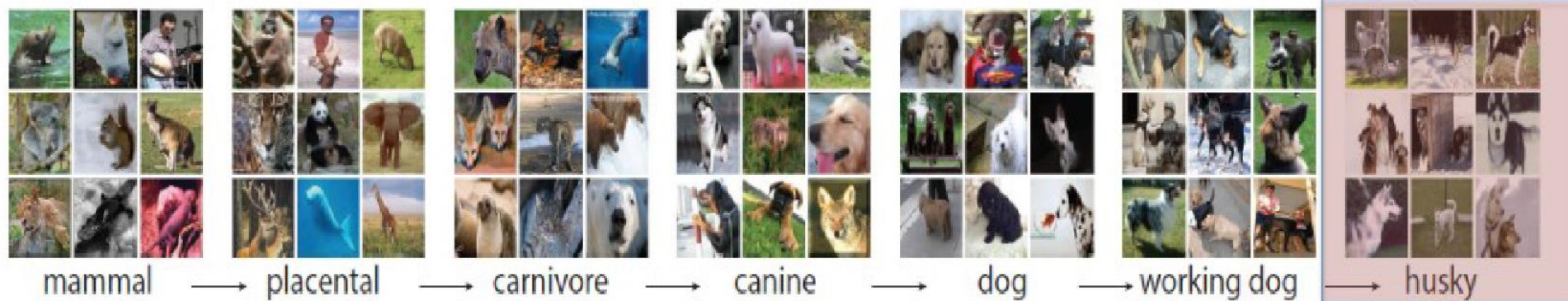


# Convolutional Neural Networks

# Dataset: ImageNet 2012



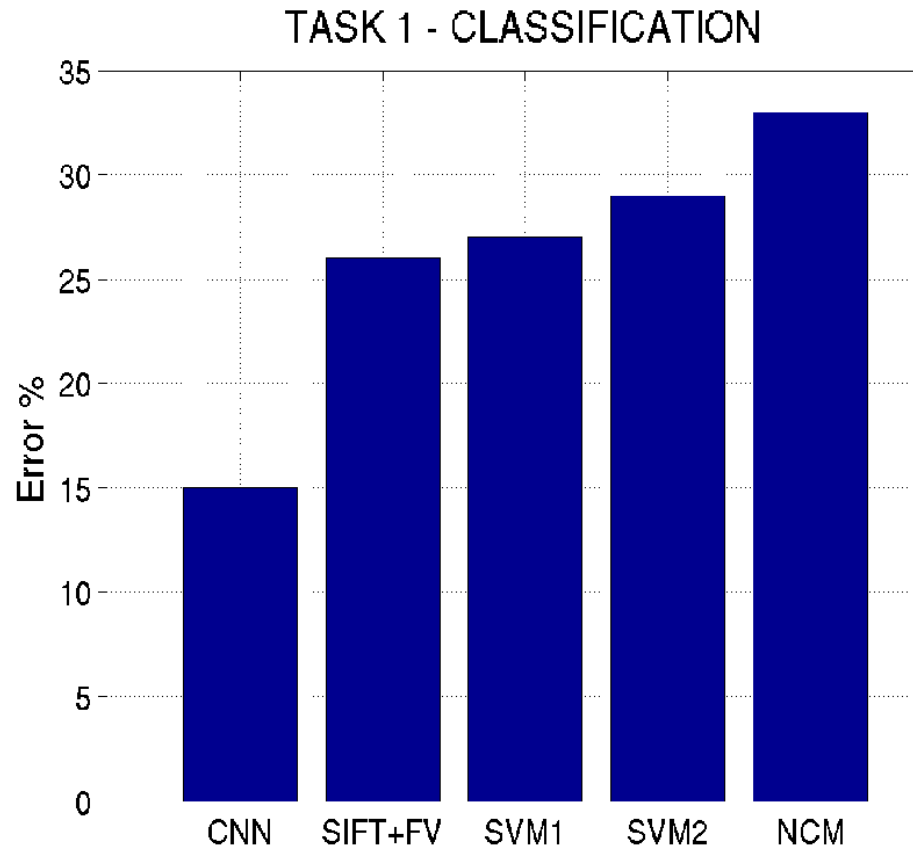
- S: (n) [Eskimo dog](#), [husky](#) (breed of heavy-coated Arctic sled dog)
  - *direct hypernym / inherited hypernym / sister term*
    - S: (n) [working dog](#) (any of several breeds of usually large powerful dogs bred to work as draft animals and guard and guide dogs)
      - S: (n) [dog](#), [domestic dog](#), [Canis familiaris](#) (a member of the genus *Canis* (probably descended from the common wolf) that has been domesticated by man since prehistoric times; occurs in many breeds) *"the dog barked all night"*
        - S: (n) [canine](#), [canid](#) (any of various fissiped mammals with nonretractile claws and typically long muzzles)
          - S: (n) [carnivore](#) (a terrestrial or aquatic flesh-eating mammal) *"terrestrial carnivores have four or five clawed digits on each limb"*
            - S: (n) [placental](#), [placental mammal](#), [eutherian](#), [eutherian mammal](#) (mammals having a placenta; all mammals except monotremes and marsupials)
              - S: (n) [mammal](#), [mammalian](#) (any warm-blooded vertebrate having the skin more or less covered with hair; young are born alive except for the small subclass of monotremes and nourished with milk)
                - S: (n) [vertebrate](#), [craniate](#) (animals having a bony or cartilaginous skeleton with a segmented spinal column and a large brain enclosed in a skull or cranium)
                  - S: (n) [chordate](#) (any animal of the phylum Chordata having a notochord or spinal column)
                    - S: (n) [animal](#), [animate being](#), [beast](#), [brute](#), [creature](#), [fauna](#) (a living organism characterized by voluntary movement)
                      - S: (n) [organism](#), [being](#) (a living thing that has (or can develop) the ability to act or function independently)
                        - S: (n) [living thing](#), [animate thing](#) (a living (or once living) entity)
                          - S: (n) [whole](#), [unit](#) (an assemblage of parts that is regarded as a single entity) *"how big is that part compared to the whole?"*; *"the team is a unit"*
                            - S: (n) [object](#), [physical object](#) (a tangible and visible entity, an entity that can cast a shadow) *"it was full of rackets, balls and other objects"*
                              - S: (n) [physical entity](#) (an entity that has physical existence)
                                - S: (n) [entity](#) (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))

# ImageNet

Examples of hammer:



# ImageNet Large Scale Visual Recognition Challenge, 2012

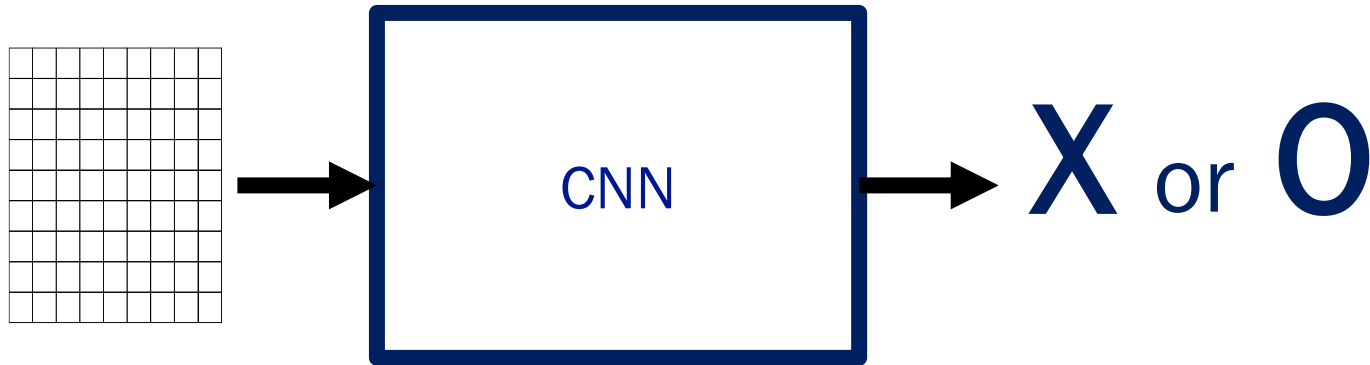


Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." Advances in Neural Information Processing Systems. 2012

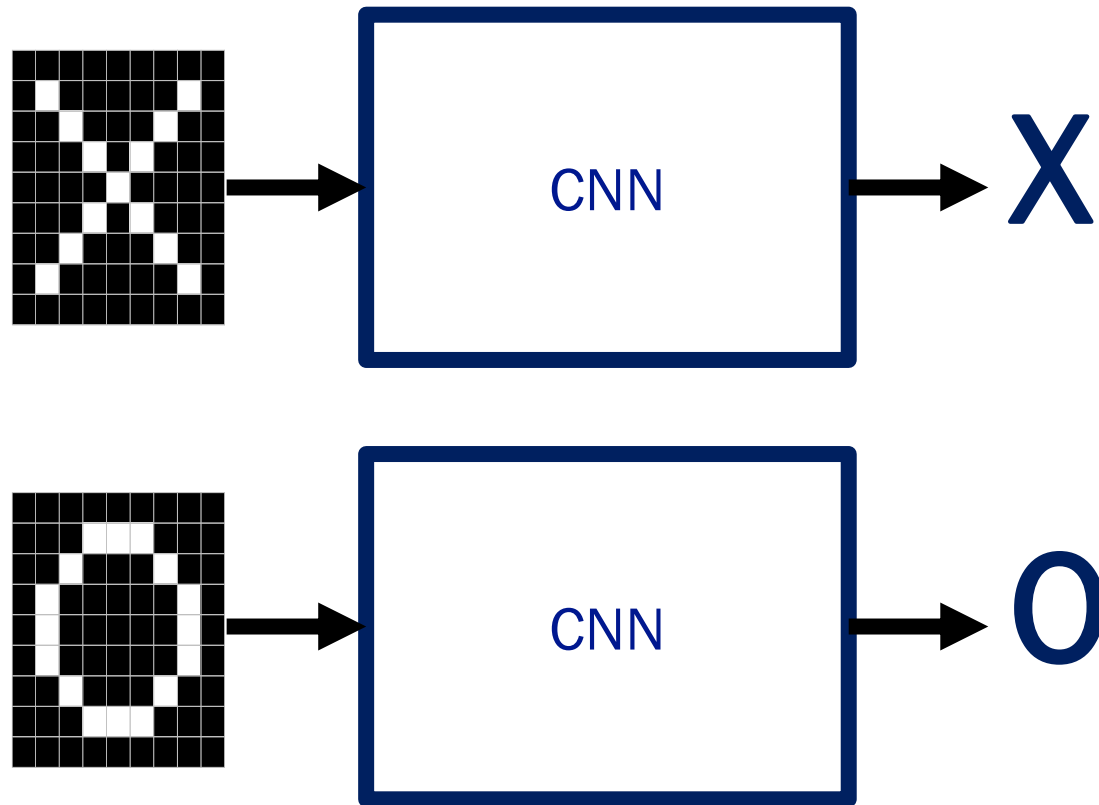
# A toy ConvNet: X's and O's

Says whether a picture is of an X or an O

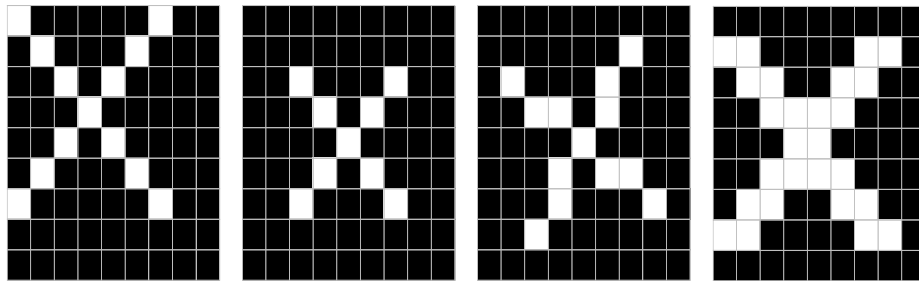
A two-dimensional  
array of pixels



# For example



# Trickier cases

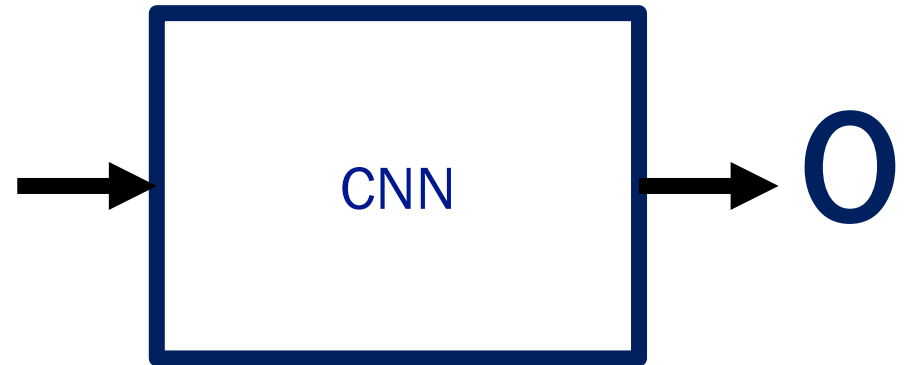
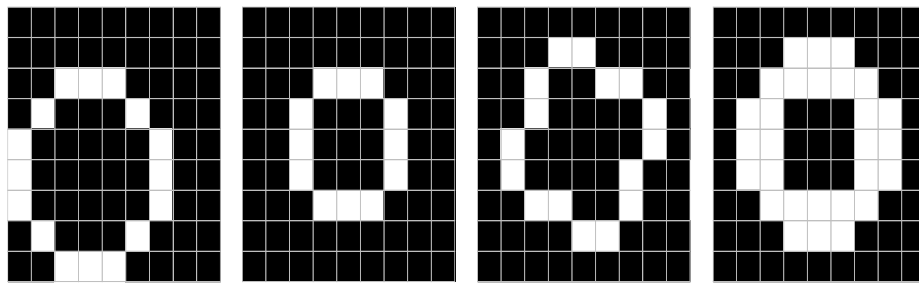


translation

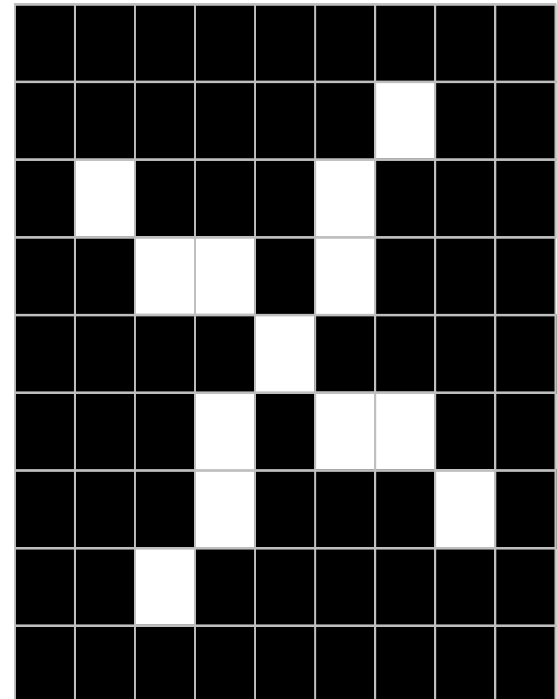
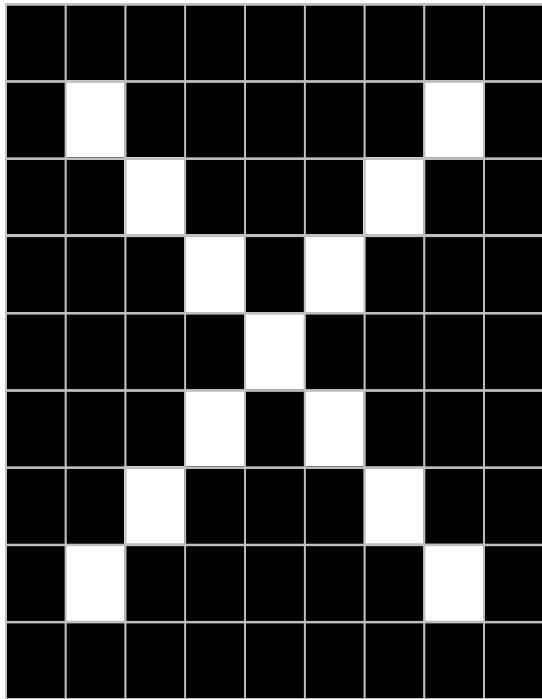
scaling

rotation

weight



# Deciding is hard



# What computers see



|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

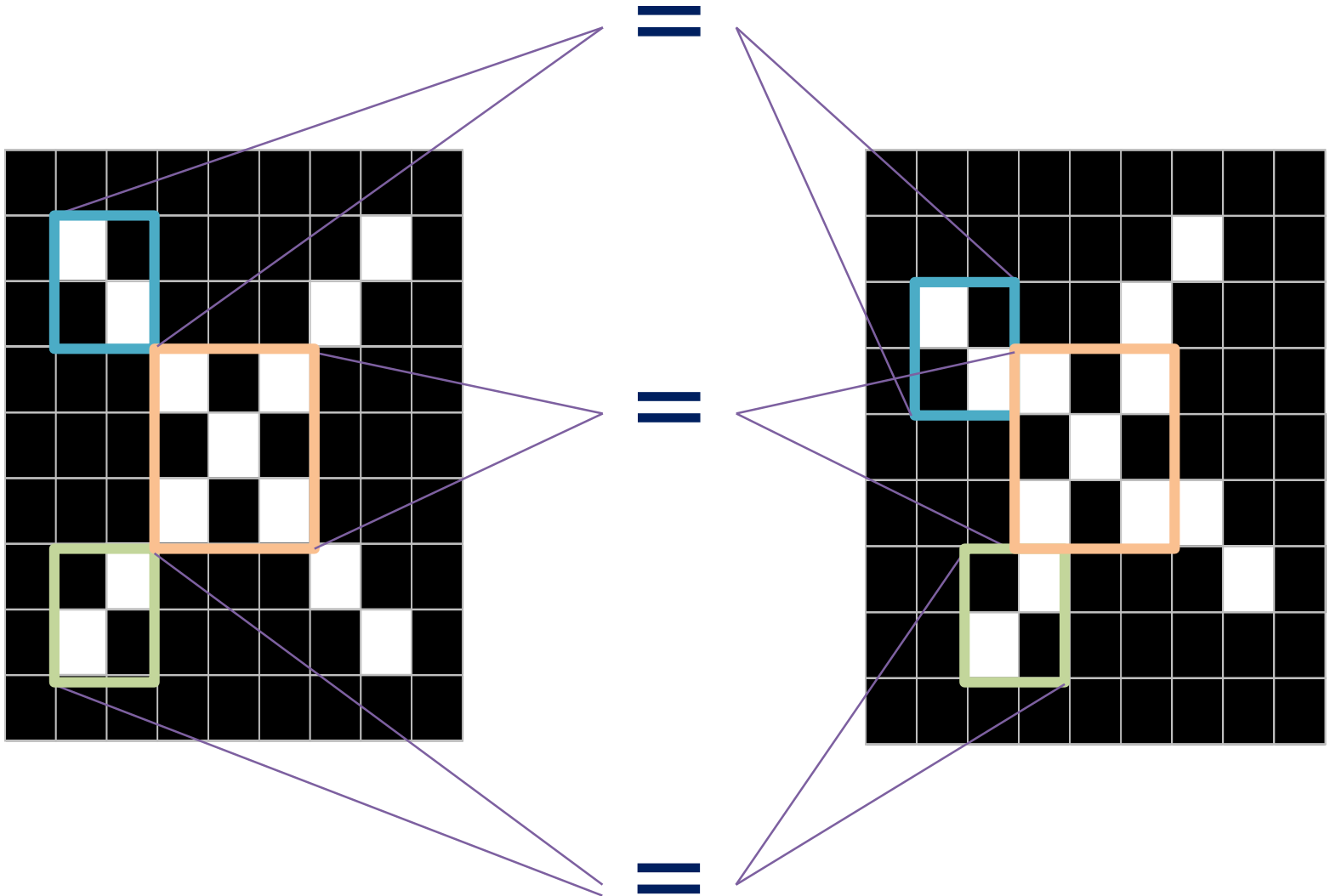


|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

# What computers see

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | X  | -1 | -1 | -1 | -1 | X  | X  | -1 |
| -1 | X  | X  | -1 | -1 | X  | X  | -1 | -1 |
| -1 | -1 | X  | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | X  | -1 | -1 |
| -1 | -1 | X  | X  | -1 | -1 | X  | X  | -1 |
| -1 | X  | X  | -1 | -1 | -1 | -1 | X  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

