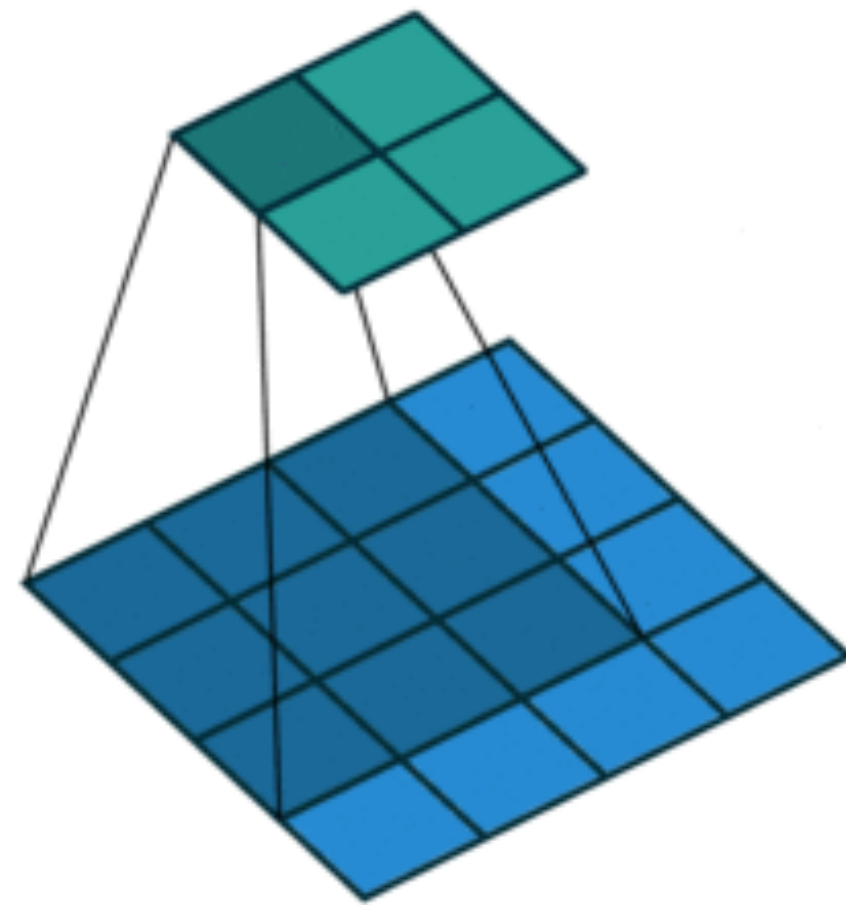


Stride=1, Padding, P=1

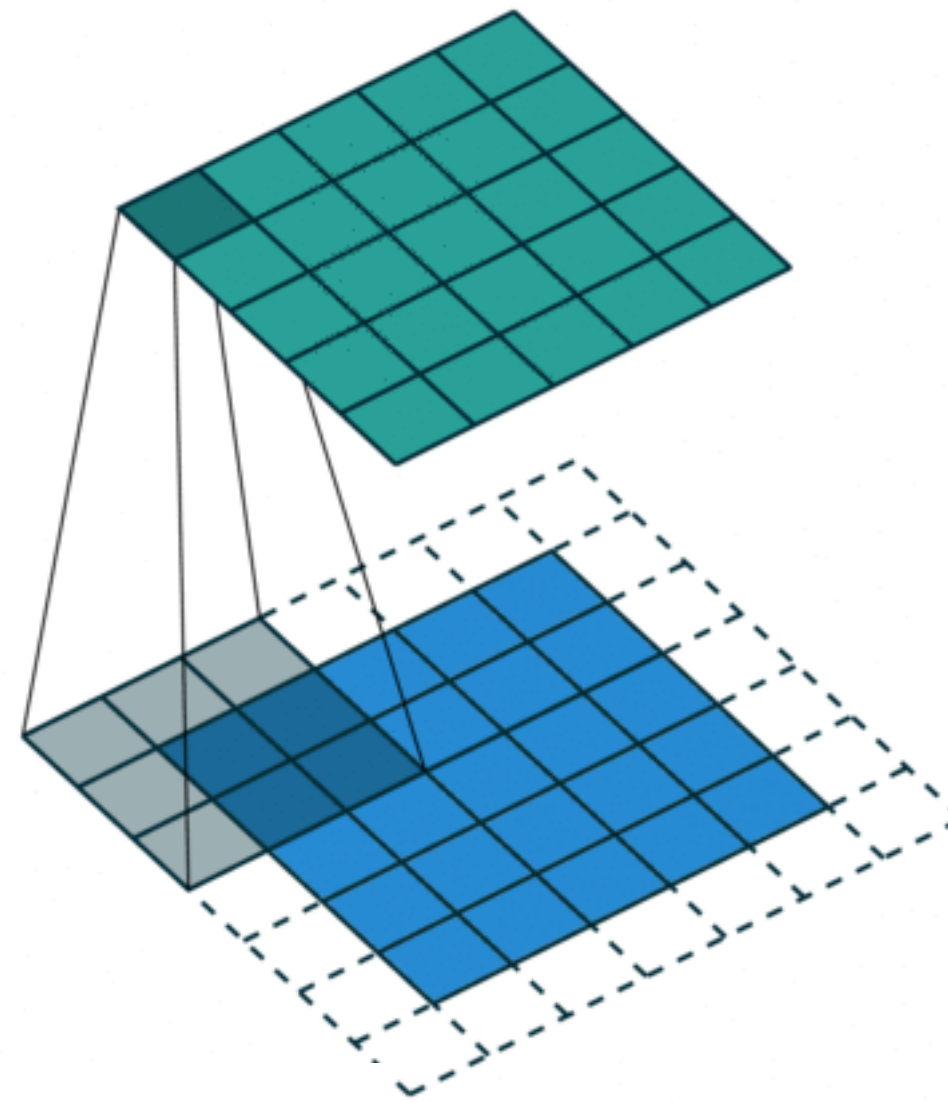


Stride=1, Padding, P=2

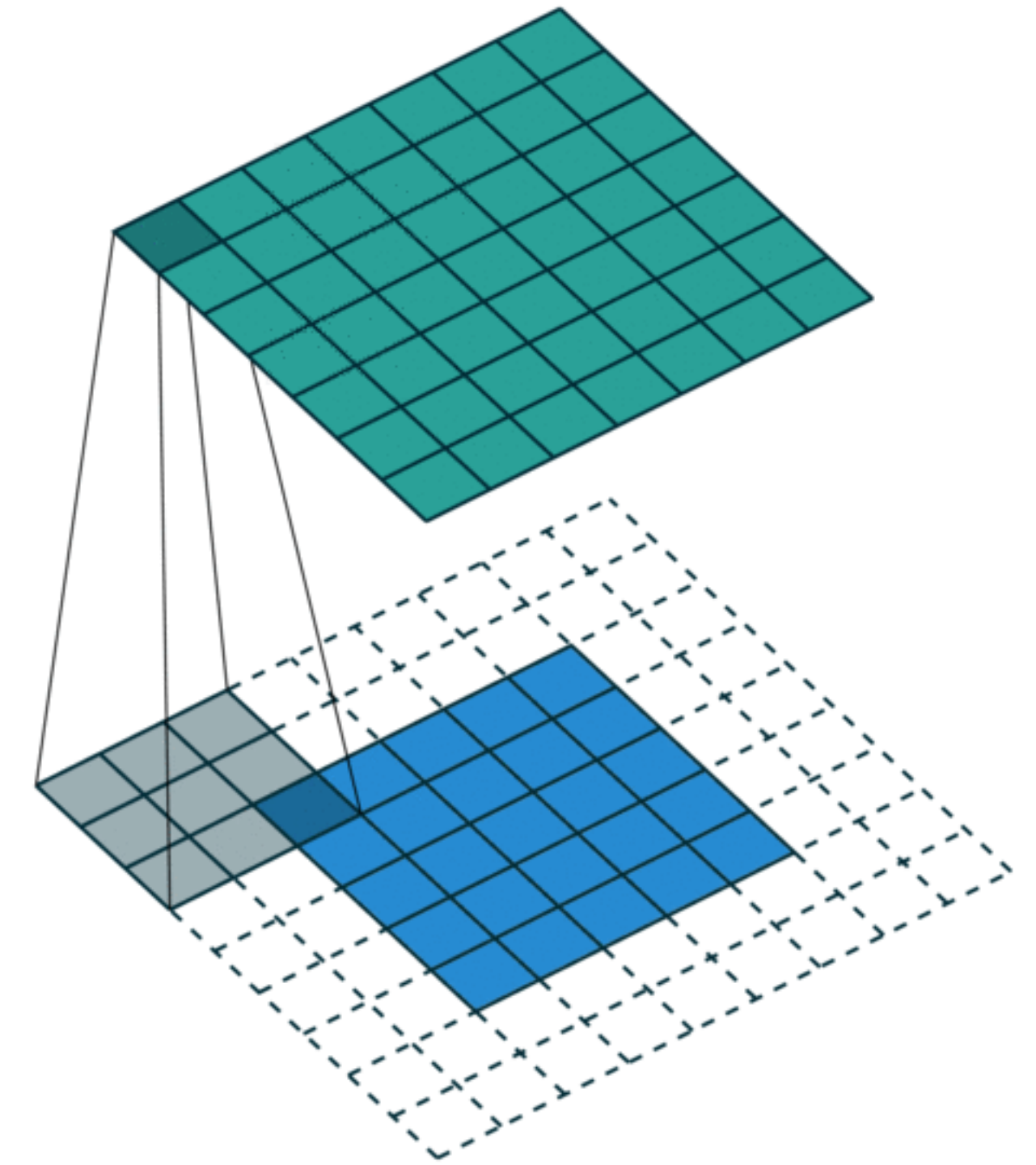
Convolutional Neural Networks (CNNs)