# ConvNets match pieces of the image



# Features match pieces of the image

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

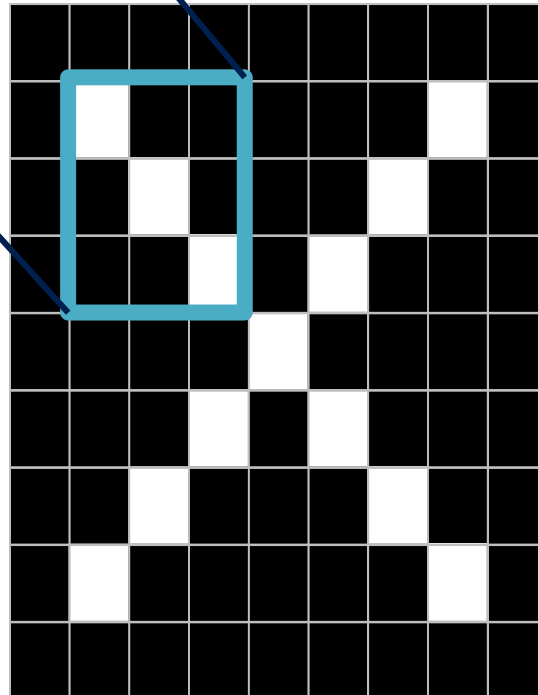
|    |    |    |
|----|----|----|
| 1  | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | 1  |

|    |    |    |
|----|----|----|
| -1 | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | -1 |

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

|    |    |    |
|----|----|----|
| 1  | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | 1  |

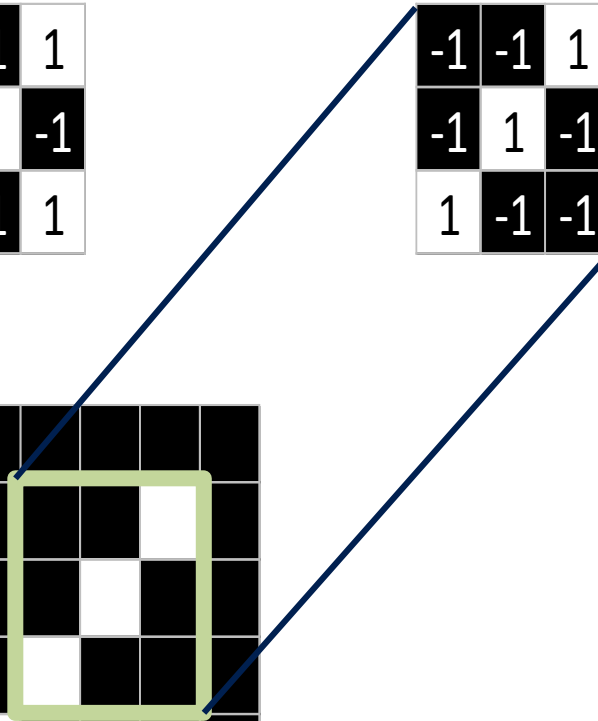
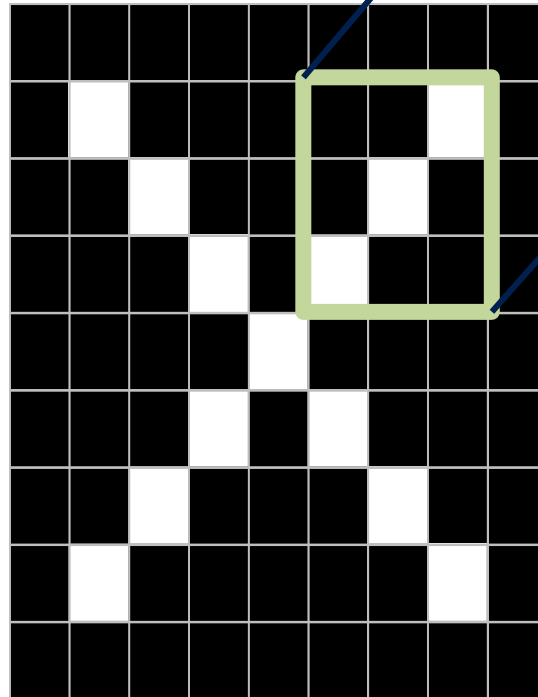
|    |    |    |
|----|----|----|
| -1 | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | -1 |



|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

|    |    |    |
|----|----|----|
| 1  | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | 1  |

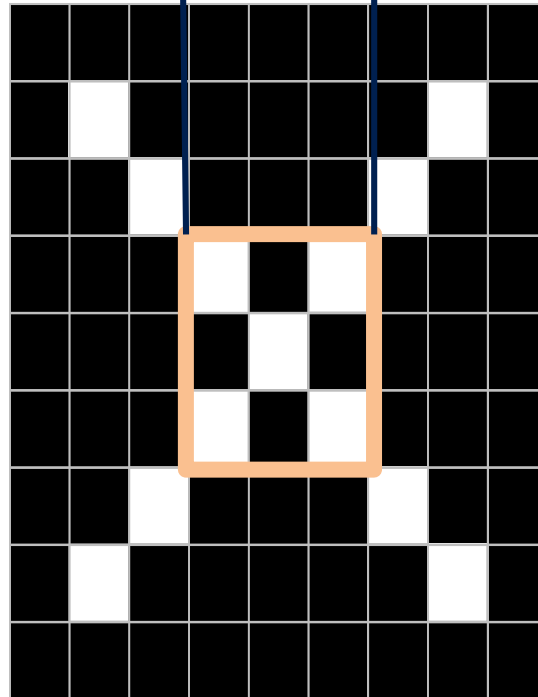
|    |    |    |
|----|----|----|
| -1 | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | -1 |



|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

|    |    |    |
|----|----|----|
| 1  | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | 1  |

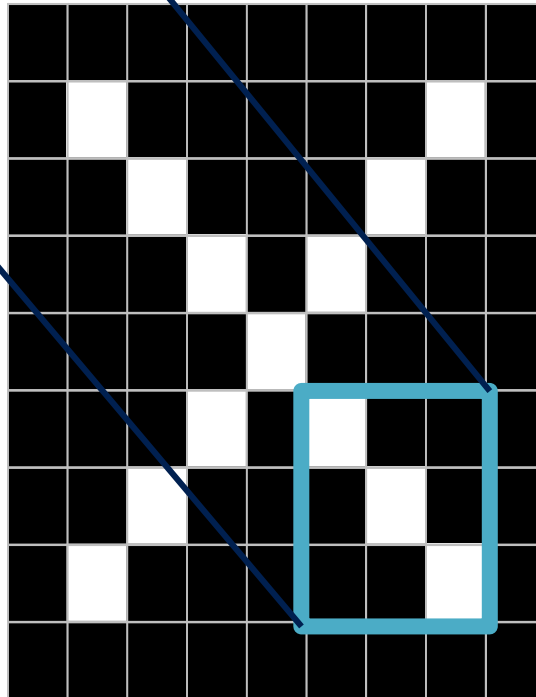
|    |    |    |
|----|----|----|
| -1 | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | -1 |



|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

|    |    |    |
|----|----|----|
| 1  | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | 1  |

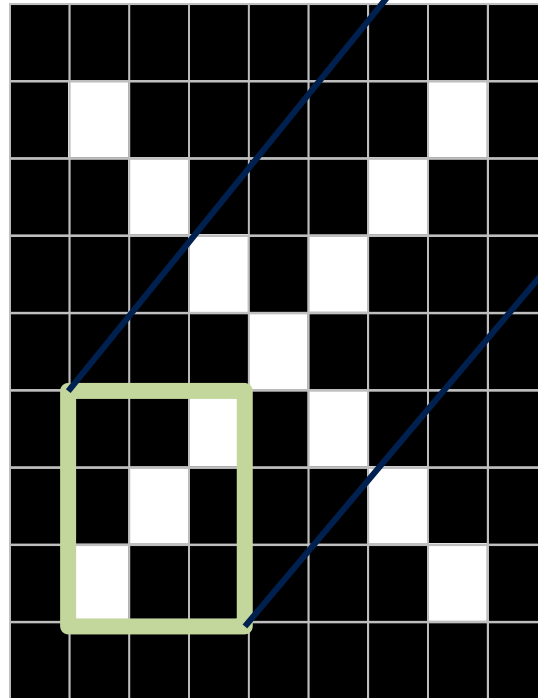
|    |    |    |
|----|----|----|
| -1 | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | -1 |



|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

|    |    |    |
|----|----|----|
| 1  | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | 1  |

|    |    |    |
|----|----|----|
| -1 | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | -1 |



|    |    |    |
|----|----|----|
| -1 | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | -1 |

# Filtering: The math behind the match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

# Filtering: The math behind the match

1. Line up the feature and the image patch.
2. Multiply each image pixel by the corresponding feature pixel.
3. Add them up.
4. Divide by the total number of pixels in the feature.

# Filtering: The math behind the match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

$$1 \times 1 = 1$$

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

# Filtering: The math behind the match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

$$1 \times 1 = 1$$

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

|   |  |  |
|---|--|--|
| 1 |  |  |
|   |  |  |
|   |  |  |

# Filtering: The math behind the match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

$$-1 \times -1 = 1$$

|   |   |  |
|---|---|--|
| 1 | 1 |  |
|   |   |  |
|   |   |  |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

# Filtering: The math behind the match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

$$-1 \times -1 = 1$$

|   |   |   |
|---|---|---|
| 1 | 1 | 1 |
|   |   |   |
|   |   |   |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

# Filtering: The math behind the match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

$$-1 \times -1 = 1$$

|   |   |   |
|---|---|---|
| 1 | 1 | 1 |
| 1 |   |   |
|   |   |   |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |



# Filtering: The math behind the match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

$$-1 \times -1 = 1$$

|   |   |   |
|---|---|---|
| 1 | 1 | 1 |
| 1 | 1 | 1 |
|   |   |   |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

# Filtering: The math behind the match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

$$\boxed{-1} \times \boxed{-1} = 1$$

|   |   |   |
|---|---|---|
| 1 | 1 | 1 |
| 1 | 1 | 1 |
| 1 |   |   |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

# Filtering: The math behind the match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

$$-1 \times -1 = 1$$

|   |   |   |
|---|---|---|
| 1 | 1 | 1 |
| 1 | 1 | 1 |
| 1 | 1 |   |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |



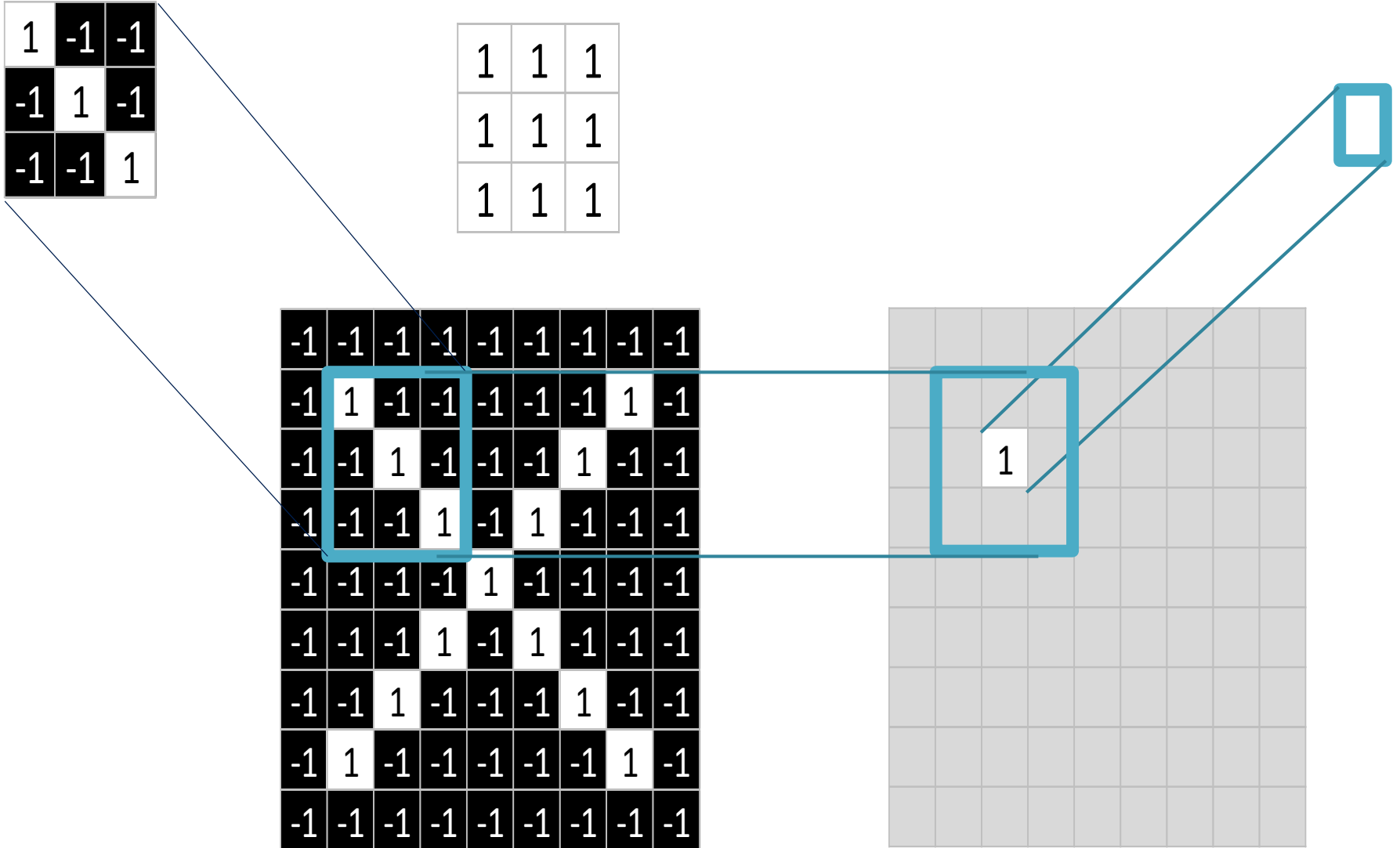
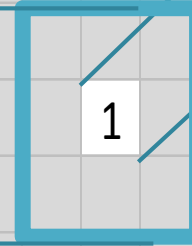
# Filtering: The math behind the match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

|   |   |   |
|---|---|---|
| 1 | 1 | 1 |
| 1 | 1 | 1 |
| 1 | 1 | 1 |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

|  |  |  |  |  |  |  |  |  |  |
|--|--|--|--|--|--|--|--|--|--|
|  |  |  |  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |  |  |  |



# Filtering: The math behind the match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

$$1 \times 1 = 1$$

|   |  |  |
|---|--|--|
| 1 |  |  |
|   |  |  |
|   |  |  |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

# Filtering: The math behind the match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

$$\begin{matrix} \boxed{-1} \\ \times \\ \boxed{1} \end{matrix} = -1$$

|   |   |    |
|---|---|----|
| 1 | 1 | -1 |
|   |   |    |
|   |   |    |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

# Filtering: The math behind the match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

|    |   |    |
|----|---|----|
| 1  | 1 | -1 |
| 1  | 1 | 1  |
| -1 | 1 | 1  |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

# Filtering: The math behind the match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

|    |   |    |
|----|---|----|
| 1  | 1 | -1 |
| 1  | 1 | 1  |
| -1 | 1 | 1  |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

|  |  |   |  |  |  |  |  |  |
|--|--|---|--|--|--|--|--|--|
|  |  |   |  |  |  |  |  |  |
|  |  |   |  |  |  |  |  |  |
|  |  | 1 |  |  |  |  |  |  |
|  |  |   |  |  |  |  |  |  |
|  |  |   |  |  |  |  |  |  |
|  |  |   |  |  |  |  |  |  |
|  |  |   |  |  |  |  |  |  |
|  |  |   |  |  |  |  |  |  |
|  |  |   |  |  |  |  |  |  |
|  |  |   |  |  |  |  |  |  |



# Convolution: Trying every possible match

|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |



|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.11  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |

# Convolution: Trying every possible match

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |



|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |



|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.11  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |



|    |    |    |
|----|----|----|
| 1  | -1 | -1 |
| -1 | 1  | -1 |
| -1 | -1 | 1  |



|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.11  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |



|    |    |    |
|----|----|----|
| 1  | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | 1  |



|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.33  | -0.55 | 0.11  | -0.11 | 0.11  | -0.55 | 0.33  |
| -0.55 | 0.55  | -0.55 | 0.33  | -0.55 | 0.55  | -0.55 |
| 0.11  | -0.55 | 0.55  | -0.77 | 0.55  | -0.55 | 0.11  |
| -0.11 | 0.33  | -0.77 | 1.00  | -0.77 | 0.33  | -0.11 |
| 0.11  | -0.55 | 0.55  | -0.77 | 0.55  | -0.55 | 0.11  |
| -0.55 | 0.55  | -0.55 | 0.33  | -0.55 | 0.55  | -0.55 |
| 0.33  | -0.55 | 0.11  | -0.11 | 0.11  | -0.55 | 0.33  |

|    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1  | -1 | 1  | -1 | -1 | -1 |
| -1 | -1 | 1  | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |



|    |    |    |
|----|----|----|
| -1 | -1 | 1  |
| -1 | 1  | -1 |
| 1  | -1 | -1 |



|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.11  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |



# Convolution layer

One image becomes a stack of filtered images

|    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | 1  | -1 | -1 | -1 | -1 | 1  |
| -1 | -1 | 1  | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | 1  | 1  | -1 | -1 |
| -1 | -1 | -1 | -1 | 1  | -1 | -1 |
| -1 | -1 | 1  | -1 | 1  | -1 | -1 |
| -1 | -1 | -1 | -1 | -1 | 1  | -1 |
| -1 | 1  | -1 | -1 | -1 | 1  | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 |



|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.11  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |
| 0.33  | -0.55 | 0.11  | -0.11 | 0.11  | -0.55 | 0.33  |
| 0.55  | 0.55  | -0.55 | 0.33  | -0.55 | 0.55  | -0.55 |
| 0.11  | -0.55 | 0.55  | -0.77 | 0.55  | -0.55 | 0.11  |
| -0.11 | 0.33  | -0.77 | 1.00  | -0.77 | 0.33  | -0.11 |
| 0.11  | -0.55 | 0.55  | -0.77 | 0.55  | -0.55 | 0.11  |
| -0.55 | 0.55  | -0.55 | 0.33  | -0.55 | 0.55  | -0.55 |
| 0.33  | -0.55 | 0.11  | -0.11 | 0.11  | -0.55 | 0.33  |
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.33  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |

# Pooling: Shrinking the image stack

1. Pick a window size (usually 2 or 3).
2. Pick a stride (usually 2).
3. Walk your window across your filtered images.
4. From each window, take the maximum value.