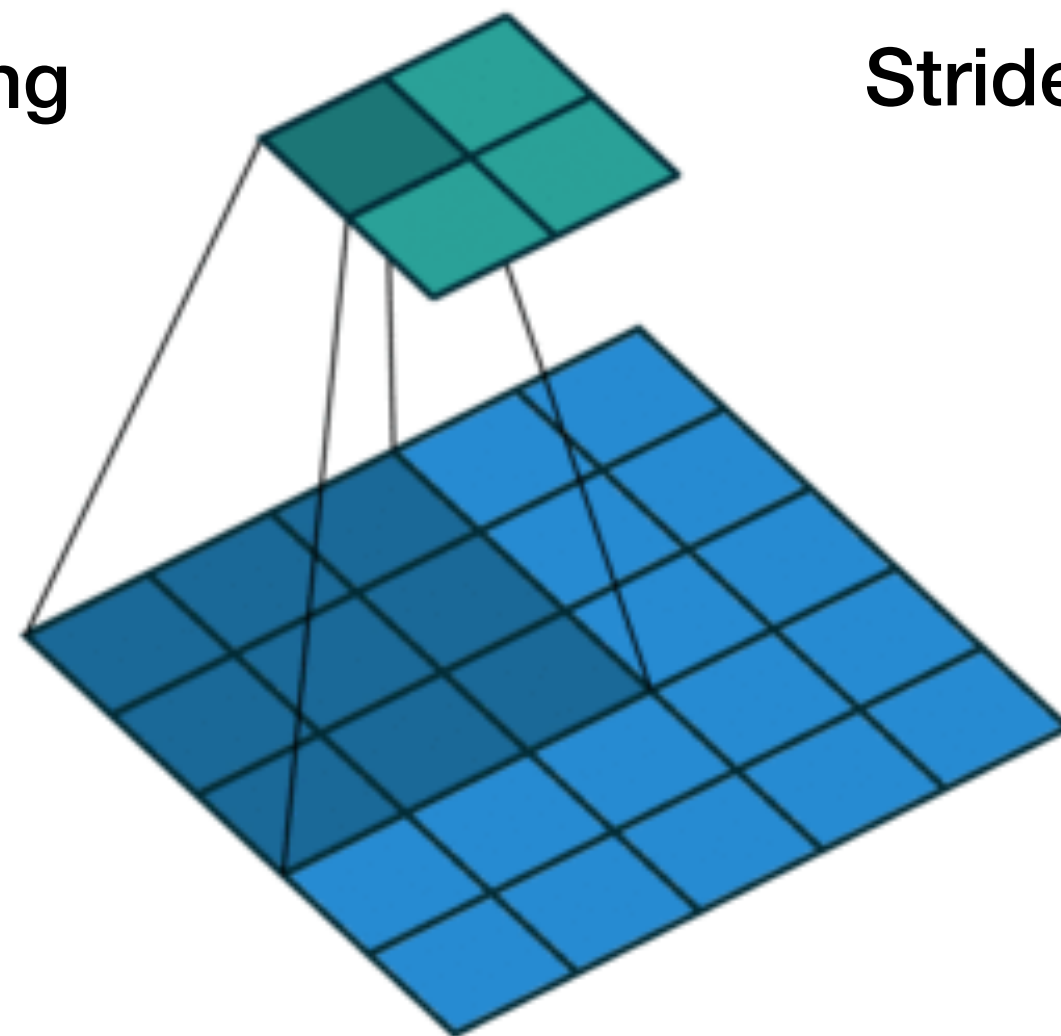
Stride=1, No padding



Stride=1, Padding, P=1

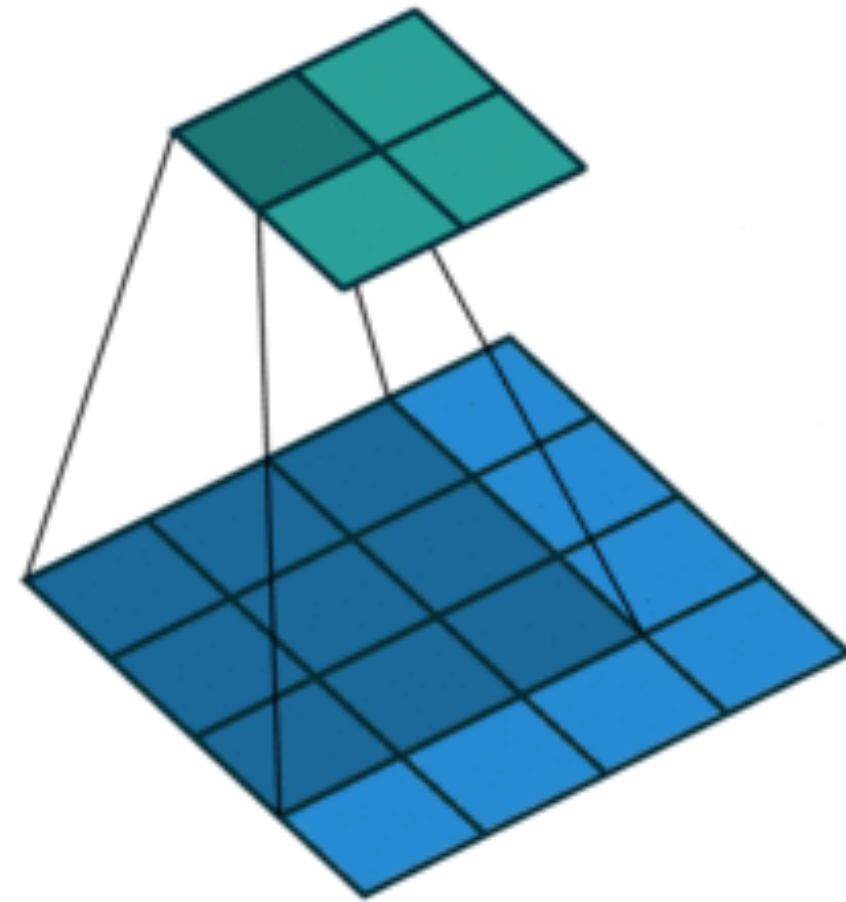


Stride=1, Padding, P=2

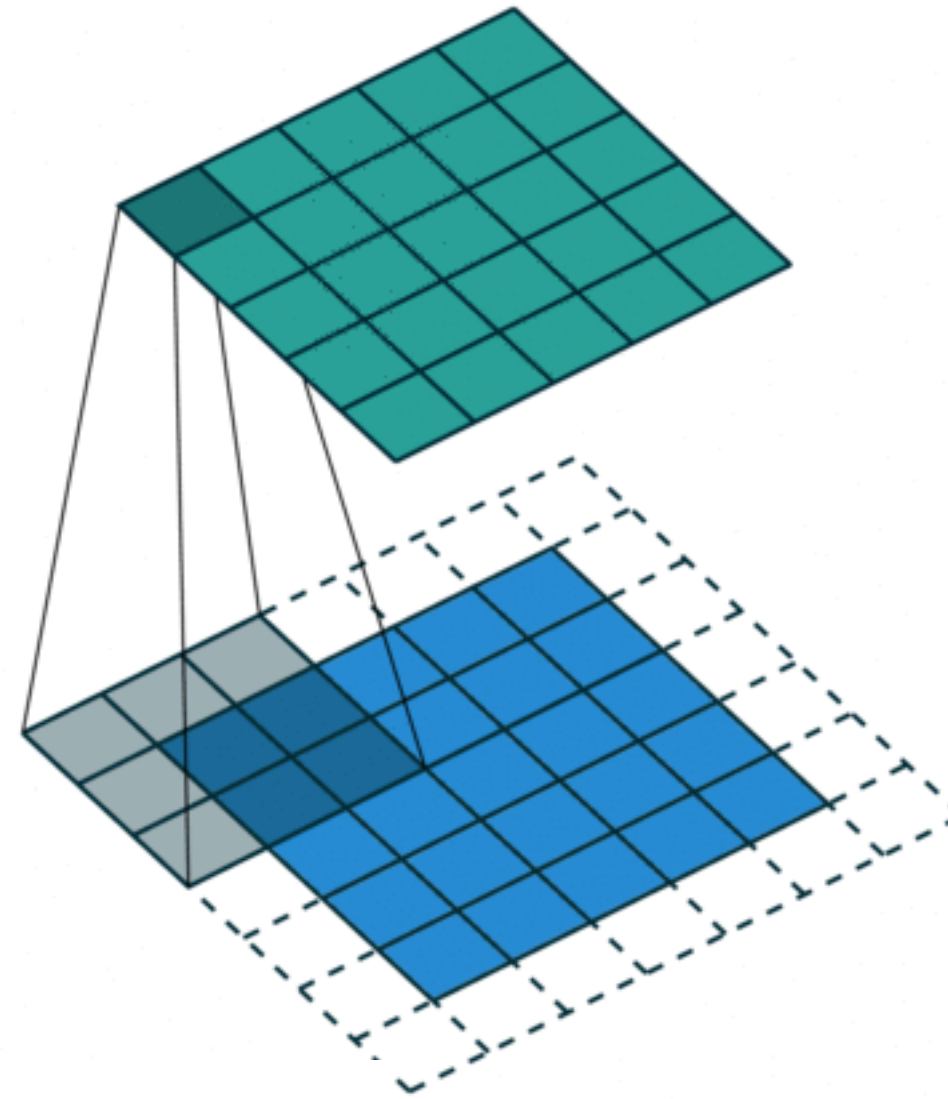


Stride=2, No padding

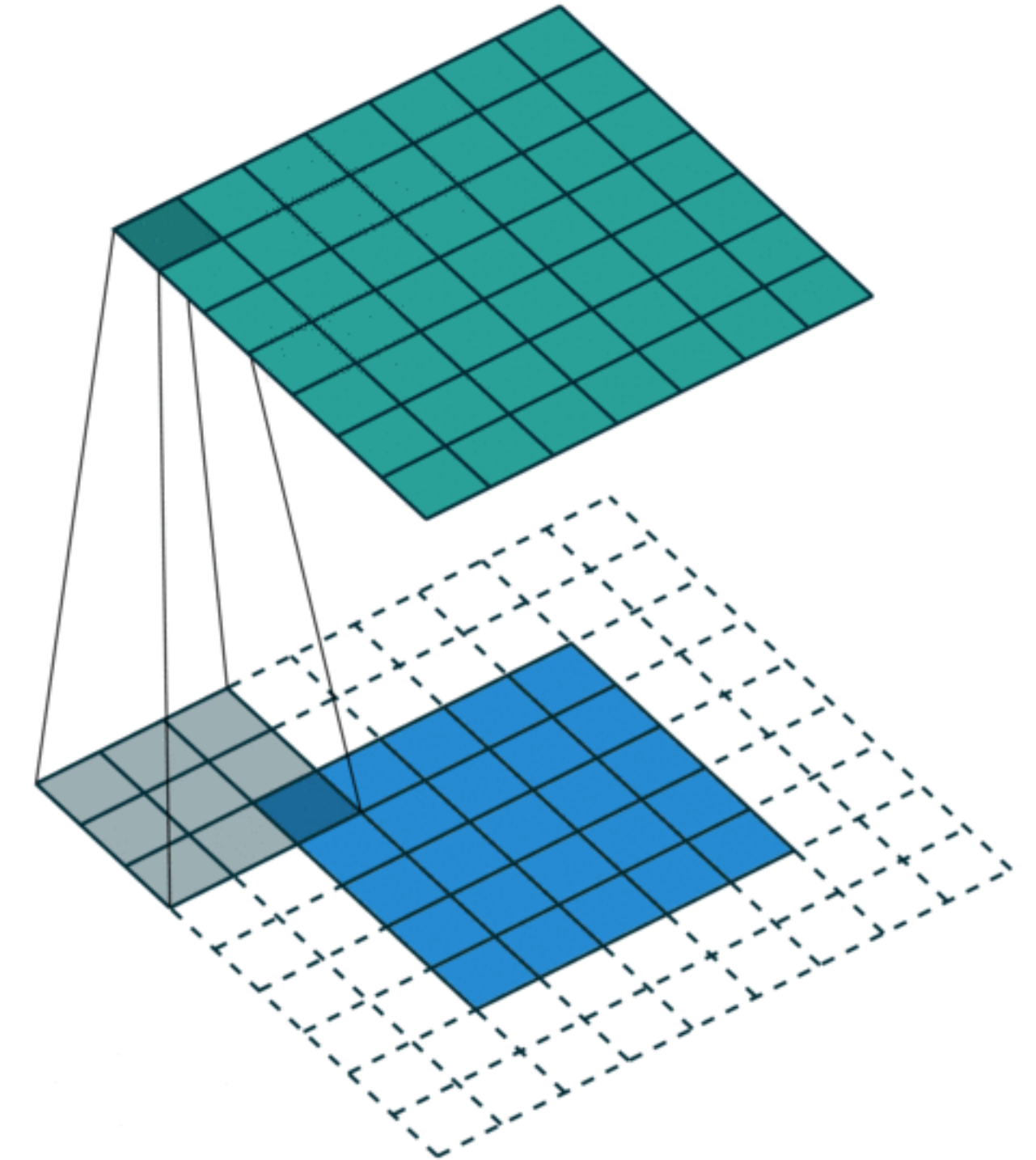
Convolutional Neural Networks (CNNs)



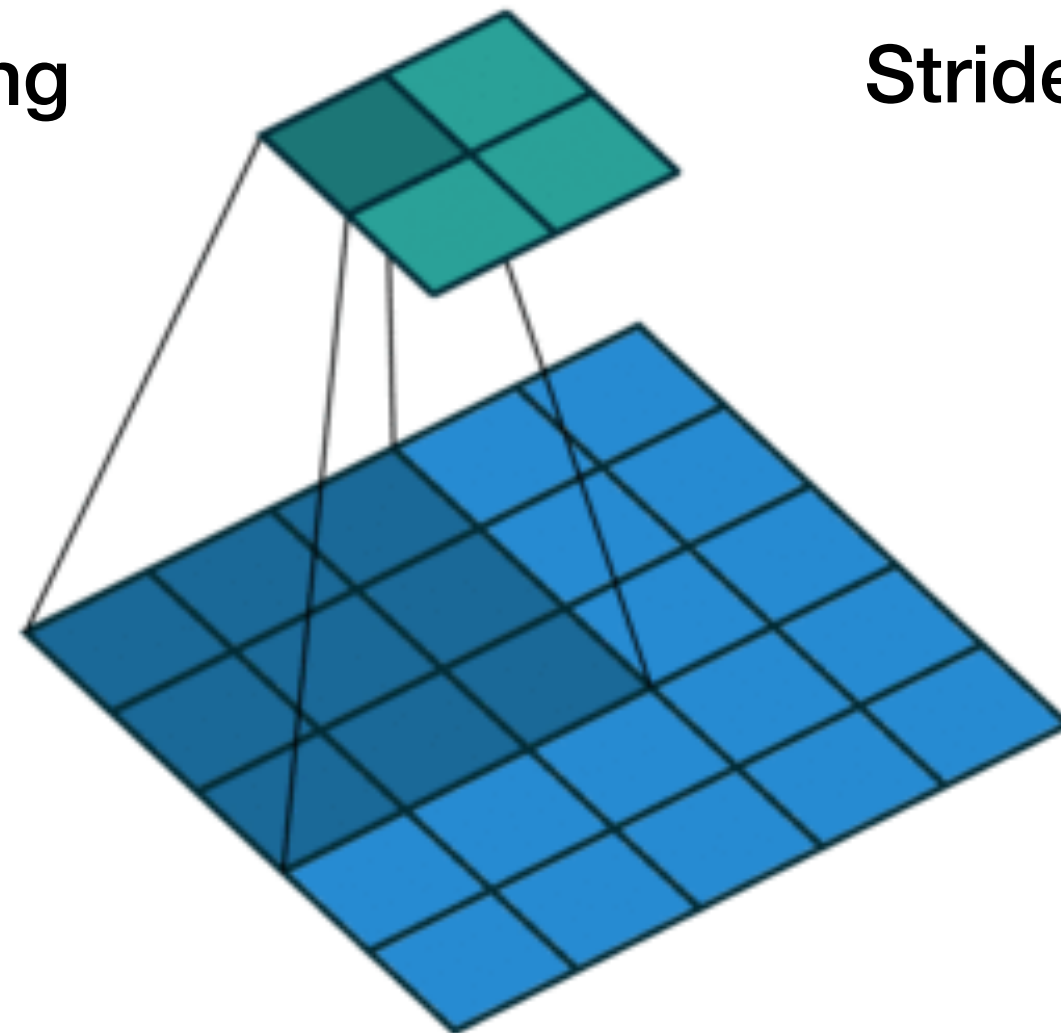
Stride=1, No padding



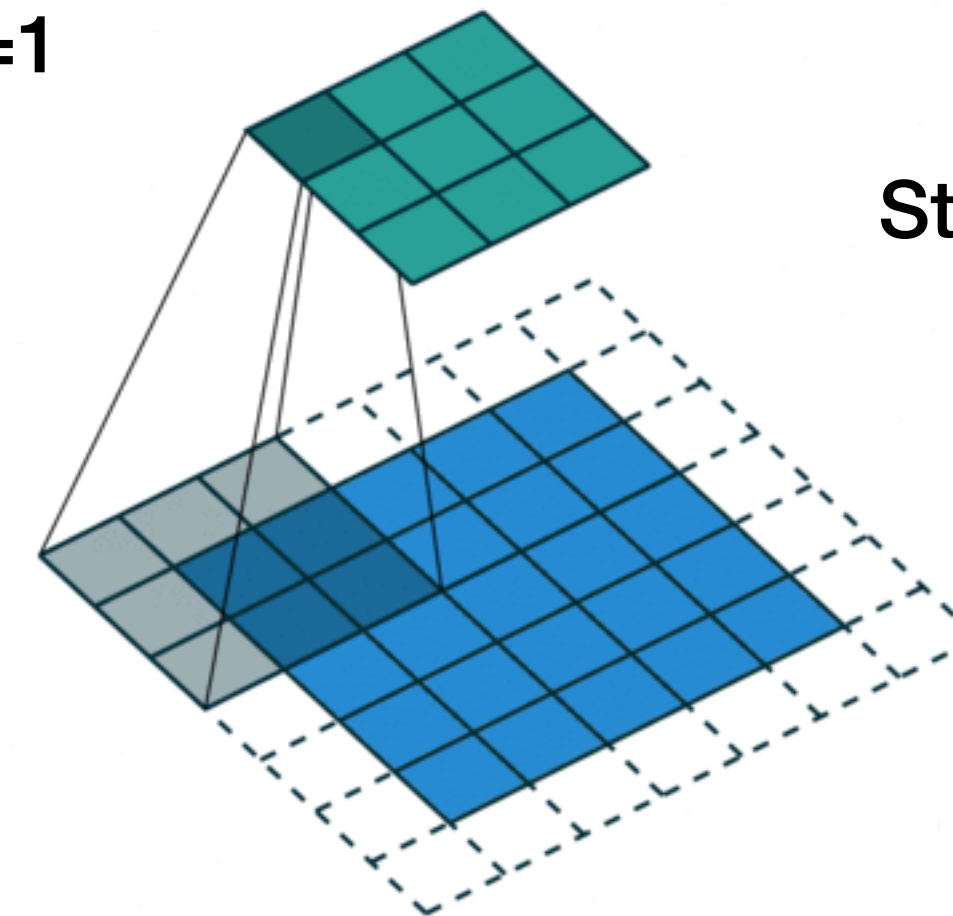
Stride=1, Padding, P=1