# Pooling

maximum

|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.11  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |

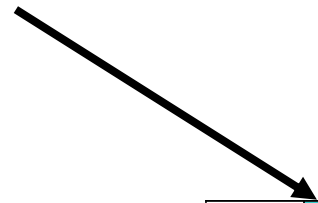
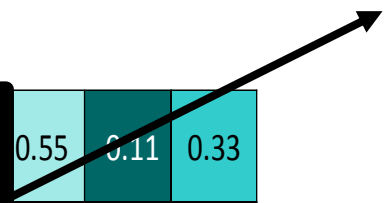
|      |  |  |  |
|------|--|--|--|
| 1.00 |  |  |  |
|      |  |  |  |
|      |  |  |  |
|      |  |  |  |
|      |  |  |  |

# Pooling

maximum

|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.11  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |

|      |      |  |  |
|------|------|--|--|
| 1.00 | 0.33 |  |  |
|      |      |  |  |
|      |      |  |  |
|      |      |  |  |
|      |      |  |  |



# Pooling

maximum

|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.11  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |

|      |      |      |  |
|------|------|------|--|
| 1.00 | 0.33 | 0.55 |  |
|      |      |      |  |
|      |      |      |  |
|      |      |      |  |

# Pooling

maximum

|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.11  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |

|      |      |      |      |
|------|------|------|------|
| 1.00 | 0.33 | 0.55 | 0.33 |
|      |      |      |      |
|      |      |      |      |
|      |      |      |      |

# Pooling

maximum

|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.11  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |

|      |      |      |      |
|------|------|------|------|
| 1.00 | 0.33 | 0.55 | 0.33 |
| 0.33 |      |      |      |
|      |      |      |      |
|      |      |      |      |

# Pooling

|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.11  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |

max pooling



|      |      |      |      |
|------|------|------|------|
| 1.00 | 0.33 | 0.55 | 0.33 |
| 0.33 | 1.00 | 0.33 | 0.55 |
| 0.55 | 0.33 | 1.00 | 0.11 |
| 0.33 | 0.55 | 0.11 | 0.77 |

|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.11  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |



|      |      |      |      |
|------|------|------|------|
| 1.00 | 0.33 | 0.55 | 0.33 |
| 0.33 | 1.00 | 0.33 | 0.55 |
| 0.55 | 0.33 | 1.00 | 0.11 |
| 0.33 | 0.55 | 0.11 | 0.77 |

|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.33  | -0.55 | 0.11  | -0.11 | 0.11  | -0.55 | 0.33  |
| -0.55 | 0.55  | -0.55 | 0.33  | -0.55 | 0.55  | -0.55 |
| 0.11  | -0.55 | 0.55  | -0.77 | 0.55  | -0.55 | 0.11  |
| -0.11 | 0.33  | -0.77 | 1.00  | -0.77 | 0.33  | -0.11 |
| 0.11  | -0.55 | 0.55  | -0.77 | 0.55  | -0.55 | 0.11  |
| -0.55 | 0.55  | -0.55 | 0.33  | -0.55 | 0.55  | -0.55 |
| 0.33  | -0.55 | 0.11  | -0.11 | 0.11  | -0.55 | 0.33  |



|      |      |      |      |
|------|------|------|------|
| 0.55 | 0.33 | 0.55 | 0.33 |
| 0.33 | 1.00 | 0.55 | 0.11 |
| 0.55 | 0.55 | 0.55 | 0.11 |
| 0.33 | 0.11 | 0.11 | 0.33 |

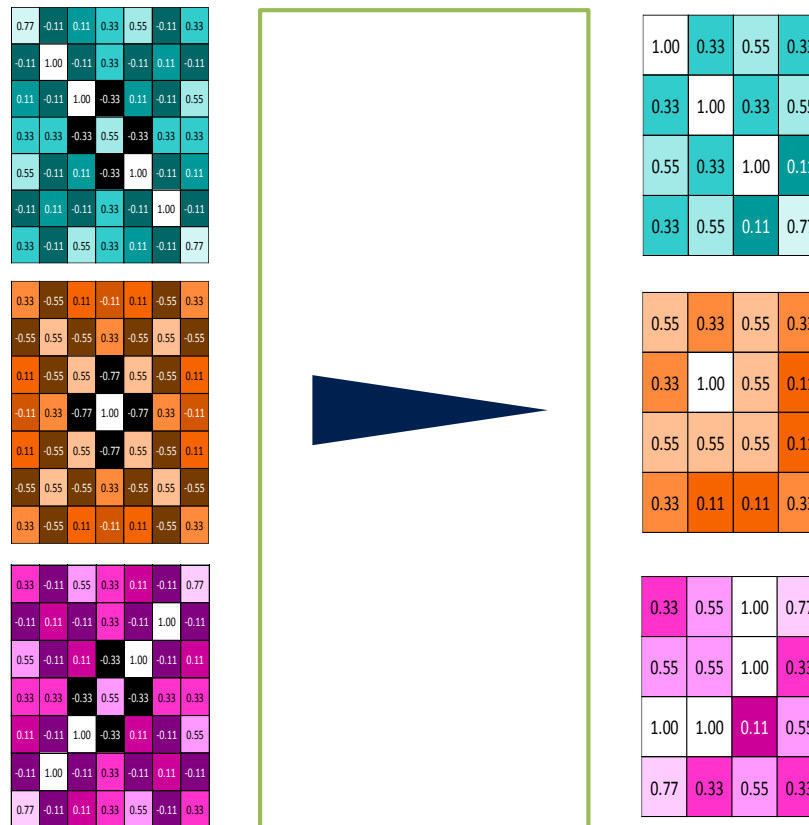
|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.11  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |



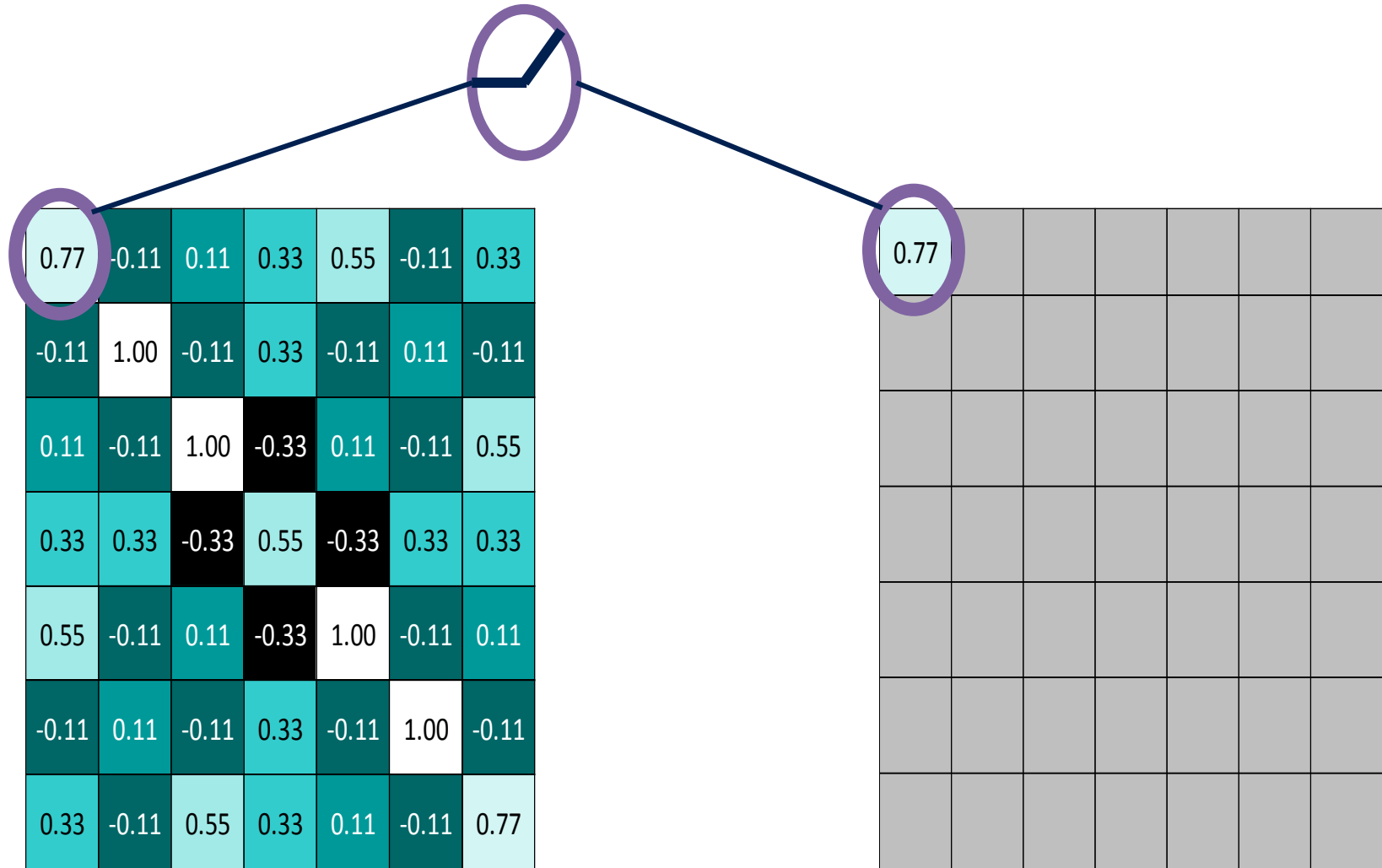
|      |      |      |      |
|------|------|------|------|
| 0.33 | 0.55 | 1.00 | 0.77 |
| 0.55 | 0.55 | 1.00 | 0.33 |
| 1.00 | 1.00 | 0.11 | 0.55 |
| 0.77 | 0.33 | 0.55 | 0.33 |

# Pooling layer

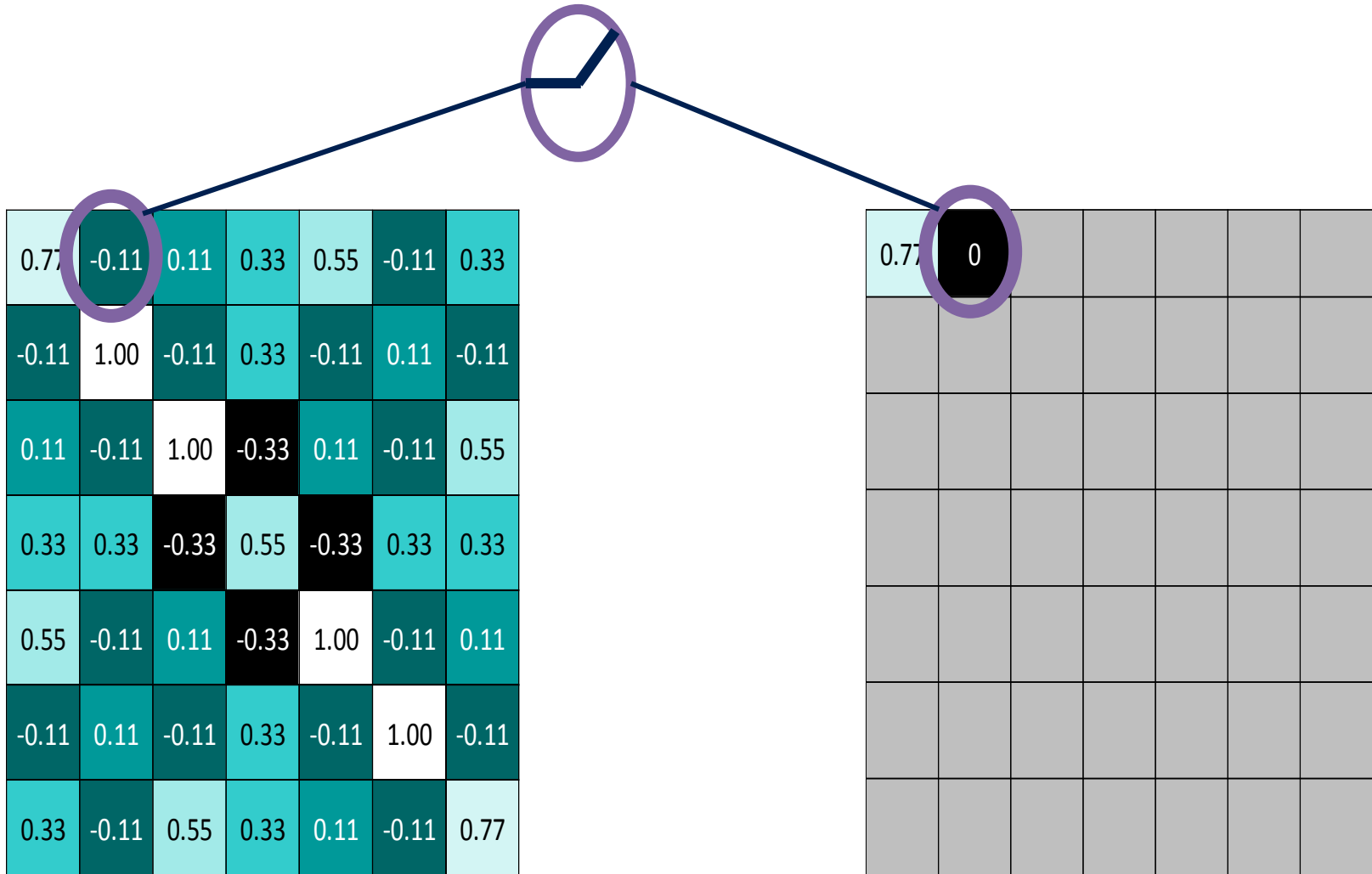
A stack of images becomes a stack of smaller images.



# Rectified Linear Units (ReLUs)



# Rectified Linear Units (ReLUs)



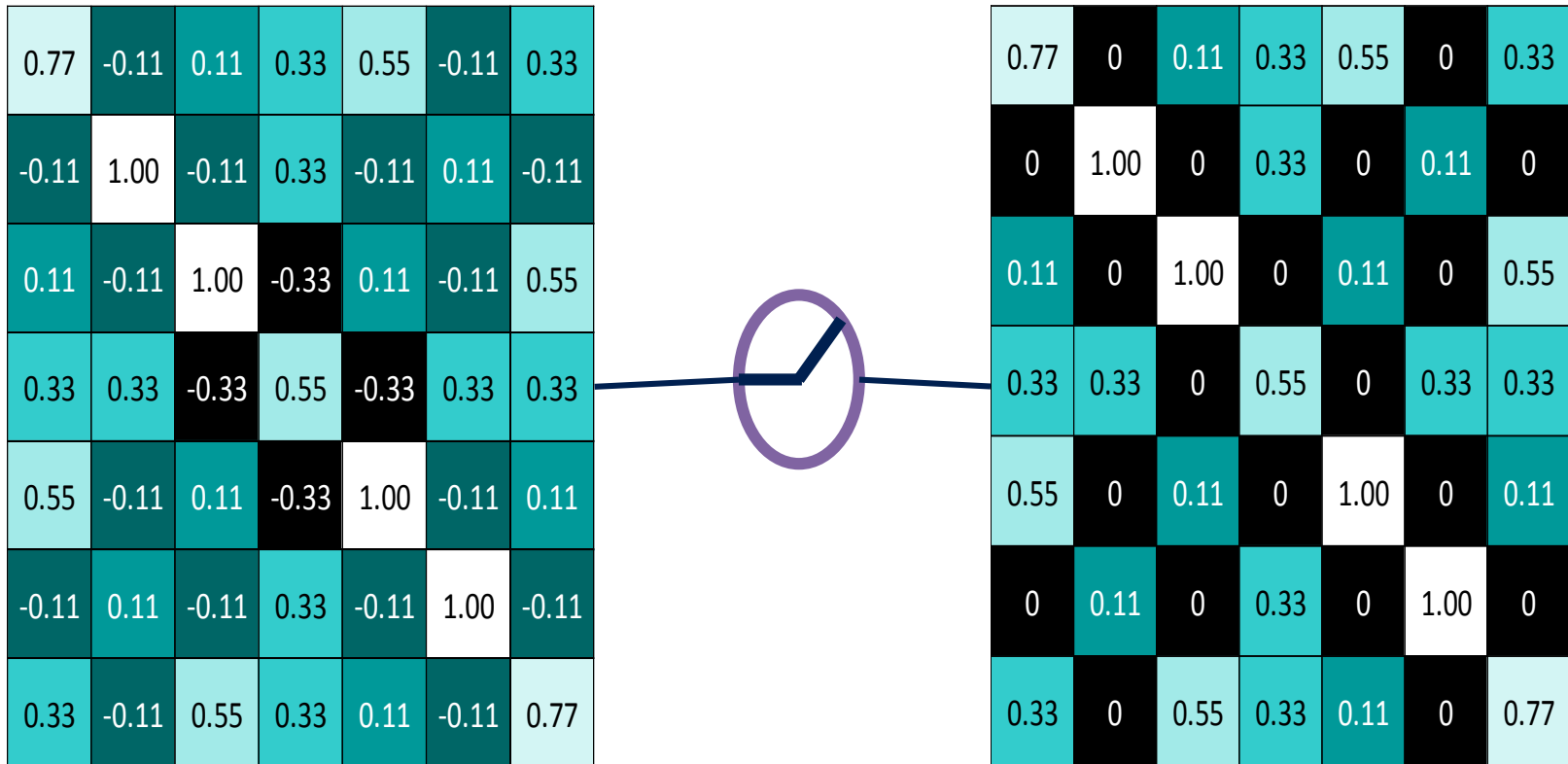
# Rectified Linear Units (ReLUs)



|       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|
| 0.77  | -0.11 | 0.11  | 0.33  | 0.55  | -0.11 | 0.33  |
| -0.11 | 1.00  | -0.11 | 0.33  | -0.11 | 0.11  | -0.11 |
| 0.11  | -0.11 | 1.00  | -0.33 | 0.11  | -0.11 | 0.55  |
| 0.33  | 0.33  | -0.33 | 0.55  | -0.33 | 0.33  | 0.33  |
| 0.55  | -0.11 | 0.11  | -0.33 | 1.00  | -0.11 | 0.11  |
| -0.11 | 0.11  | -0.11 | 0.33  | -0.11 | 1.00  | -0.11 |
| 0.33  | -0.11 | 0.55  | 0.33  | 0.11  | -0.11 | 0.77  |

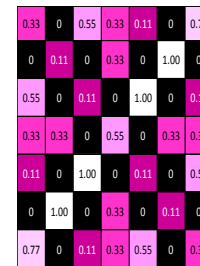
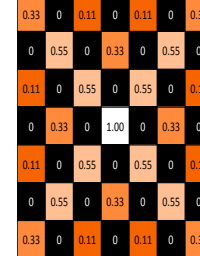
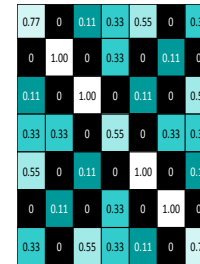
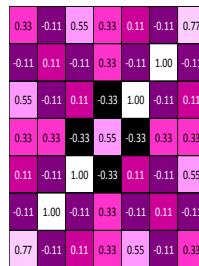
|      |   |      |      |      |   |      |
|------|---|------|------|------|---|------|
| 0.77 | 0 | 0.11 | 0.33 | 0.55 | 0 | 0.33 |
|      |   |      |      |      |   |      |
|      |   |      |      |      |   |      |
|      |   |      |      |      |   |      |
|      |   |      |      |      |   |      |
|      |   |      |      |      |   |      |
|      |   |      |      |      |   |      |

# Rectified Linear Units (ReLUs)



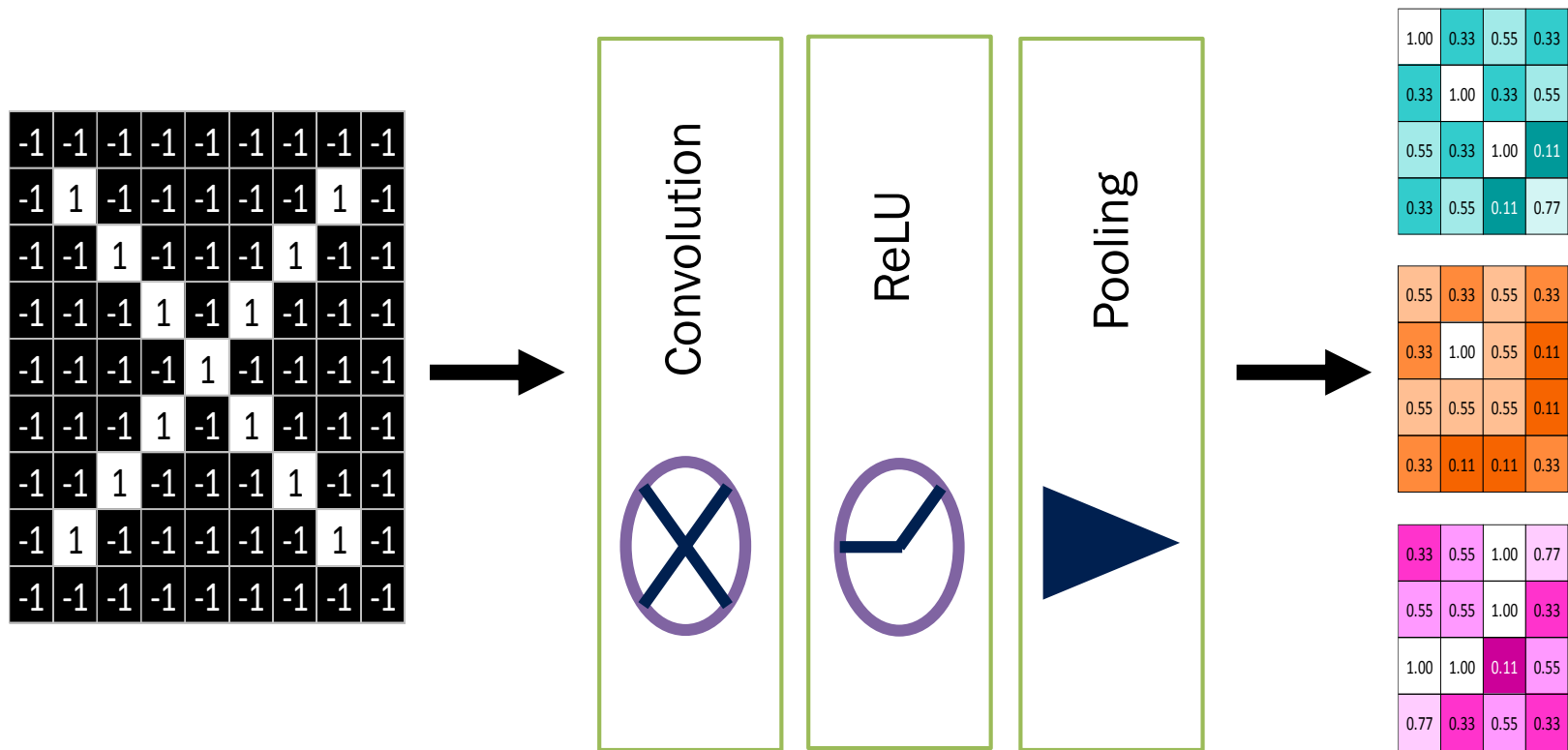
# ReLU layer

A stack of images becomes a stack of images with no negative values.



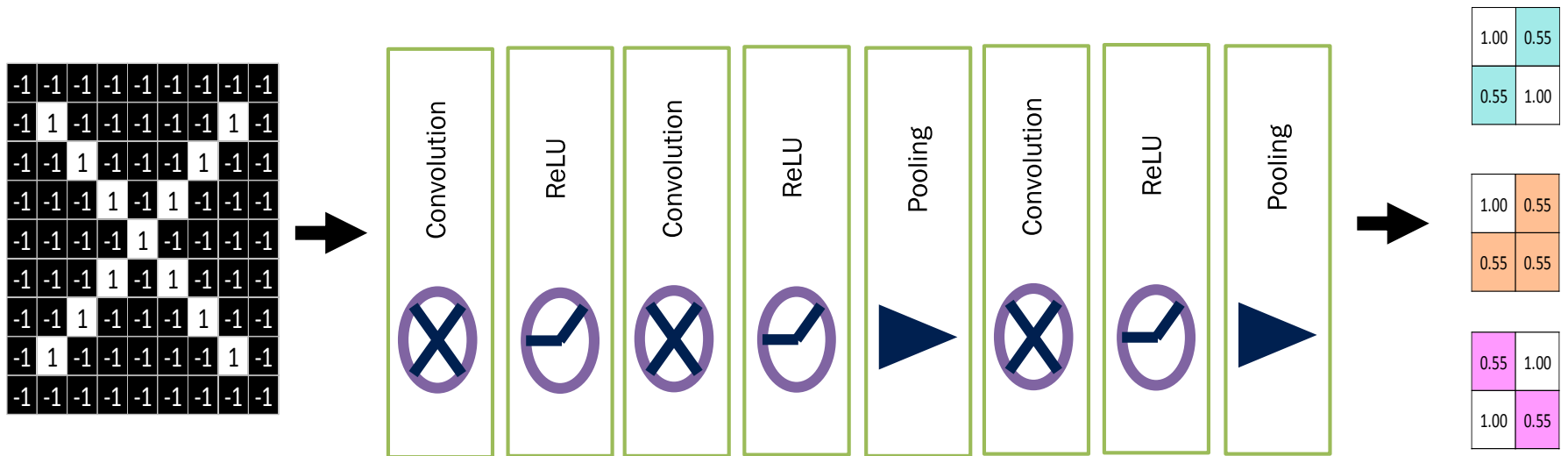
# Layers get stacked

The output of one becomes the input of the next.



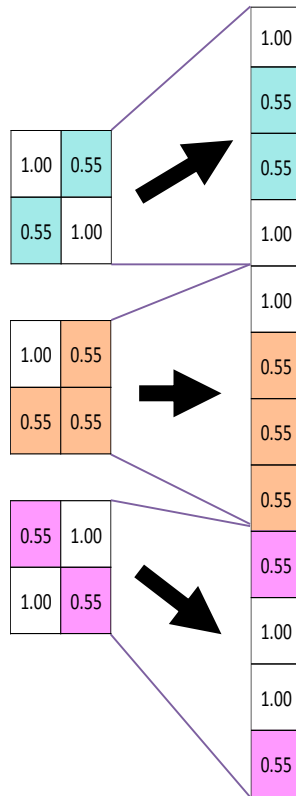
# Deep stacking

Layers can be repeated several (or many) times.