Stride=1, Padding, P=2



Stride=2, No padding

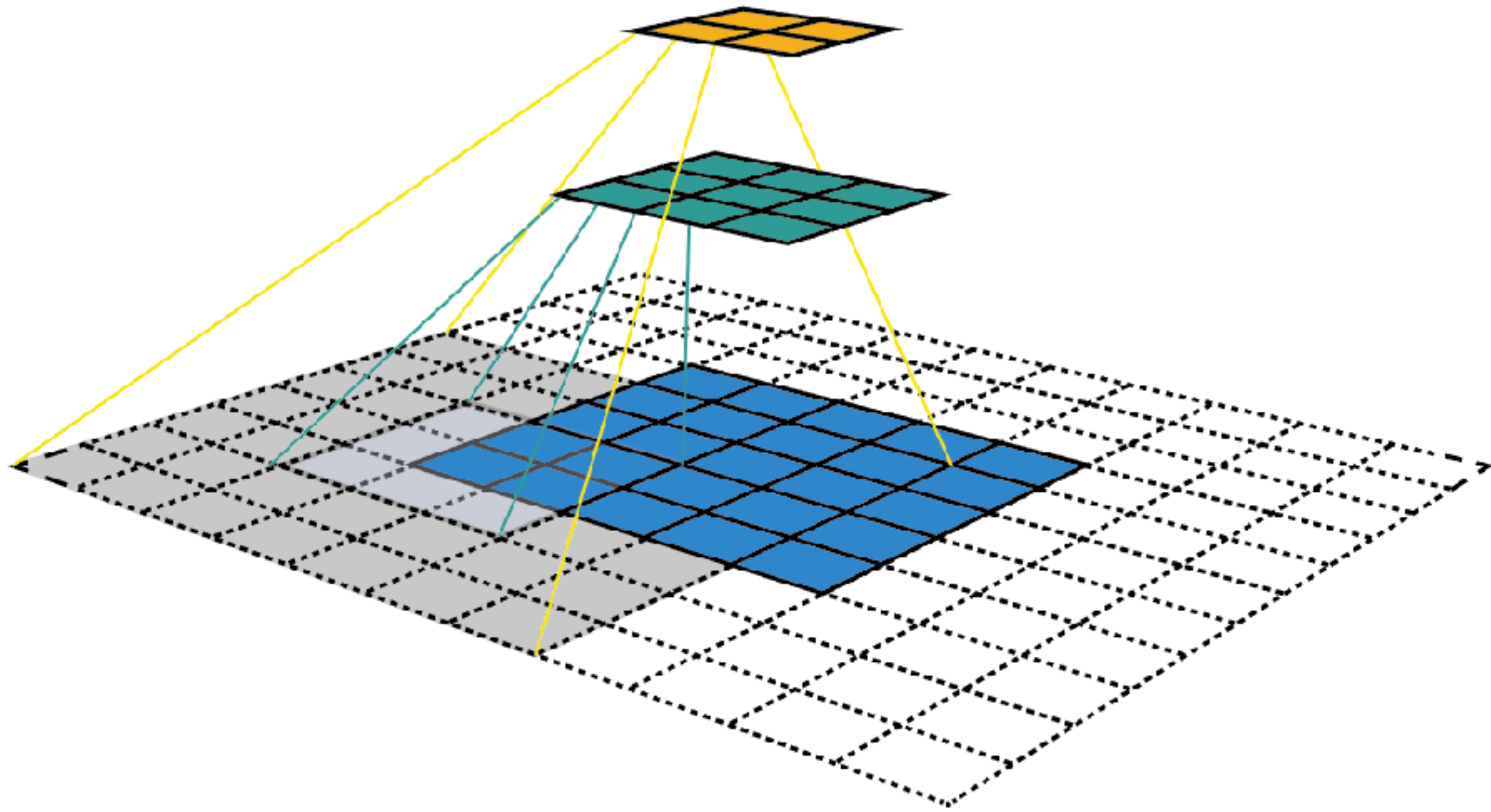
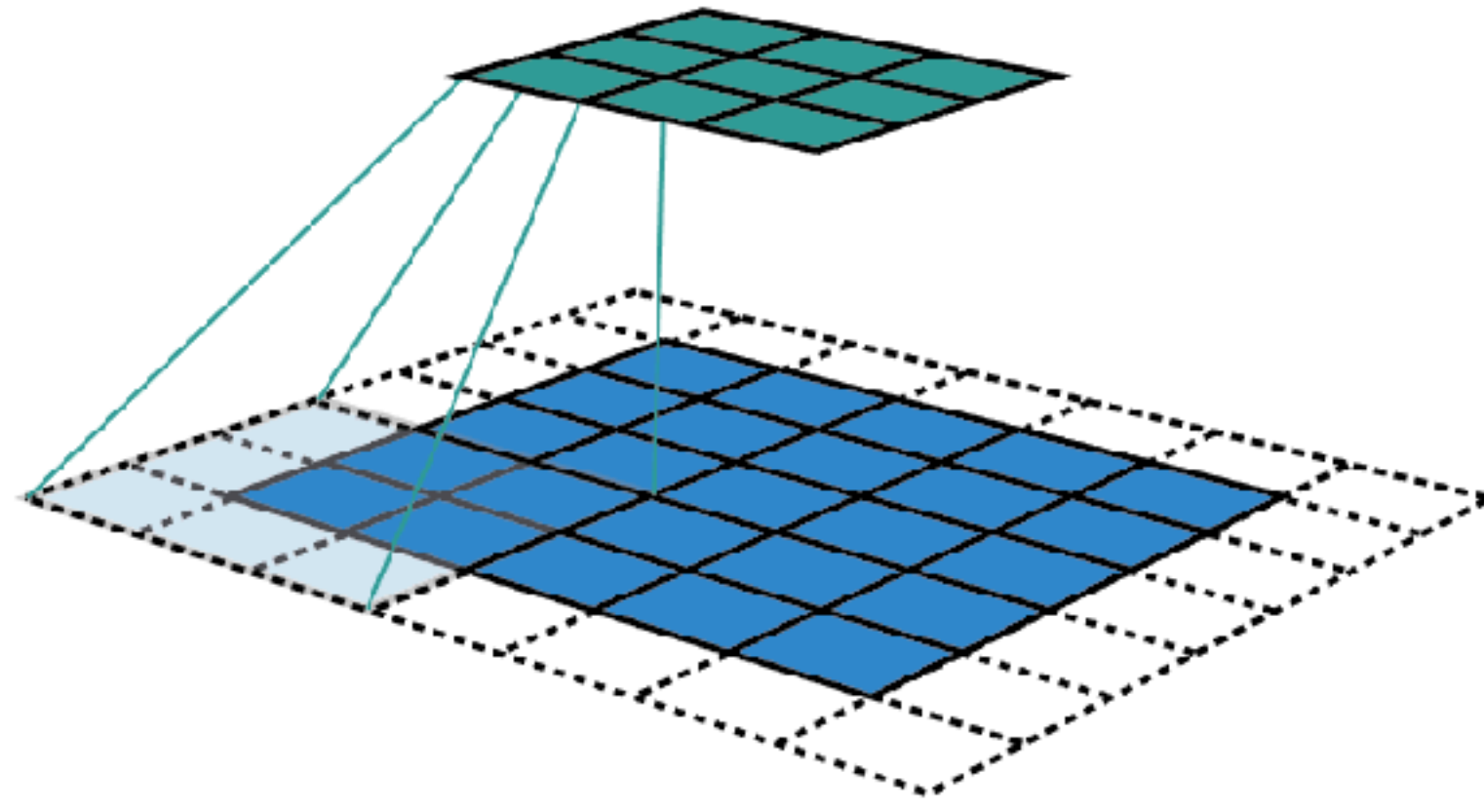