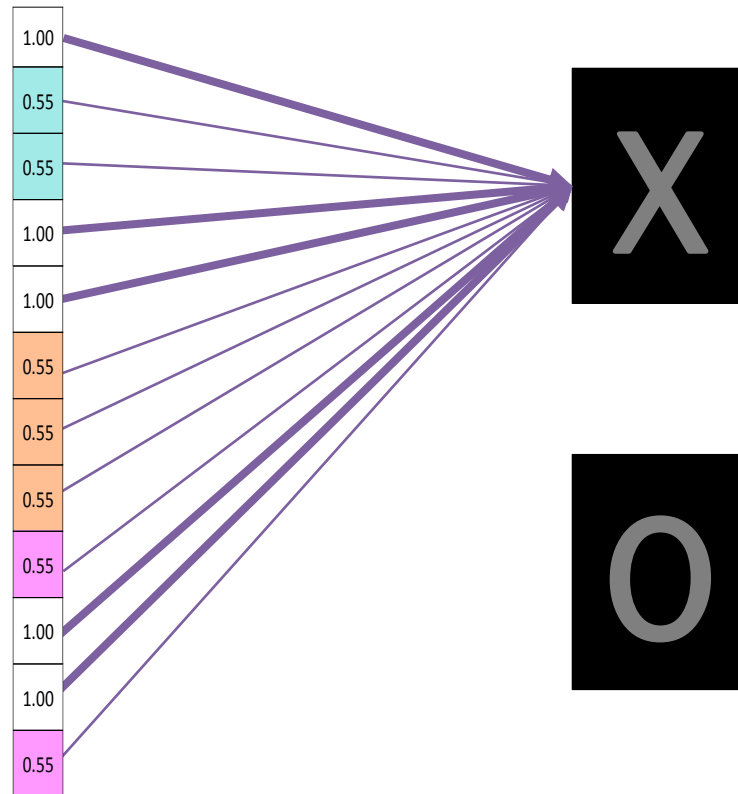
# Fully connected layer

Every value gets a vote



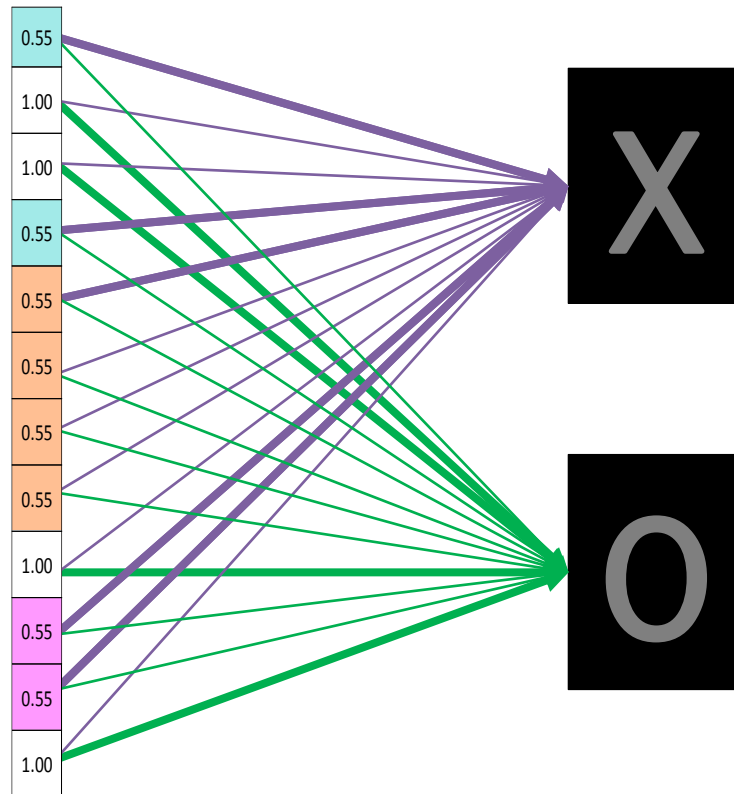
# Fully connected layer

Vote depends on how strongly a value predicts X or O



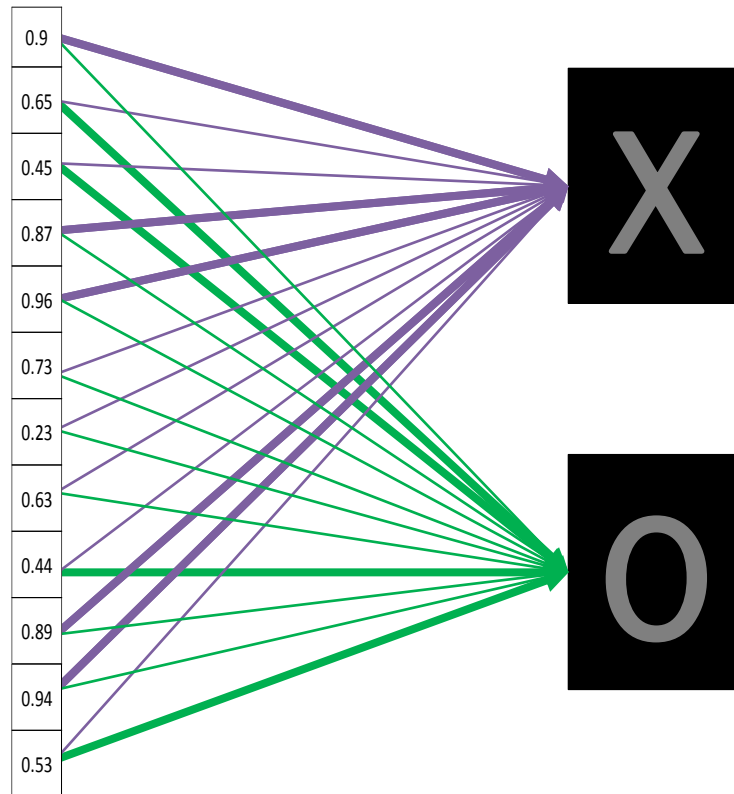
# Fully connected layer

Vote depends on how strongly a value predicts X or O



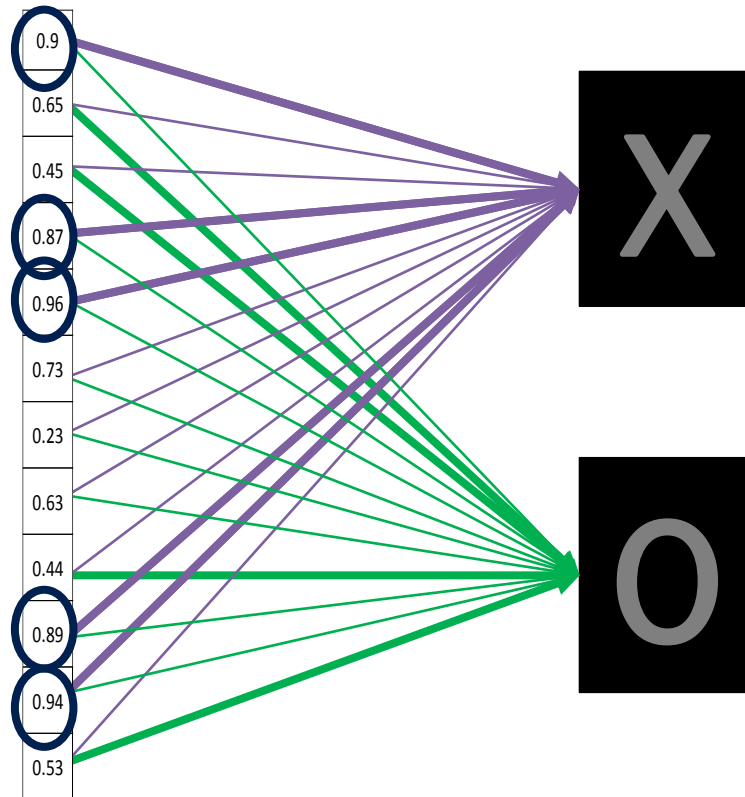
# Fully connected layer

Future values vote on X or O



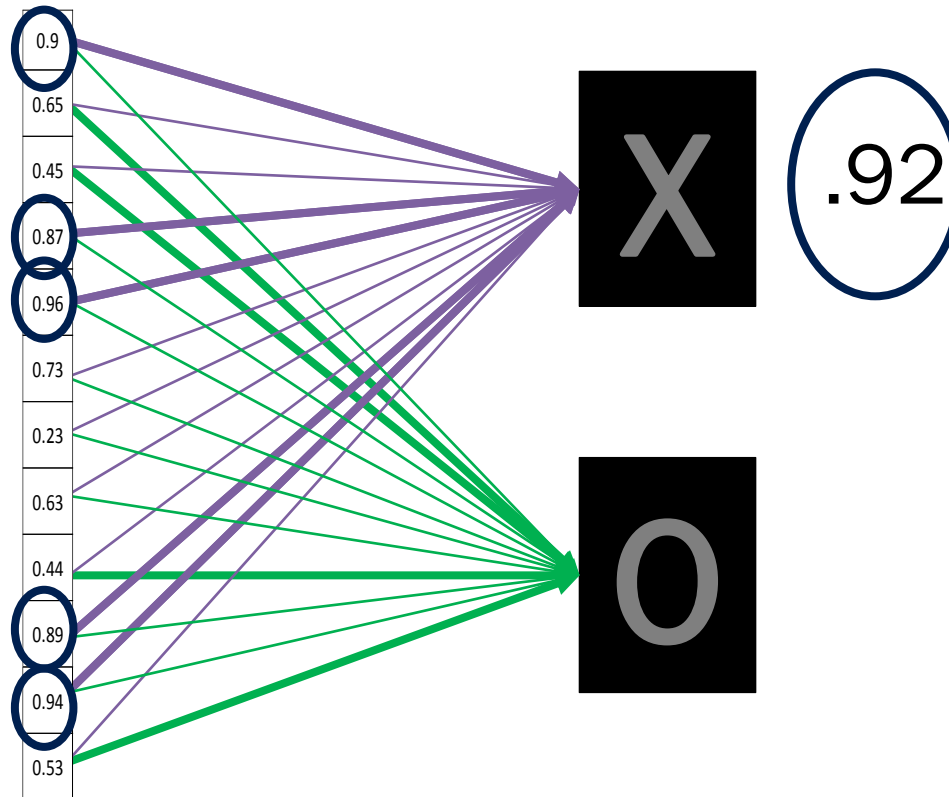
# Fully connected layer

Future values vote on X or O



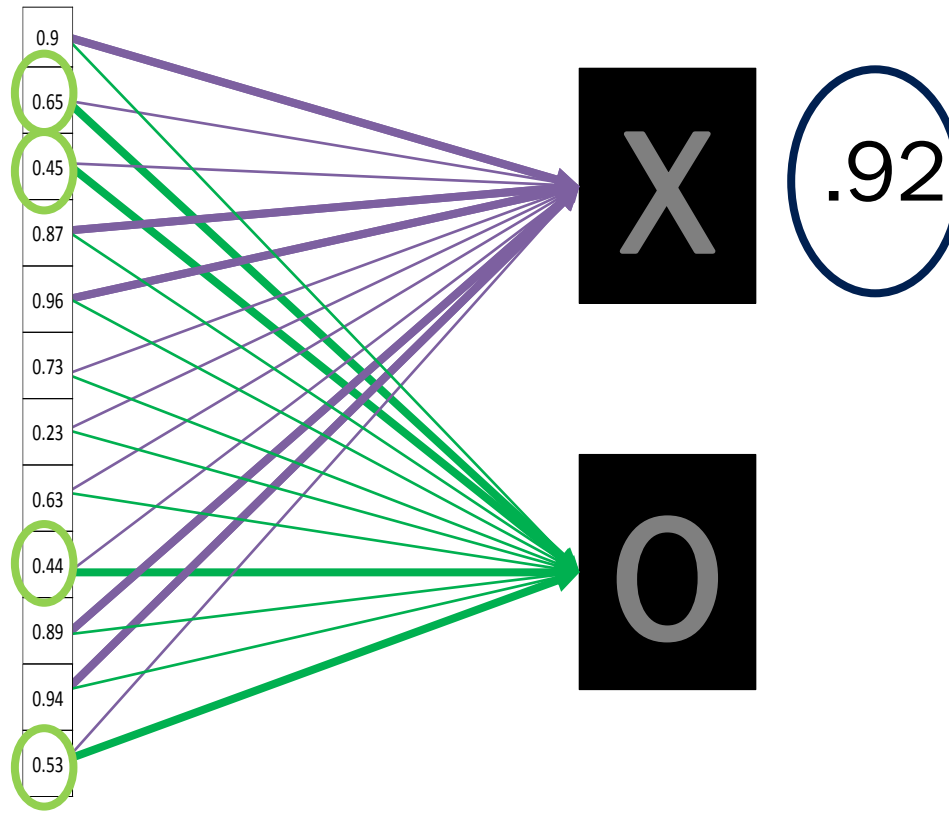
# Fully connected layer

Future values vote on X or O



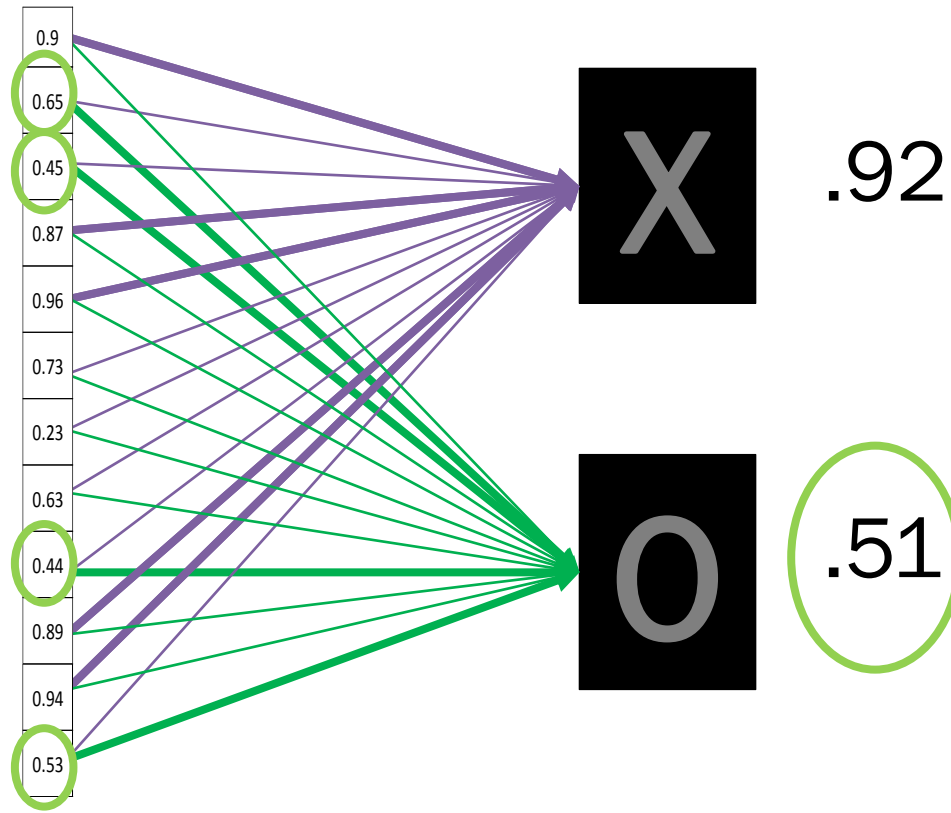
# Fully connected layer

Future values vote on X or O



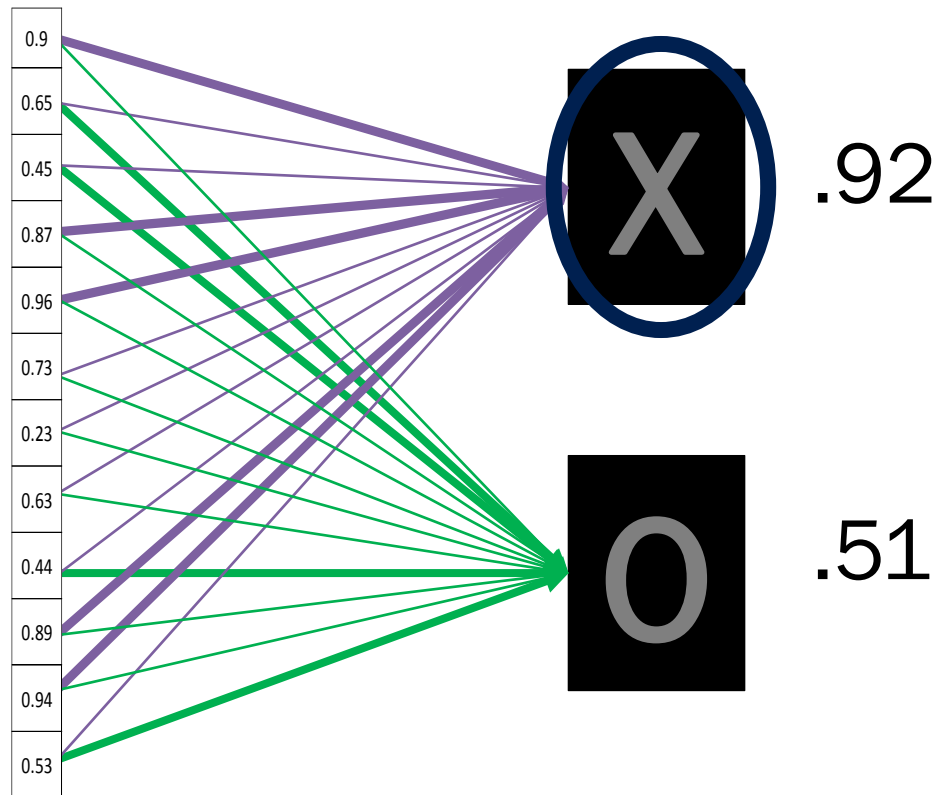
# Fully connected layer

Future values vote on X or O



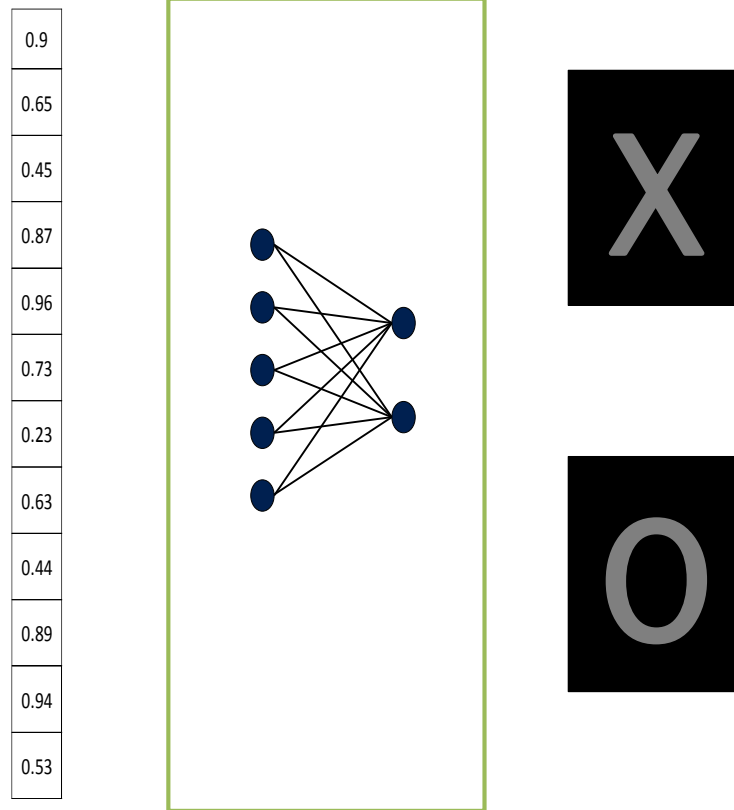
# Fully connected layer

Future values vote on X or O



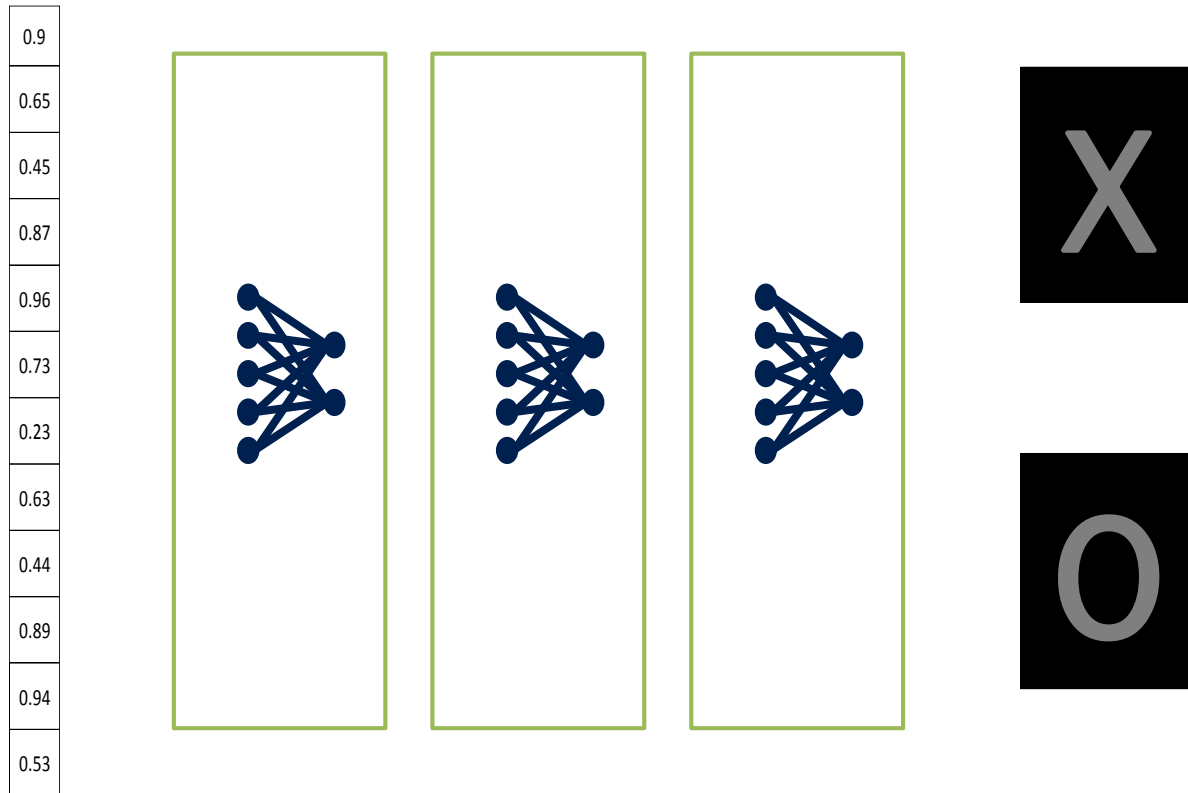
# Fully connected layer

A list of feature values becomes a list of votes.



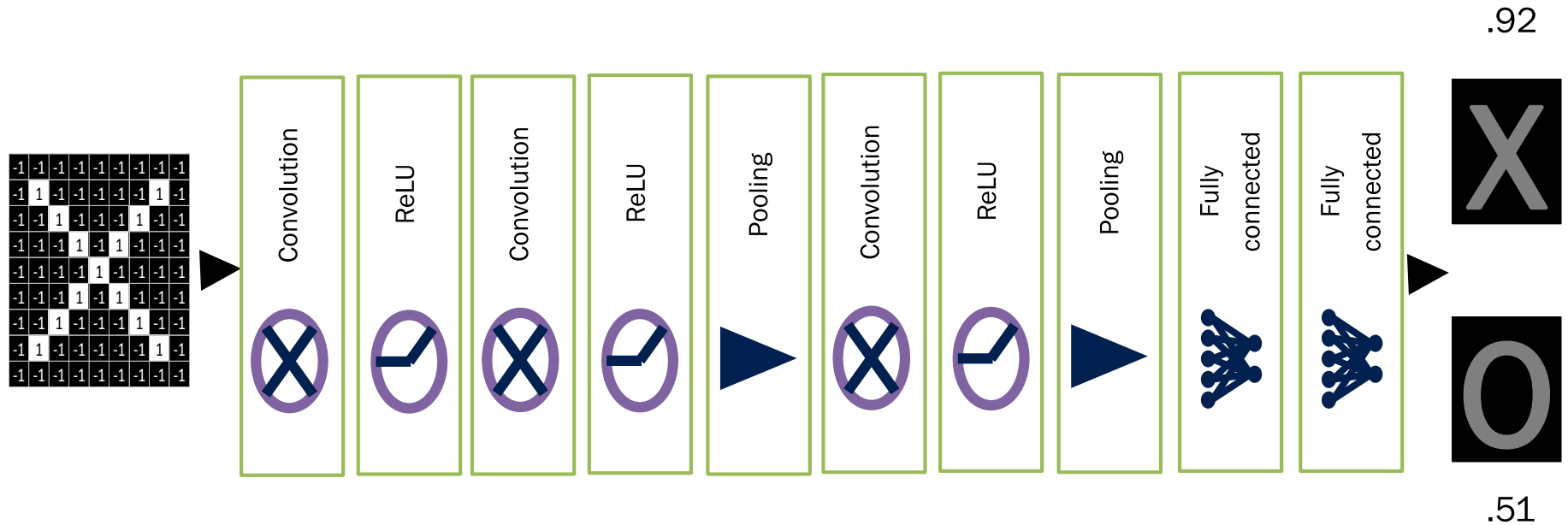
# Fully connected layer

These can also be stacked.



# Putting it all together

A set of pixels becomes a set of votes.



# Learning

Q: Where do all the magic numbers come from?

Features in convolutional layers

Voting weights in fully connected layers

A: Backpropagation

# Hyperparameters (knobs)

## Convolution

- Number of features

- Size of features

## Pooling

- Window size

- Window stride

## Fully Connected

- Number of neurons

# Architecture

How many of each type of layer?

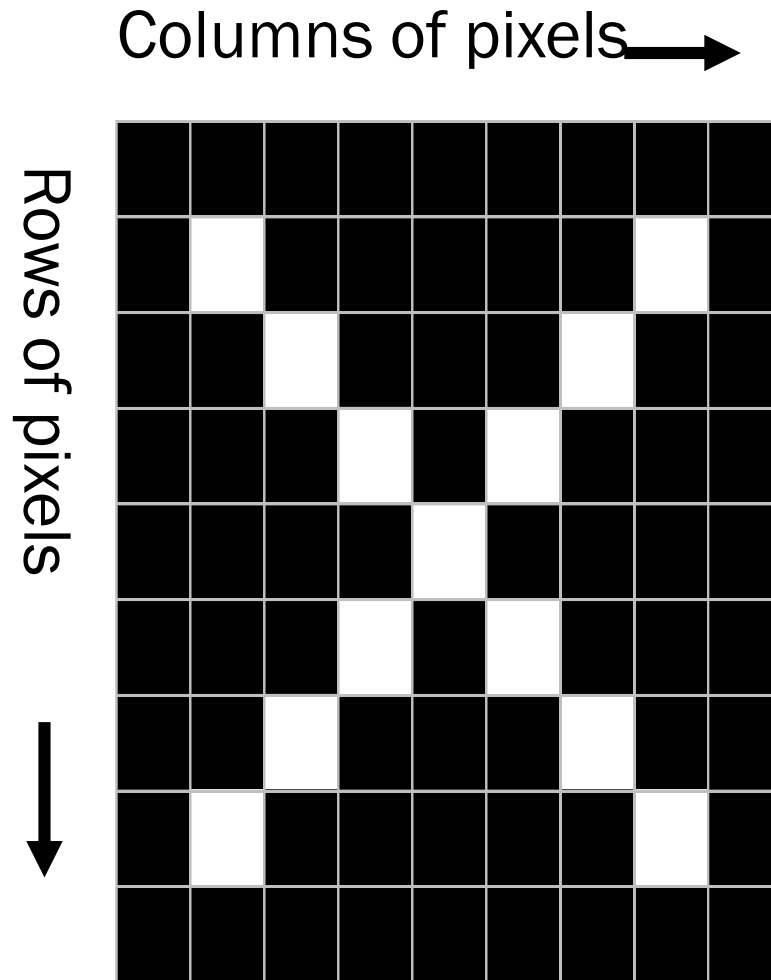
In what order?

# Not just images

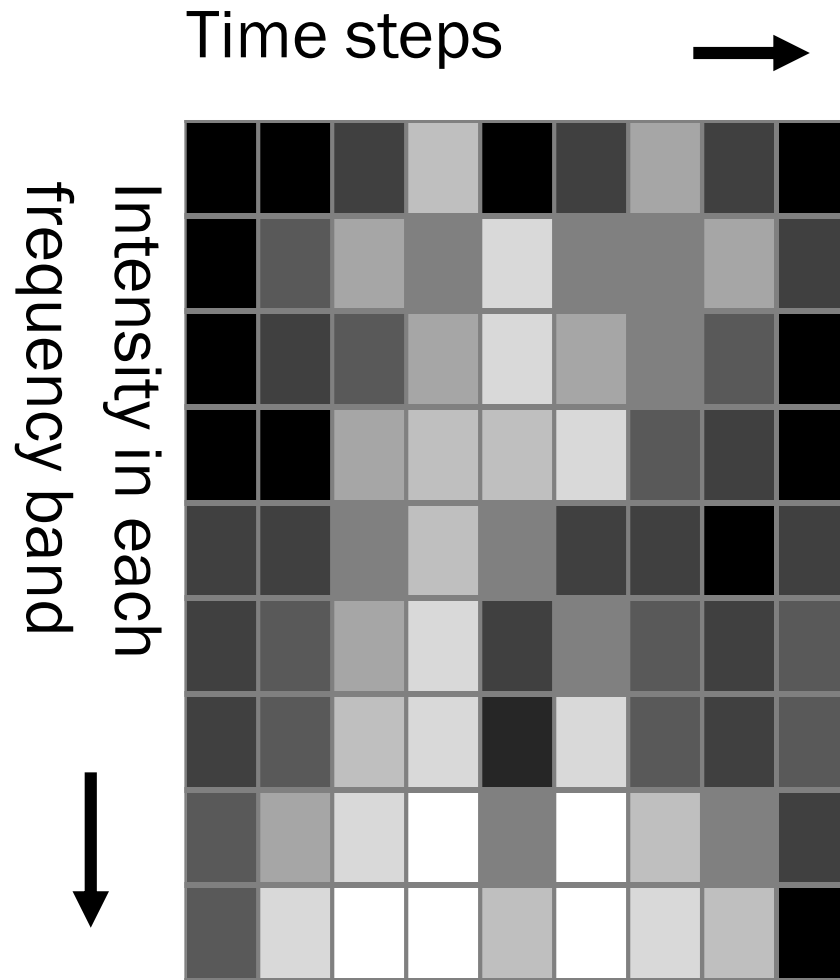
Any 2D (or 3D) data.

Things closer together are more closely related than things far away.

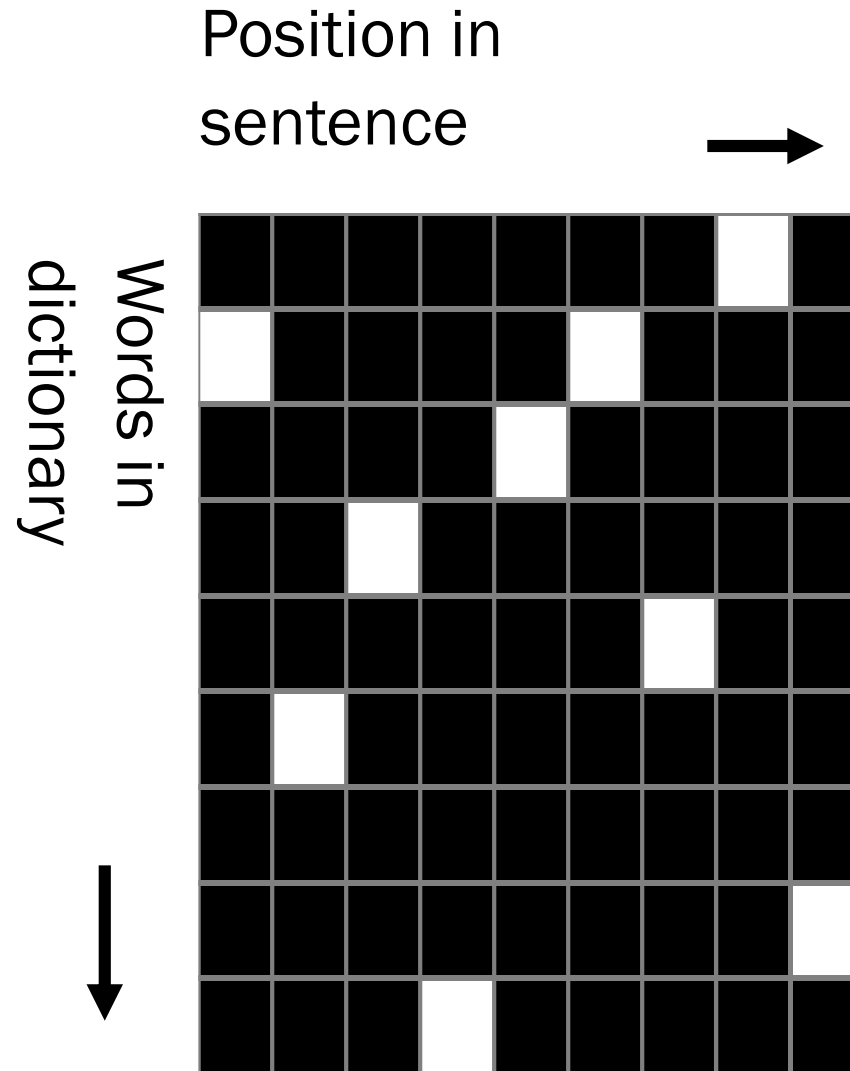
# Images



# Sound



# Text



# Limitations

ConvNets only capture local “spatial” patterns in data.

If the data can't be made to look like an image, ConvNets are less useful.

# Customer data

Name, age,  
address, email,  
purchases,  
browsing activity,...

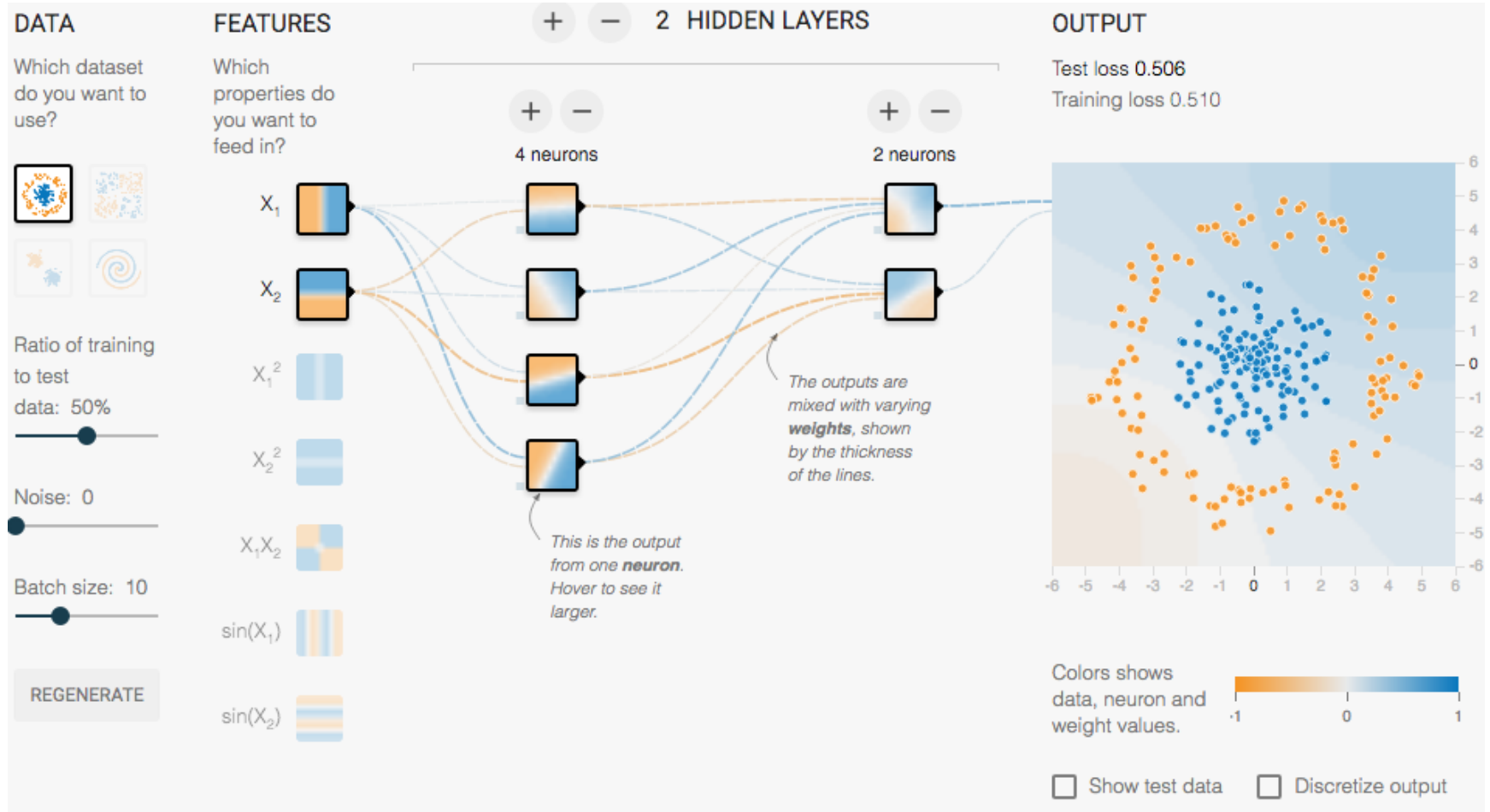
Customers

|   |     |    |                     |   |    |      |     |     |
|---|-----|----|---------------------|---|----|------|-----|-----|
| A | 22  | 1A | <a href="#">a@a</a> | 1 | aa | a1.a | 123 | aa1 |
| B | 33  | 2B | <a href="#">b@b</a> | 2 | bb | b2.b | 234 | bb2 |
| C | 44  | 3C | <a href="#">c@c</a> | 3 | cc | c3.c | 345 | cc3 |
| D | 55  | 4D | <a href="#">d@d</a> | 4 | dd | d4.d | 456 | dd4 |
| E | 66  | 5E | <a href="#">e@e</a> | 5 | ee | e5.e | 567 | ee5 |
| F | 77  | 6F | <a href="#">f@f</a> | 6 | ff | f6.f | 678 | ff6 |
| G | 88  | 7G | <a href="#">g@g</a> | 7 | gg | g7.g | 789 | gg7 |
| H | 99  | 8H | <a href="#">h@h</a> | 8 | hh | h8.h | 890 | hh8 |
| I | 111 | 9I | <a href="#">i@i</a> | 9 | ii | i9.i | 901 | ii9 |