Stride=2, Padding, P=1

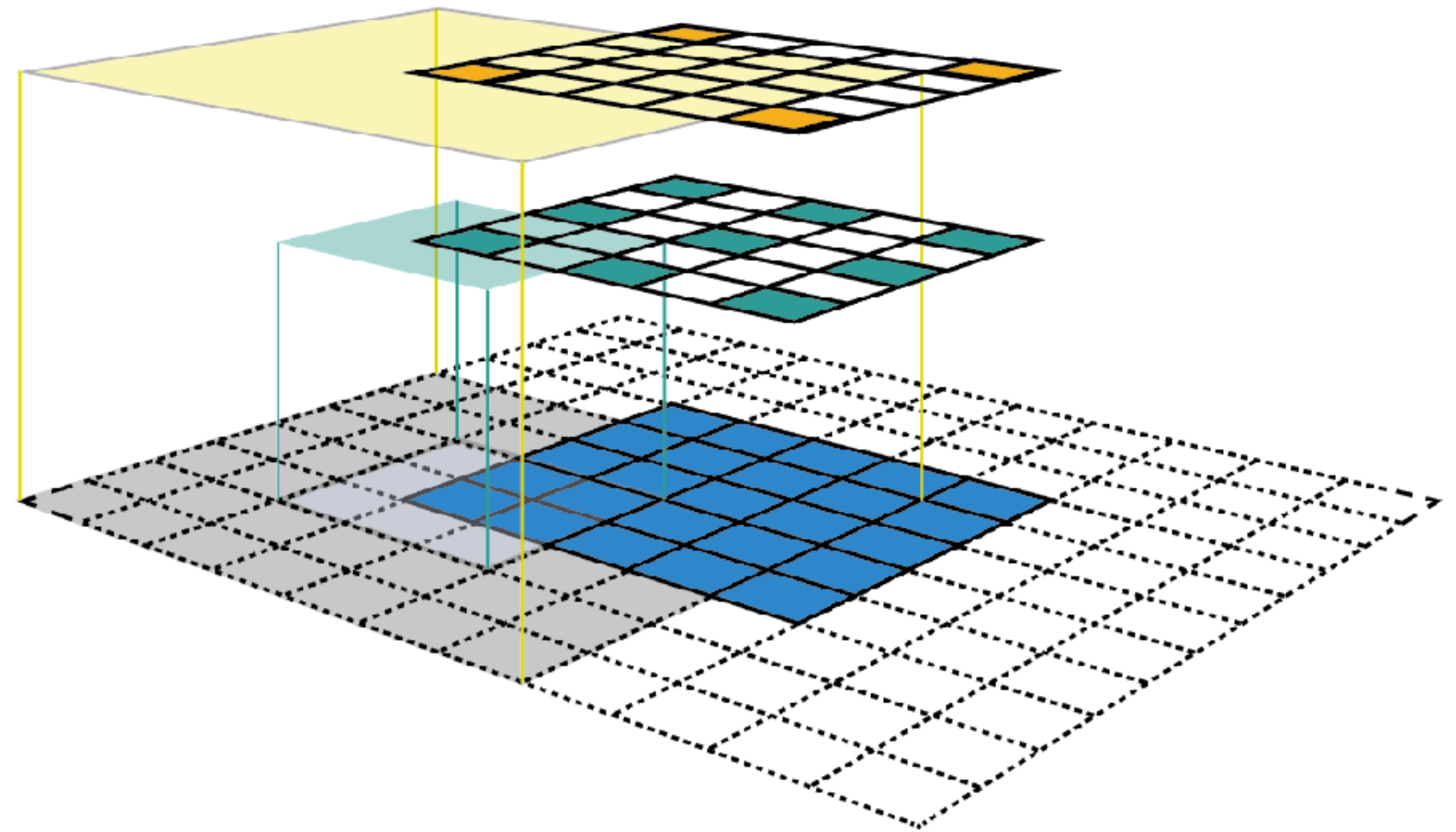
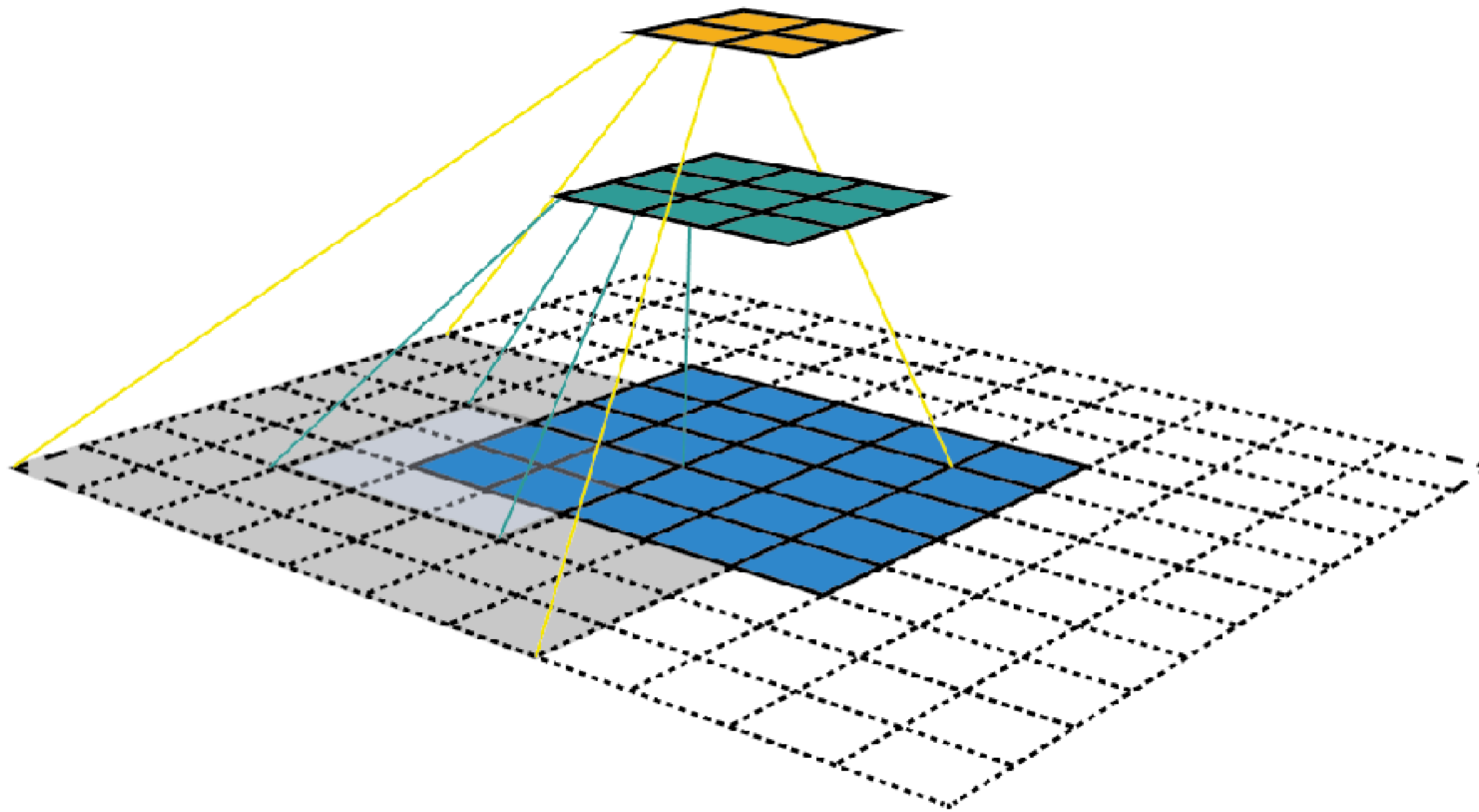
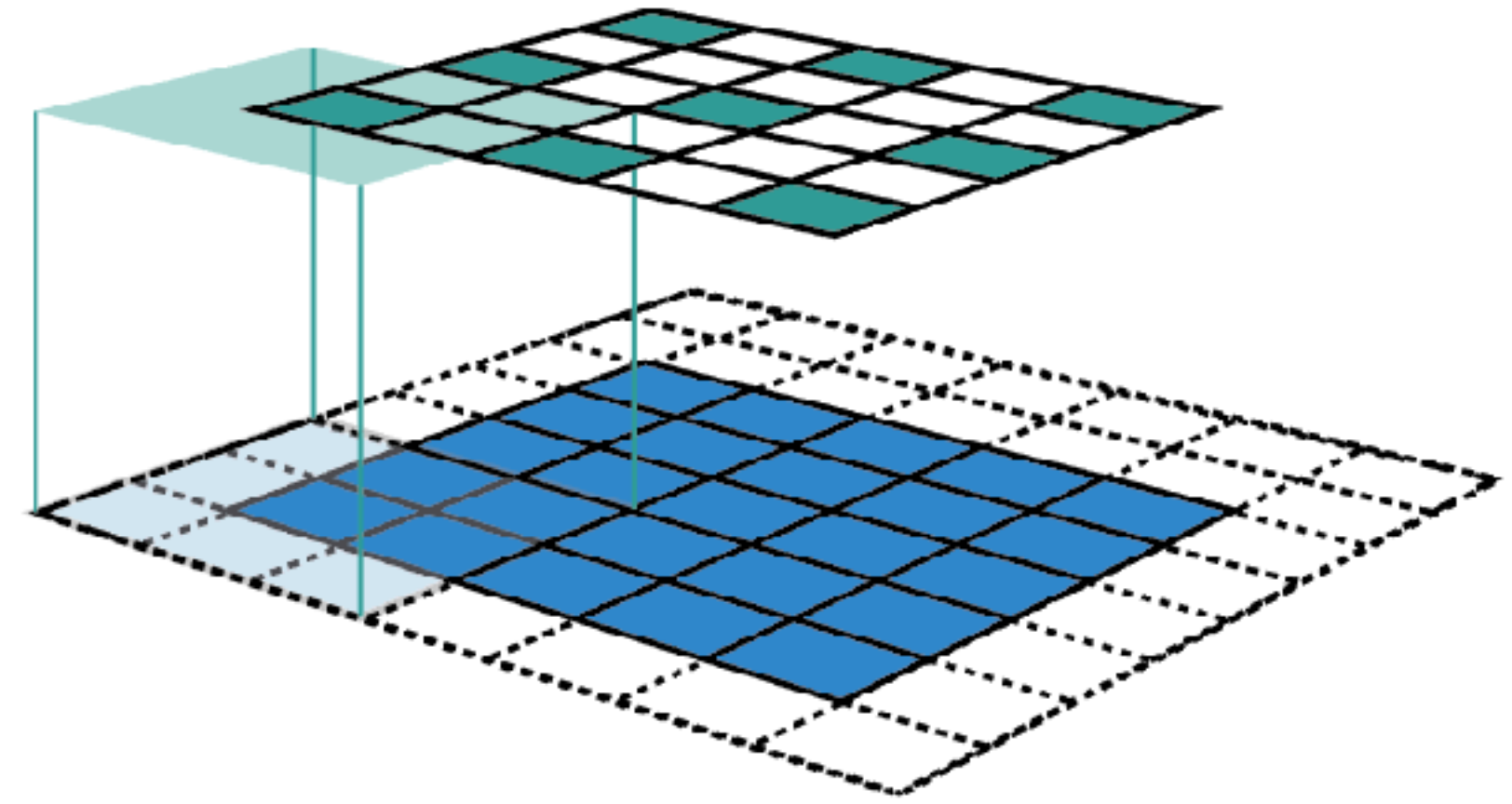
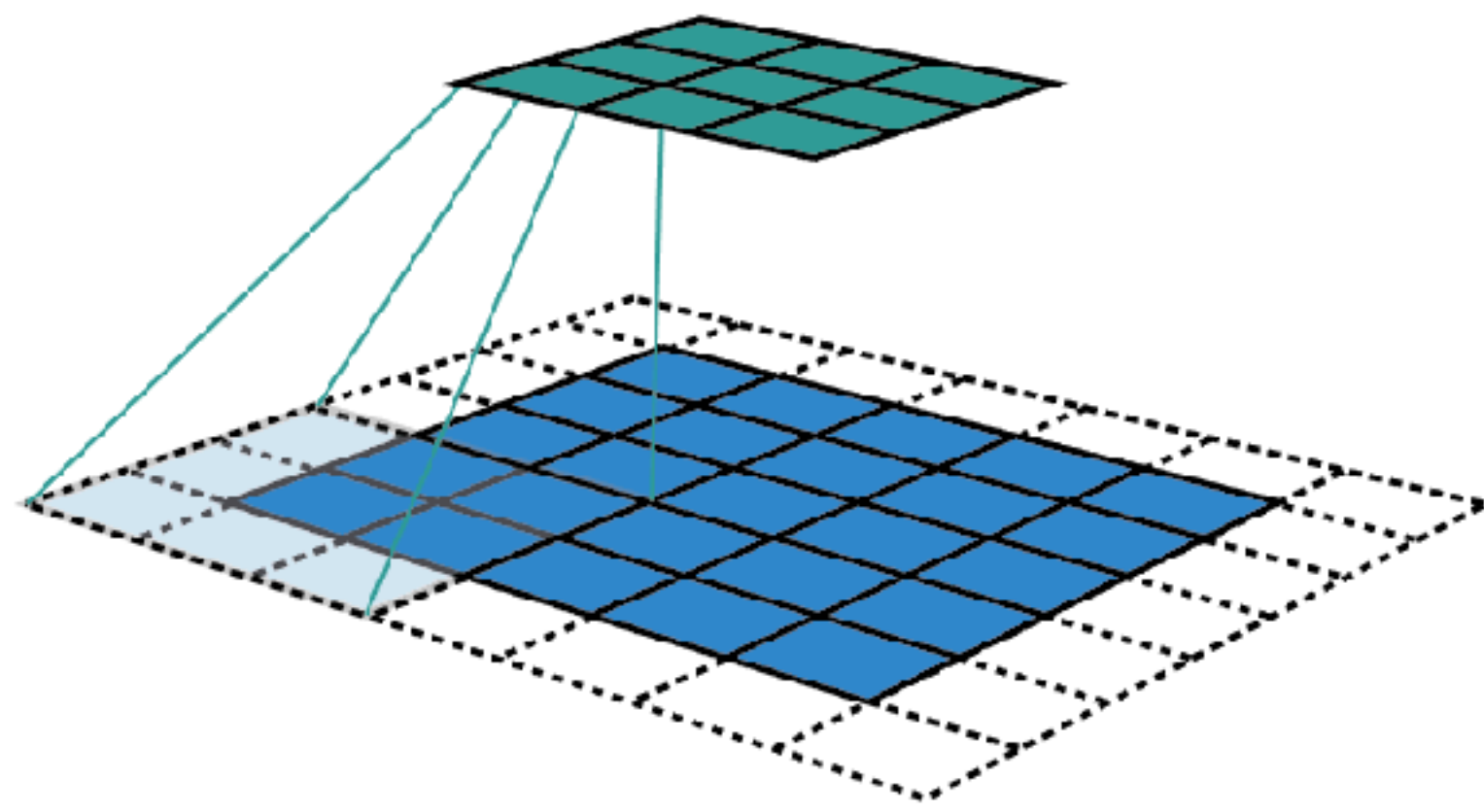
Convolution Layers: Summary