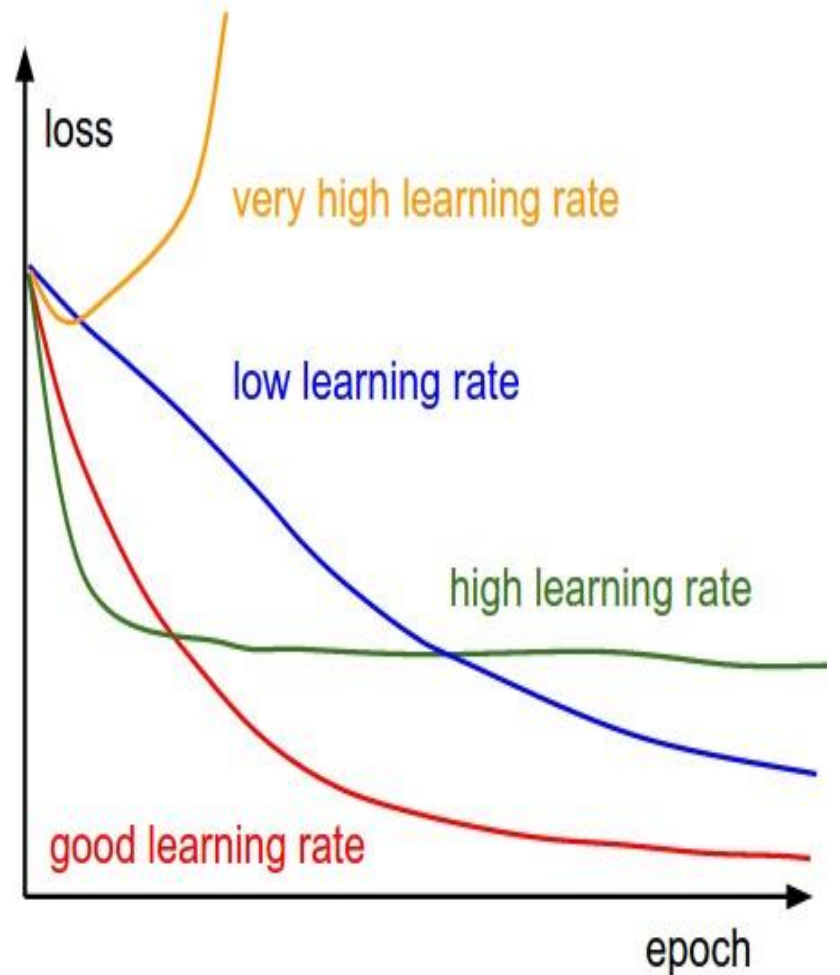
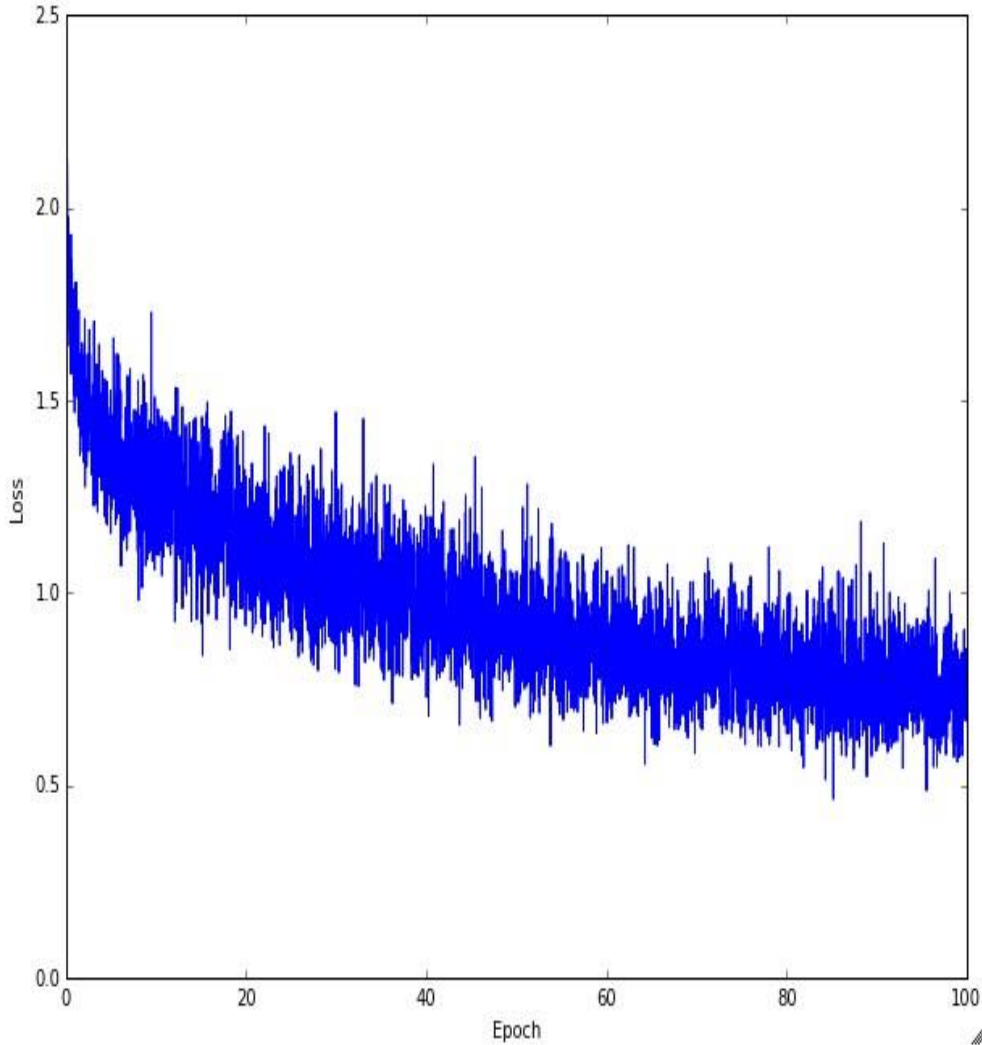


# Demo

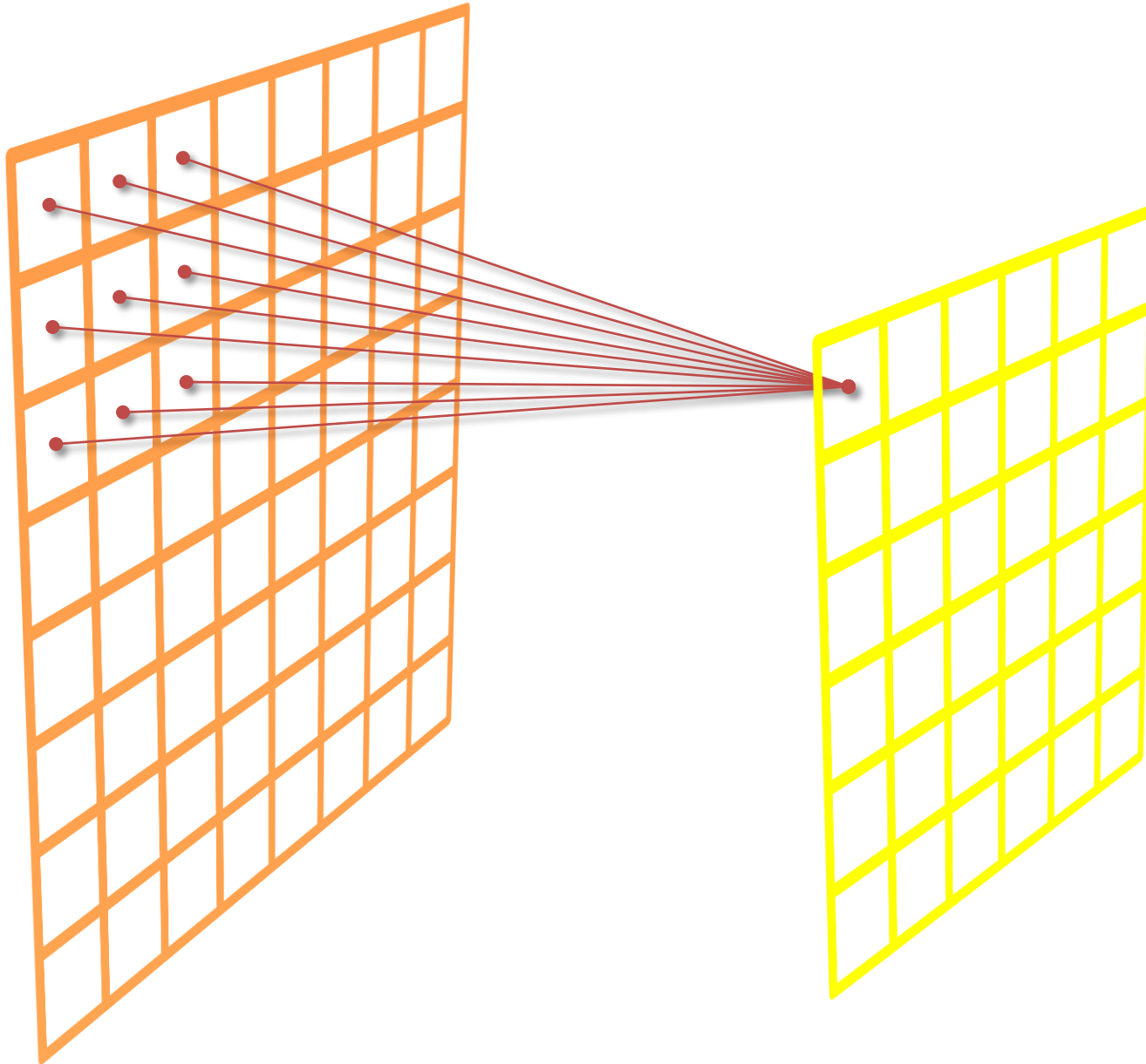
<http://playground.tensorflow.org/>



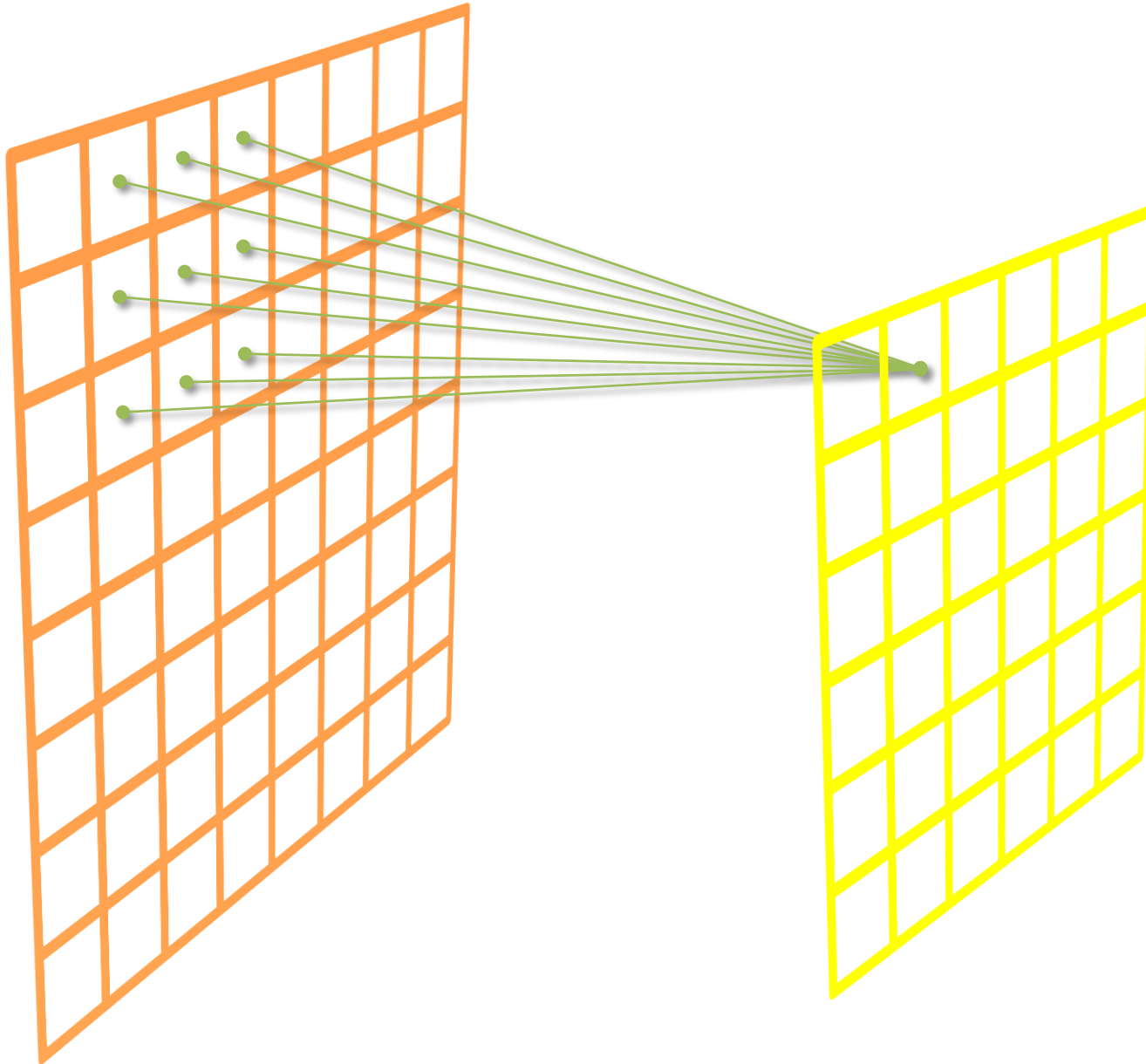
# Monitor and visualize the loss curve



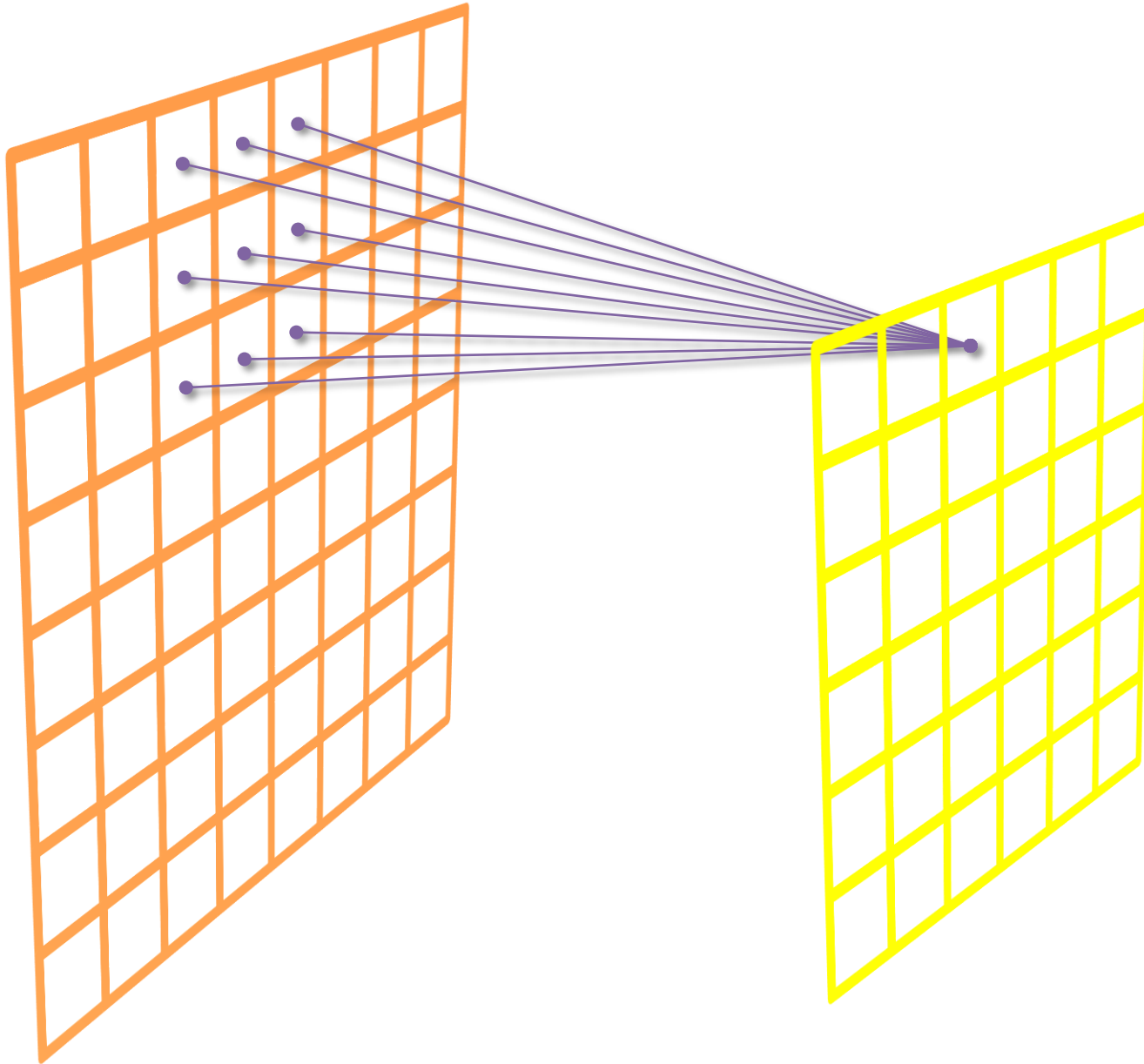
8x8 image, 3x3 filter, Stride 1



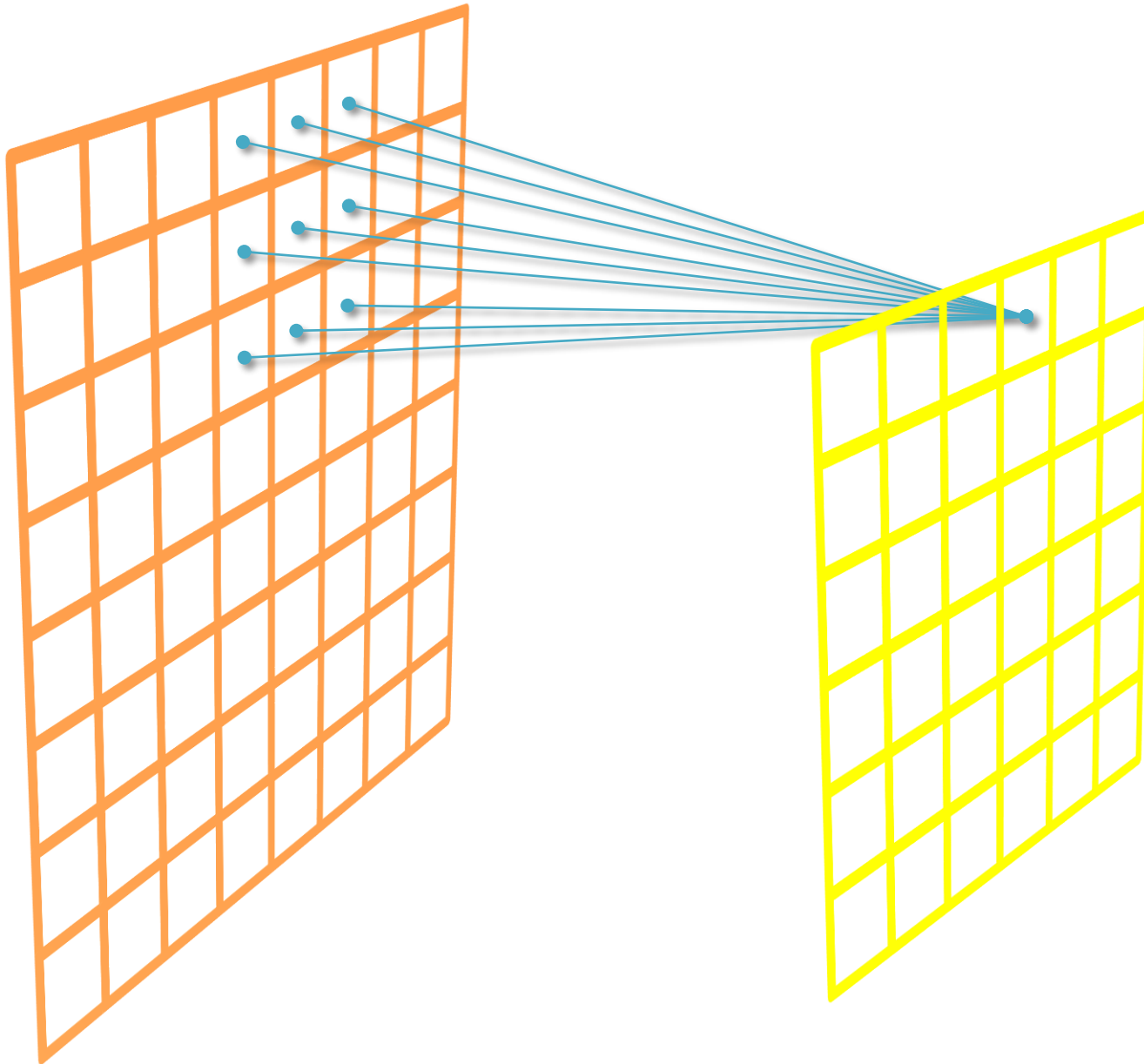
8x8 image, 3x3 filter, Stride 1



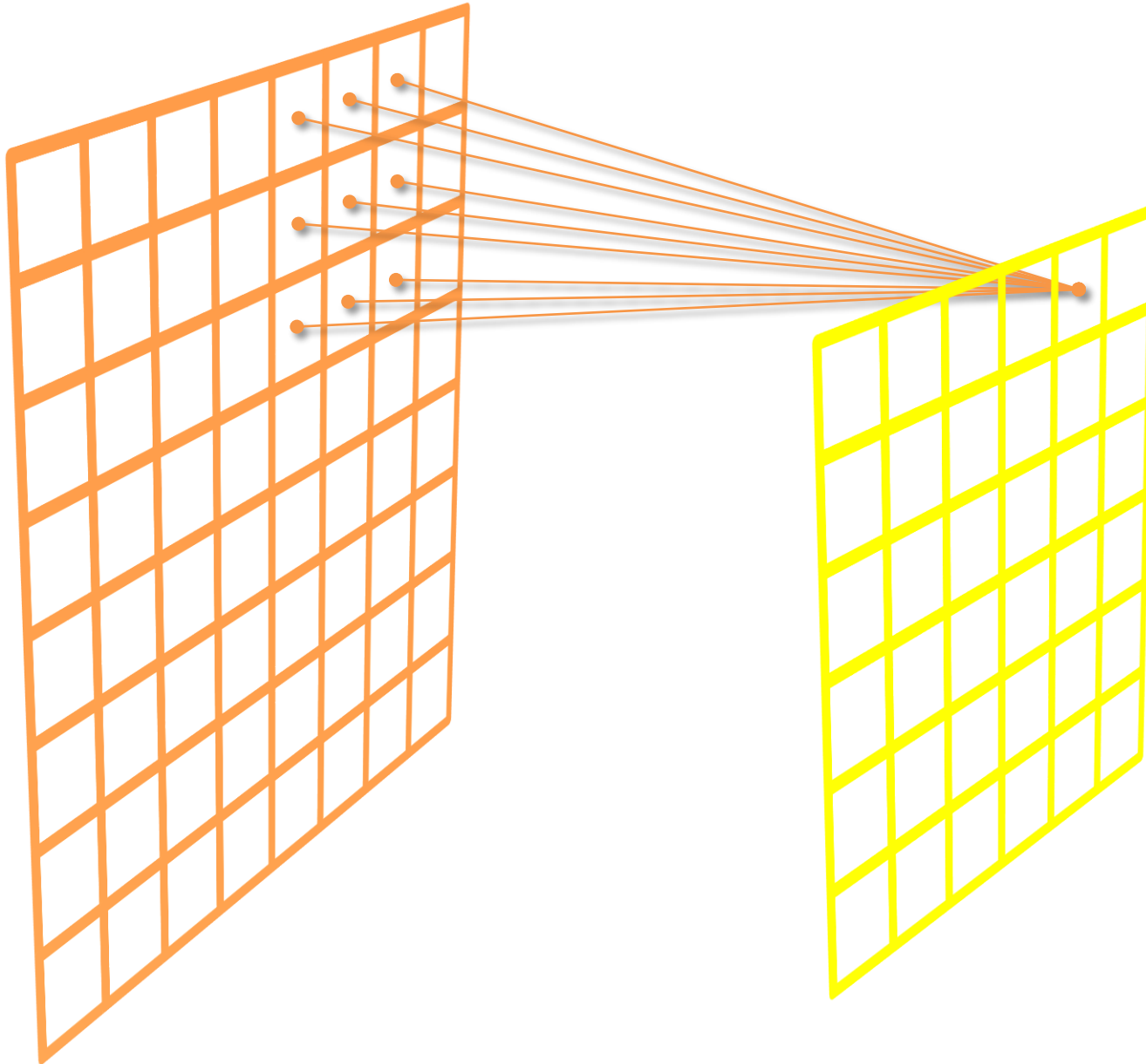
8x8 image, 3x3 filter, Stride 1



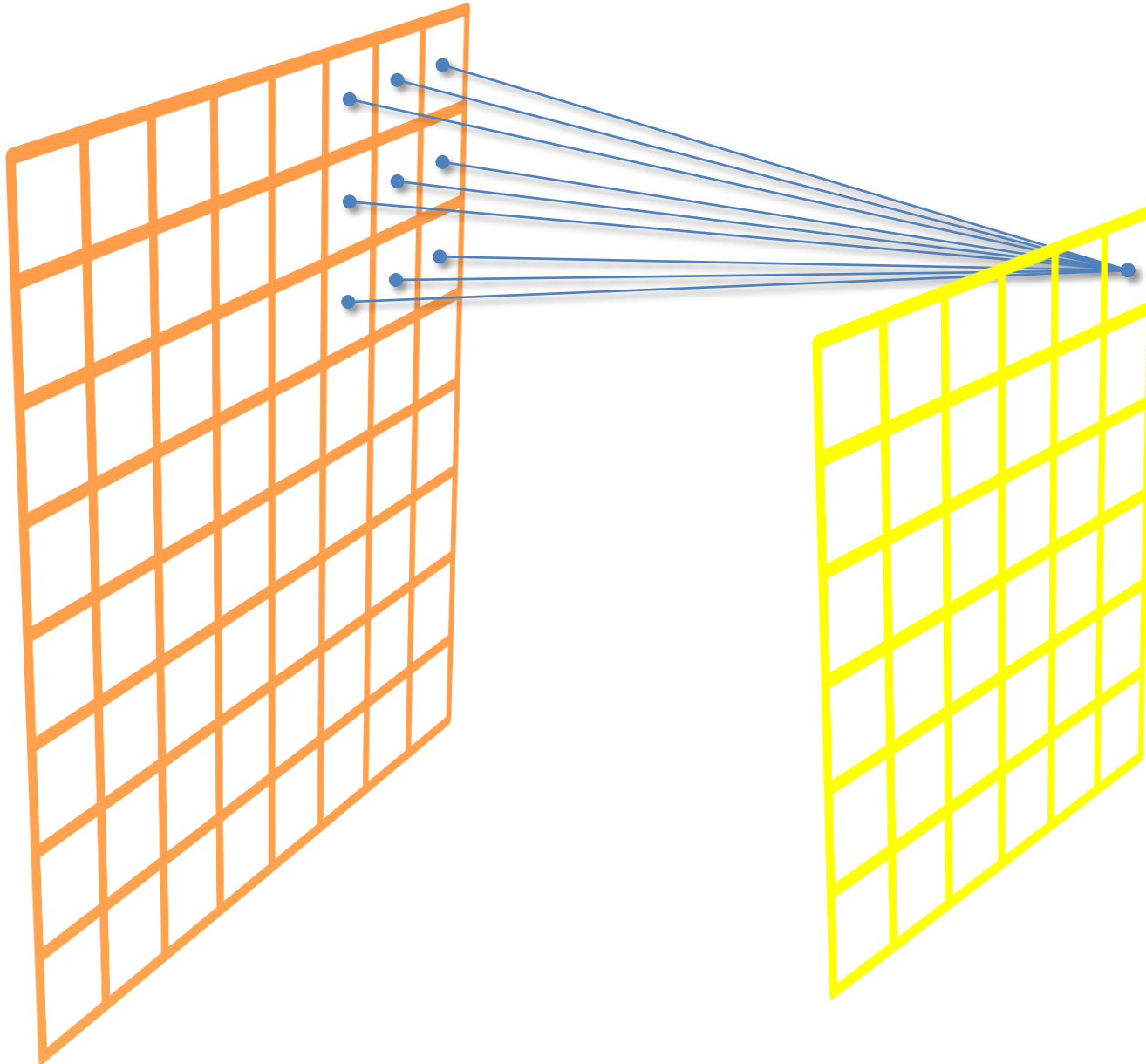
8x8 image, 3x3 filter, Stride 1



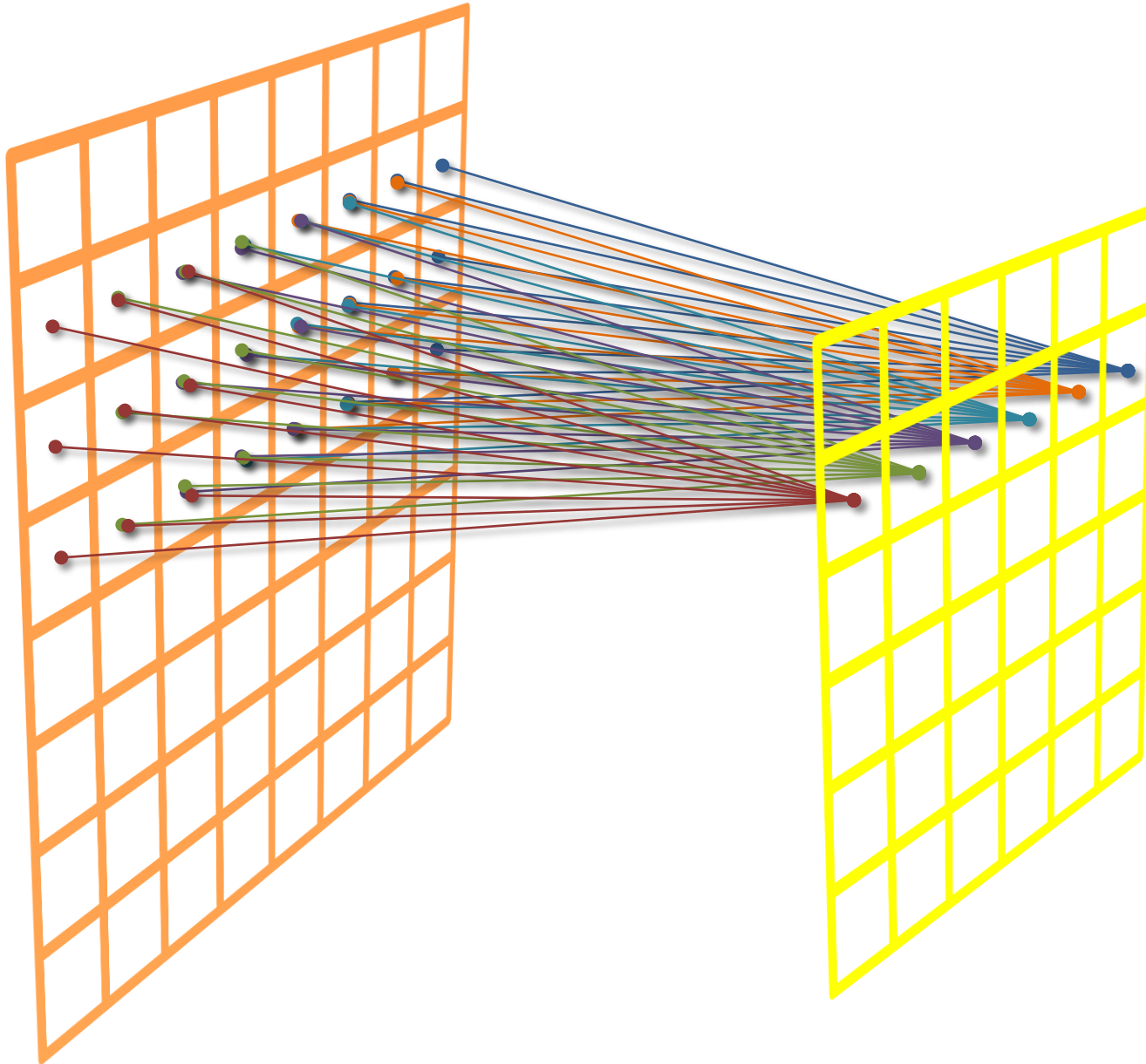
8x8 image, 3x3 filter, Stride 1



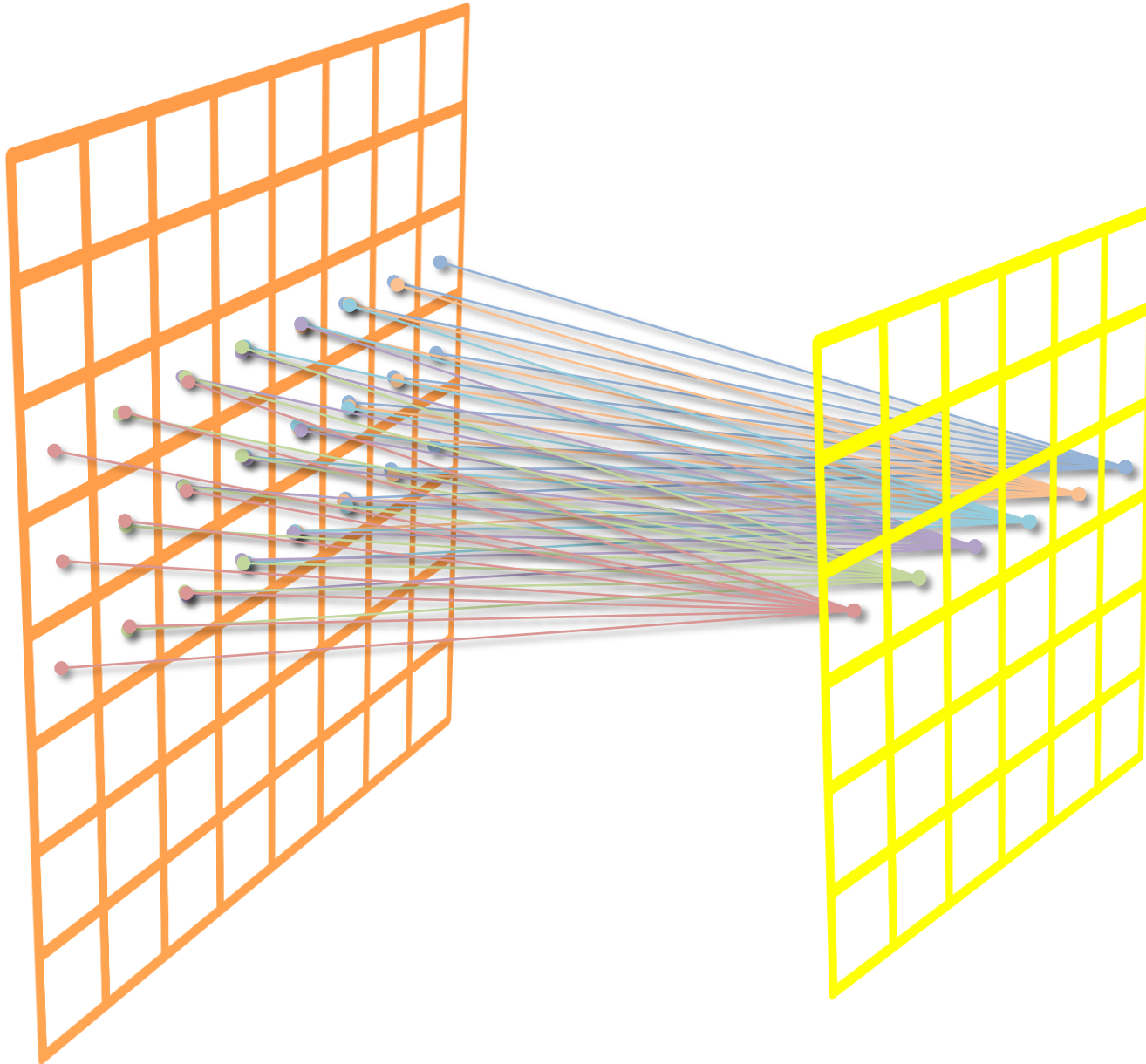
8x8 image, 3x3 filter, Stride 1

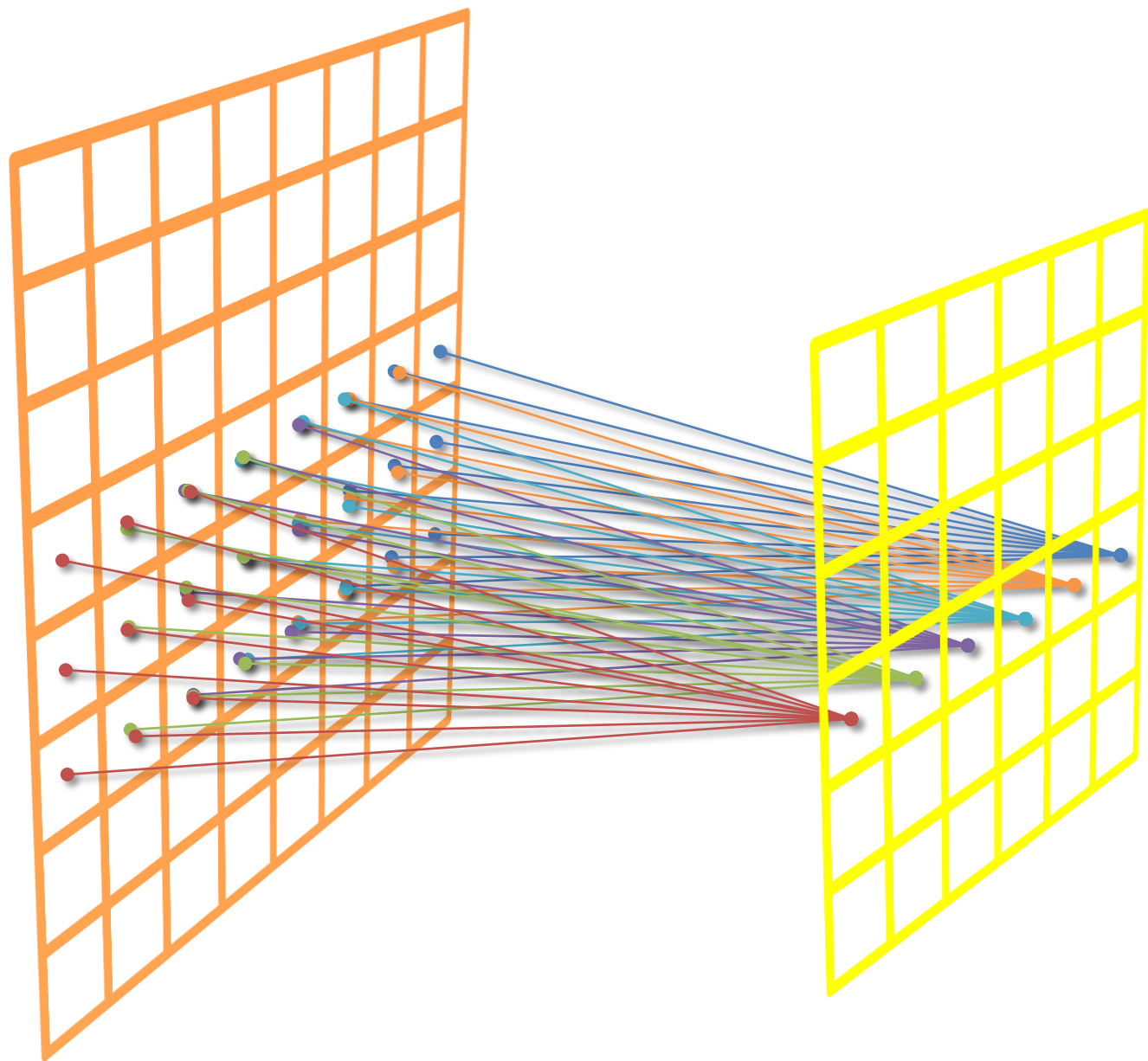


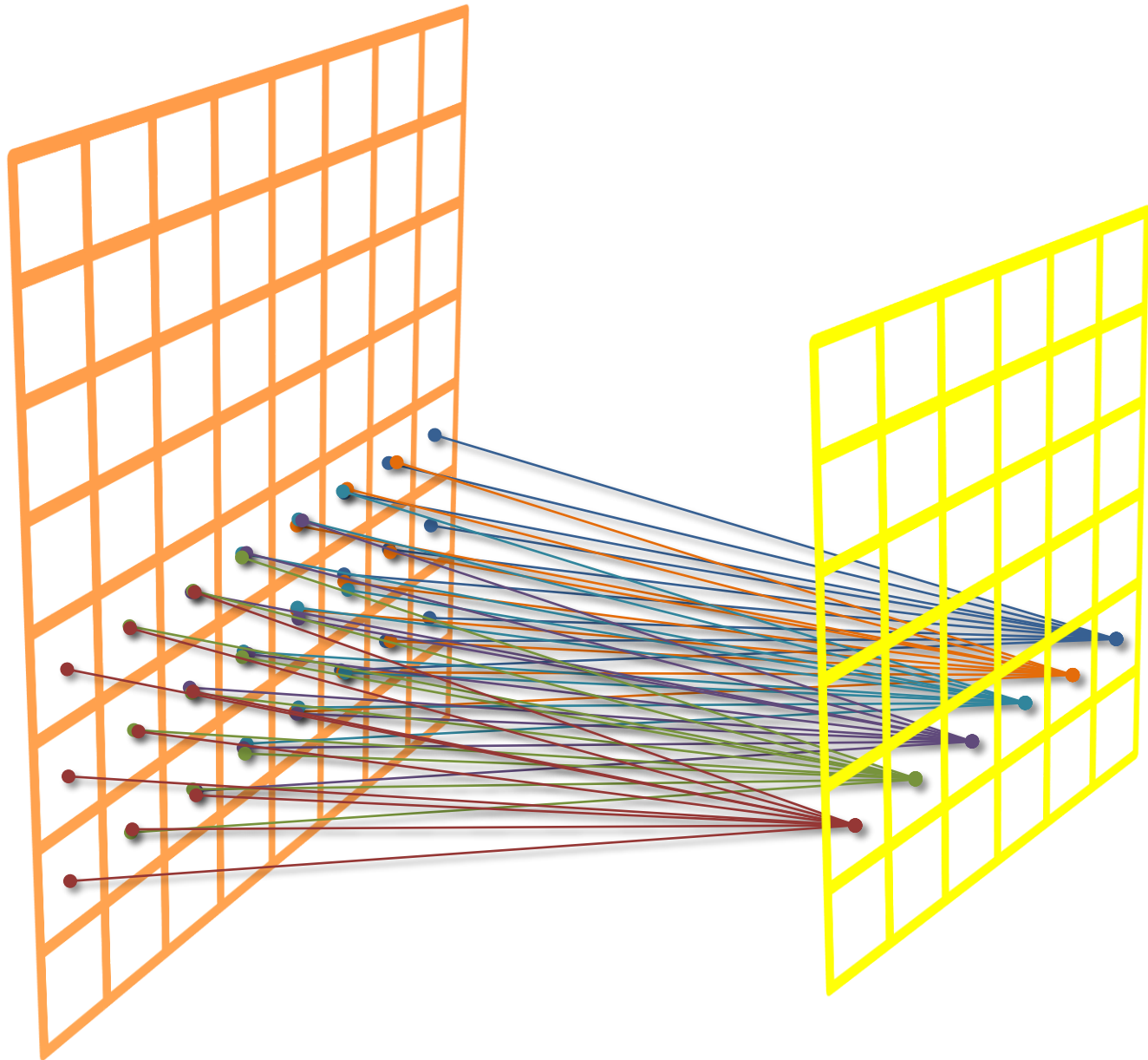
8x8 image, 3x3 filter, Stride 1



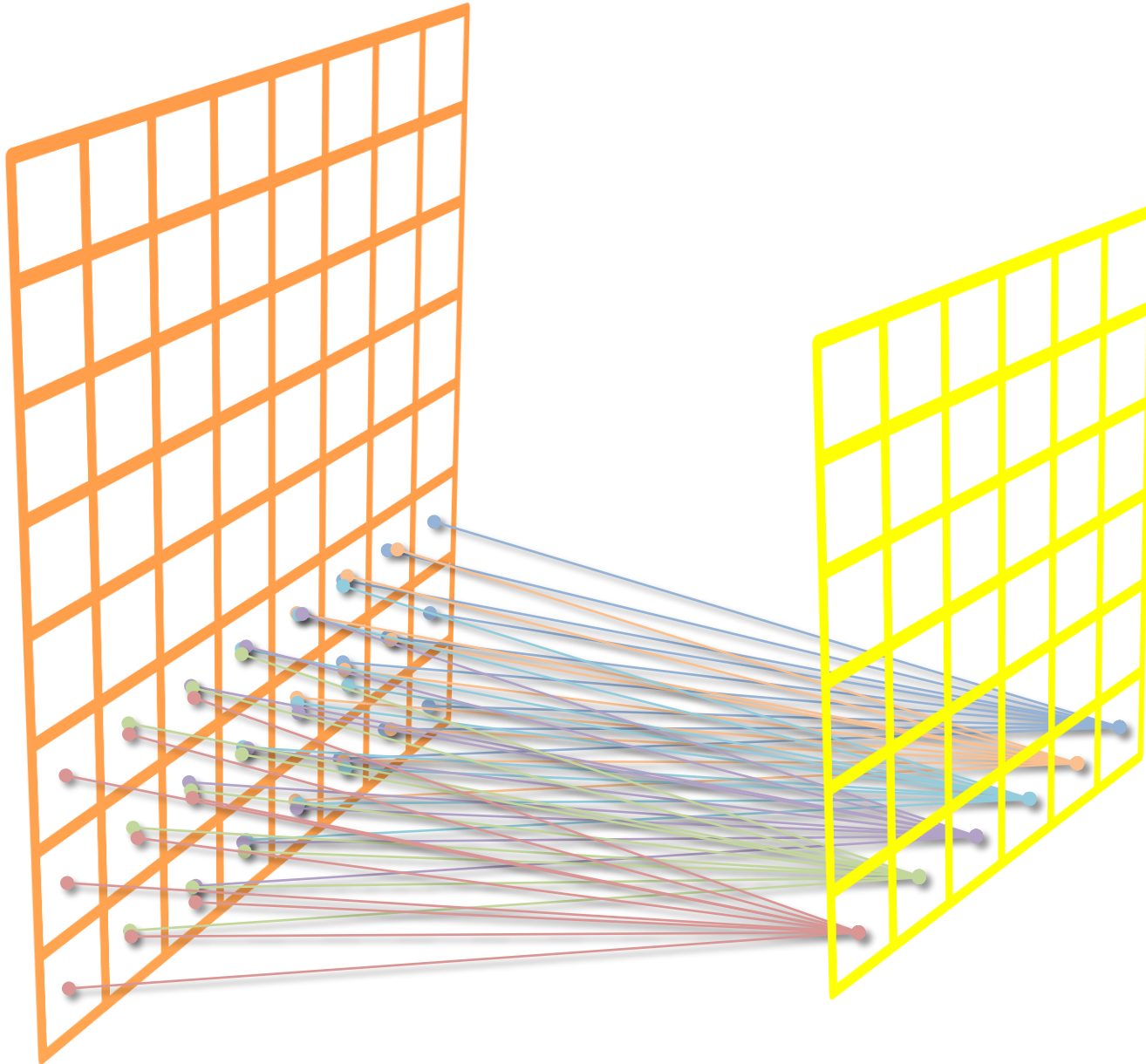
8x8 image, 3x3 filter, Stride 1