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K ,
 - their spatial extent F ,
 - the stride S ,
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry)
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.

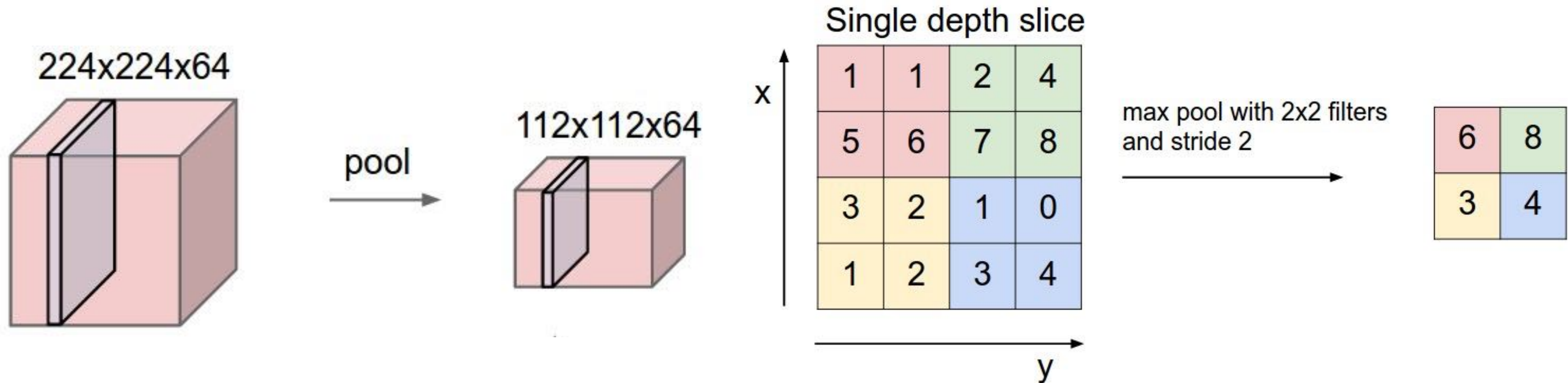
Receptive Field



Receptive Field



Pooling Layer



- Why pooling?
Reduce the size of the representation, speed up the computations and make the features a little more robust.
- Max pooling is popularly used in CNNs.

Pooling Layer

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires two hyperparameters:
 - their spatial extent F ,
 - the stride S ,
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F)/S + 1$
 - $H_2 = (H_1 - F)/S + 1$
 - $D_2 = D_1$

Batch Normalization

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots x_m\}$;
Parameters to be learned: γ, β

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

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

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

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

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

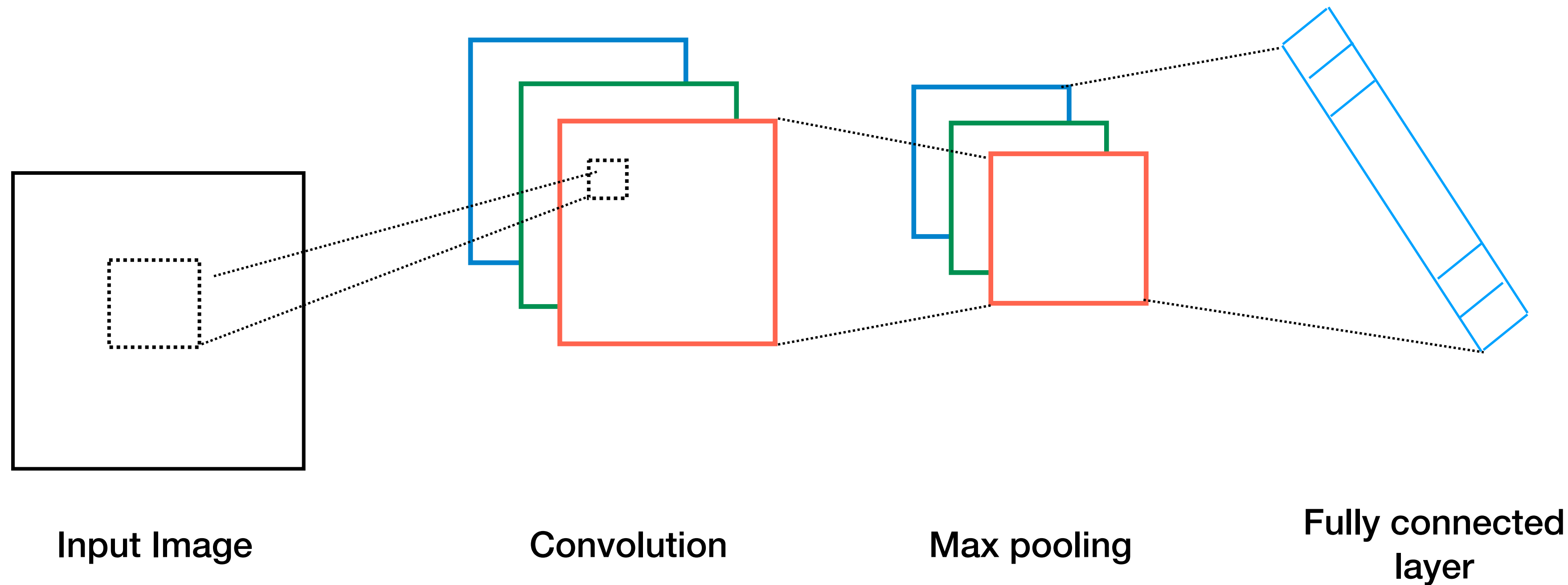
- γ and β are learned parameters used across all batches
- At test time: Individual inputs, no mini-batch.
 - First, normalize inputs using training population statistics.

$$\hat{x} \leftarrow \frac{x - \mu_{\text{pop}}}{\sqrt{\sigma_{\text{pop}}^2 + \epsilon}}$$

- Then, scale and shift.

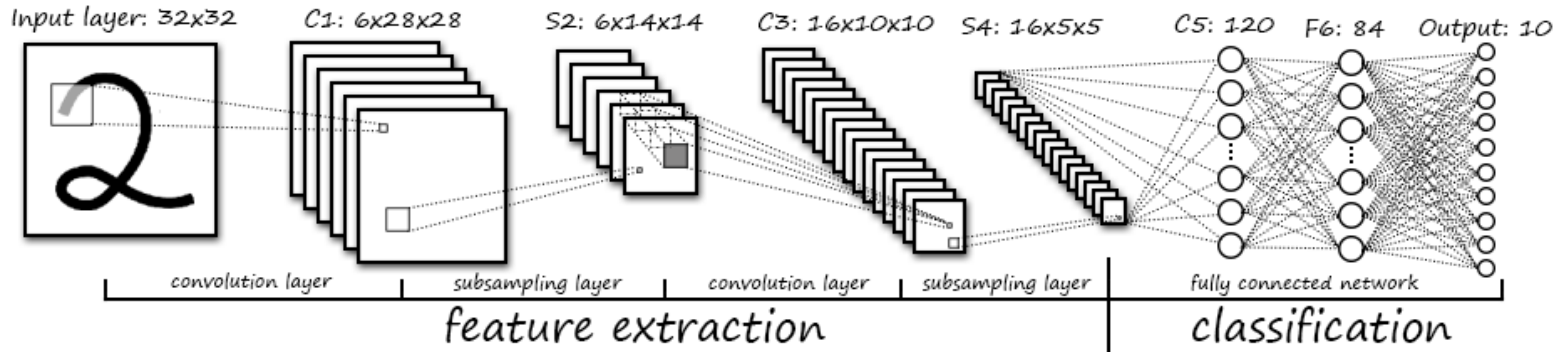
$$\hat{y} \leftarrow \gamma \hat{x} + \beta$$

Convolutional Architectures



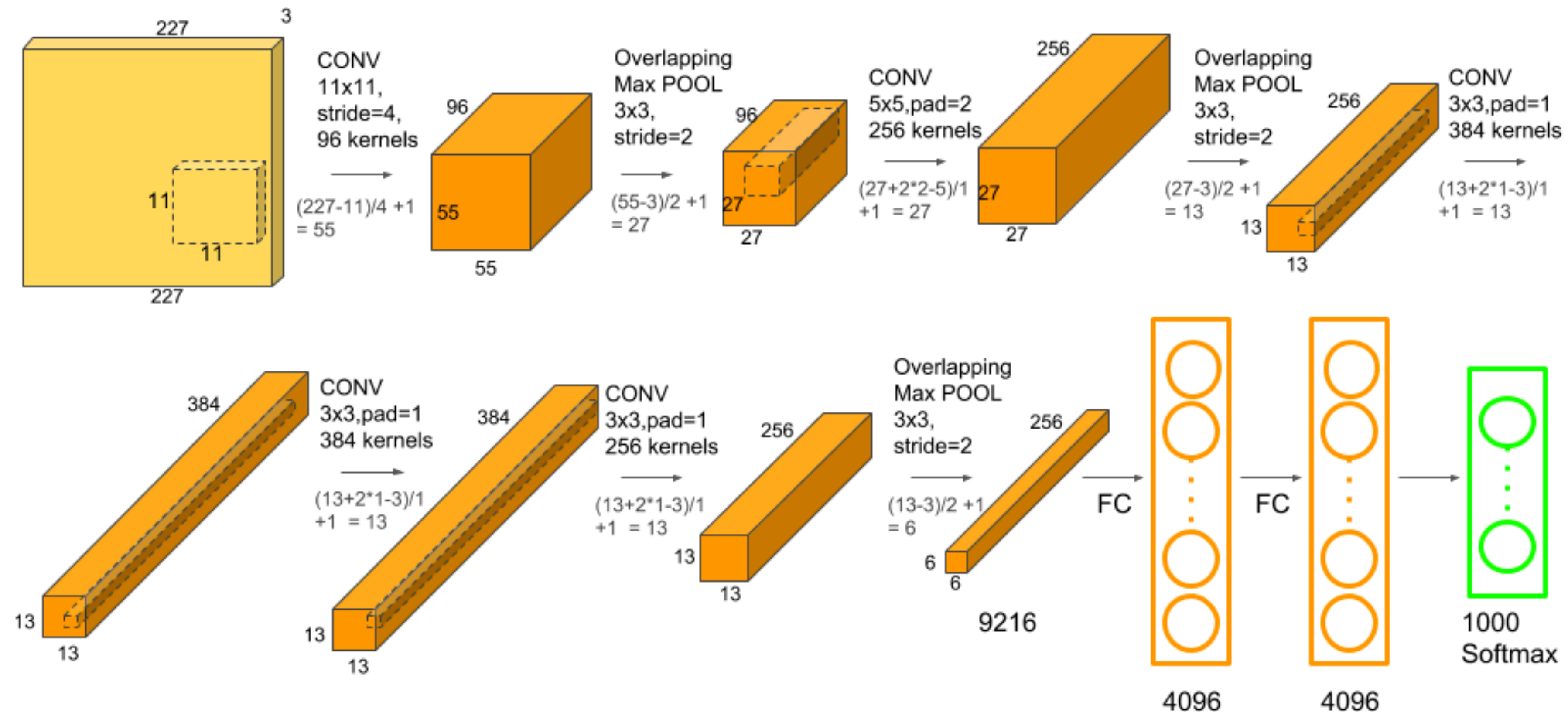
- Block that can be repeated: Convolutional layer, followed by non-linearity (e.g. ReLU) + Max pooling
- Fully connected layers before classification

LeNet-5



- One of the first successful CNN architectures
- Used to classify images of hand-written digits

AlexNet



- Winner (by a large margin) of the ImageNet challenge in 2012.
- Much larger than previous architectures.