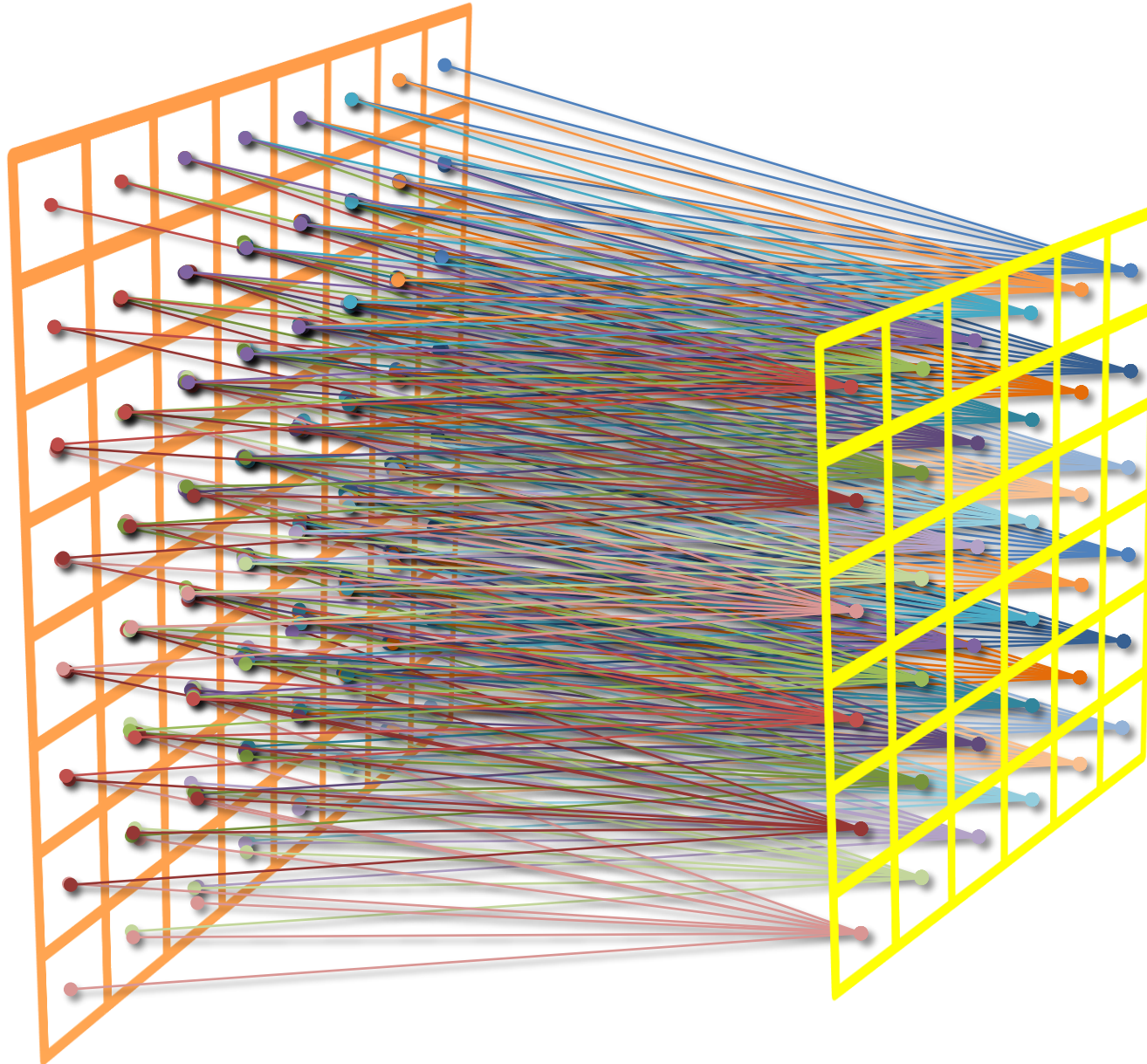




8x8 image, 3x3 filter, Stride 1



8x8 image, 3x3 filter, Stride 1



Total Connections?

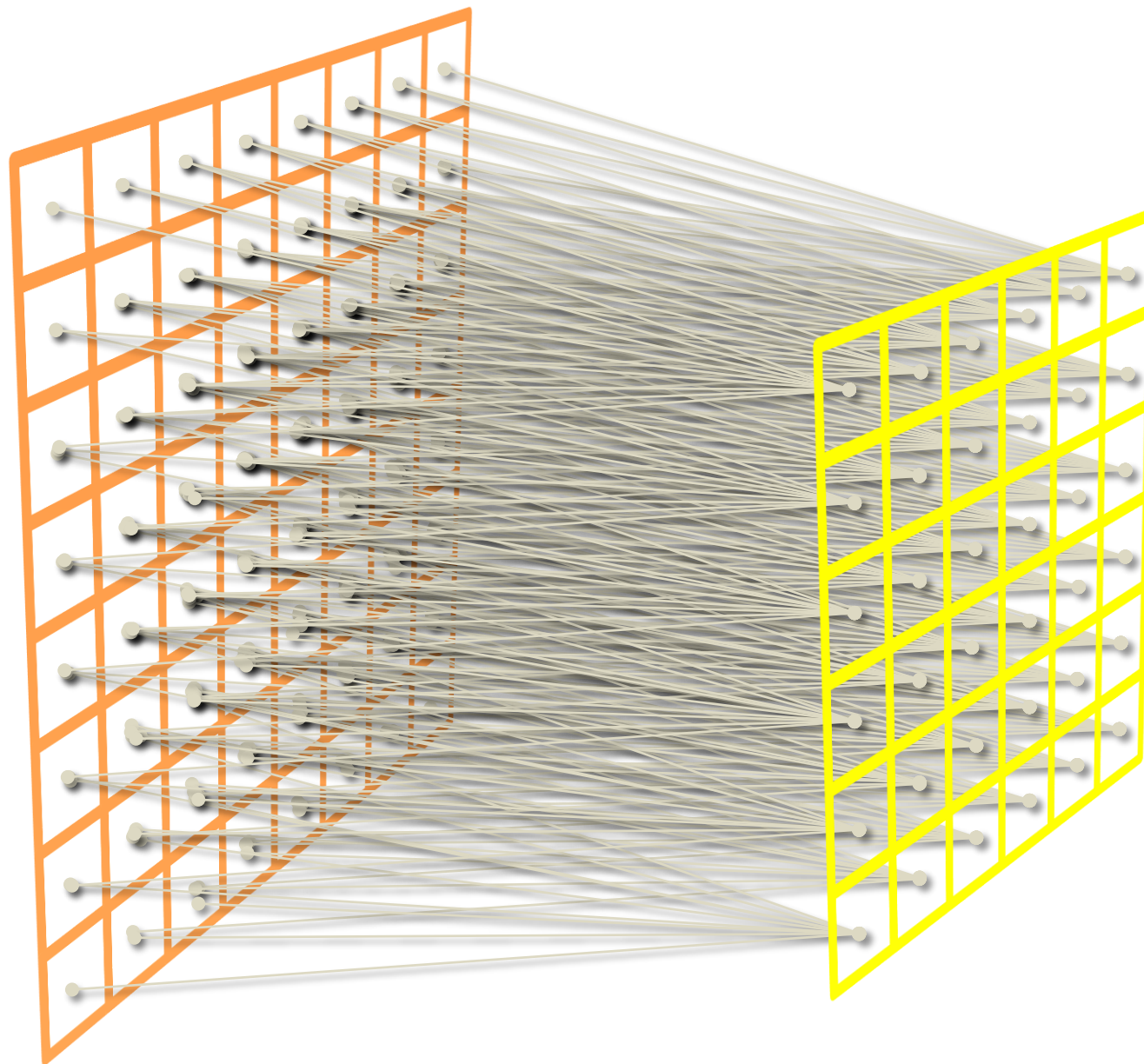
6x6x1 x 3x3x1

What would have been the number of connections if this had been a fully-connected layer?

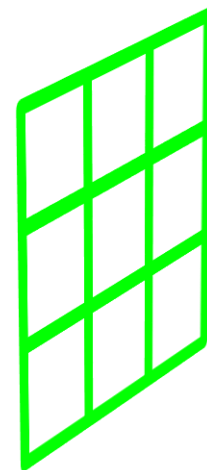
Total Unique Parameters?

3x3x1  
+ bias

8x8 image, 3x3 filter, Stride 1

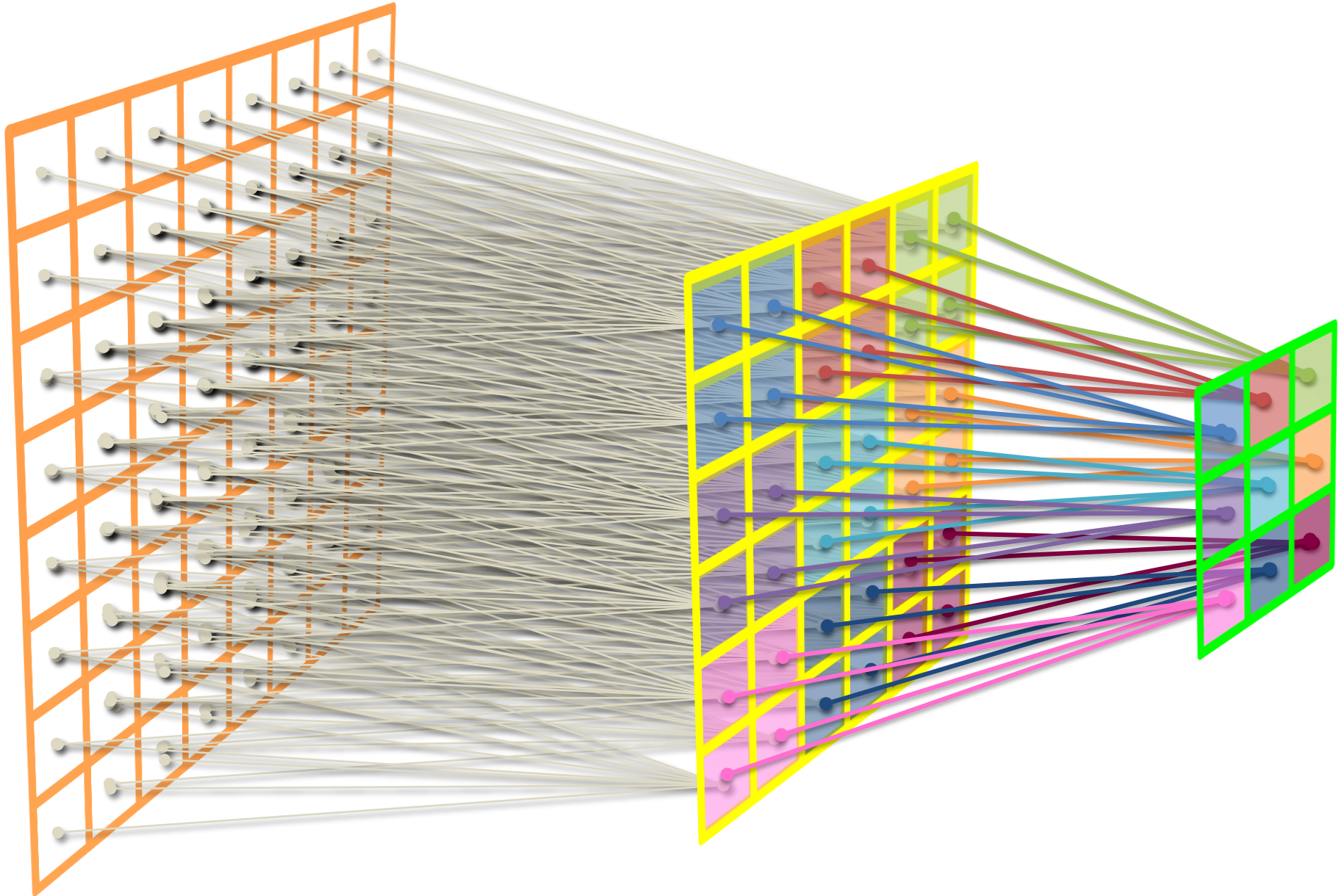


2x2 pooling, Stride 2



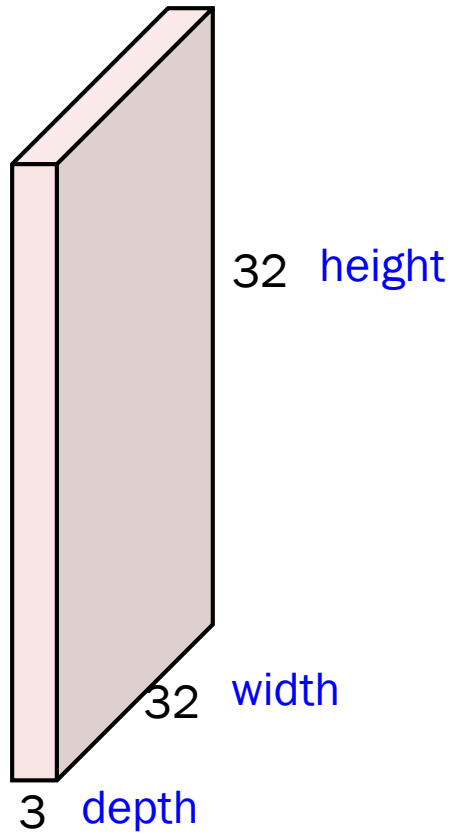
8x8 image, 3x3 filter, Stride 1

2x2 pooling, Stride 2



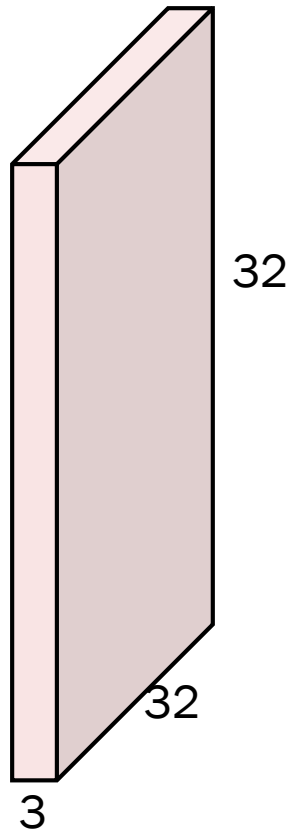
# Convolution Layer

32x32x3 image



# Convolution Layer

32x32x3 image



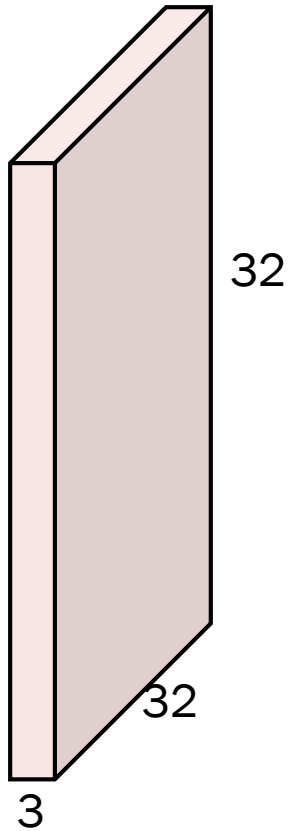
5x5x3 filter



**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# Convolution Layer

32x32x3 image



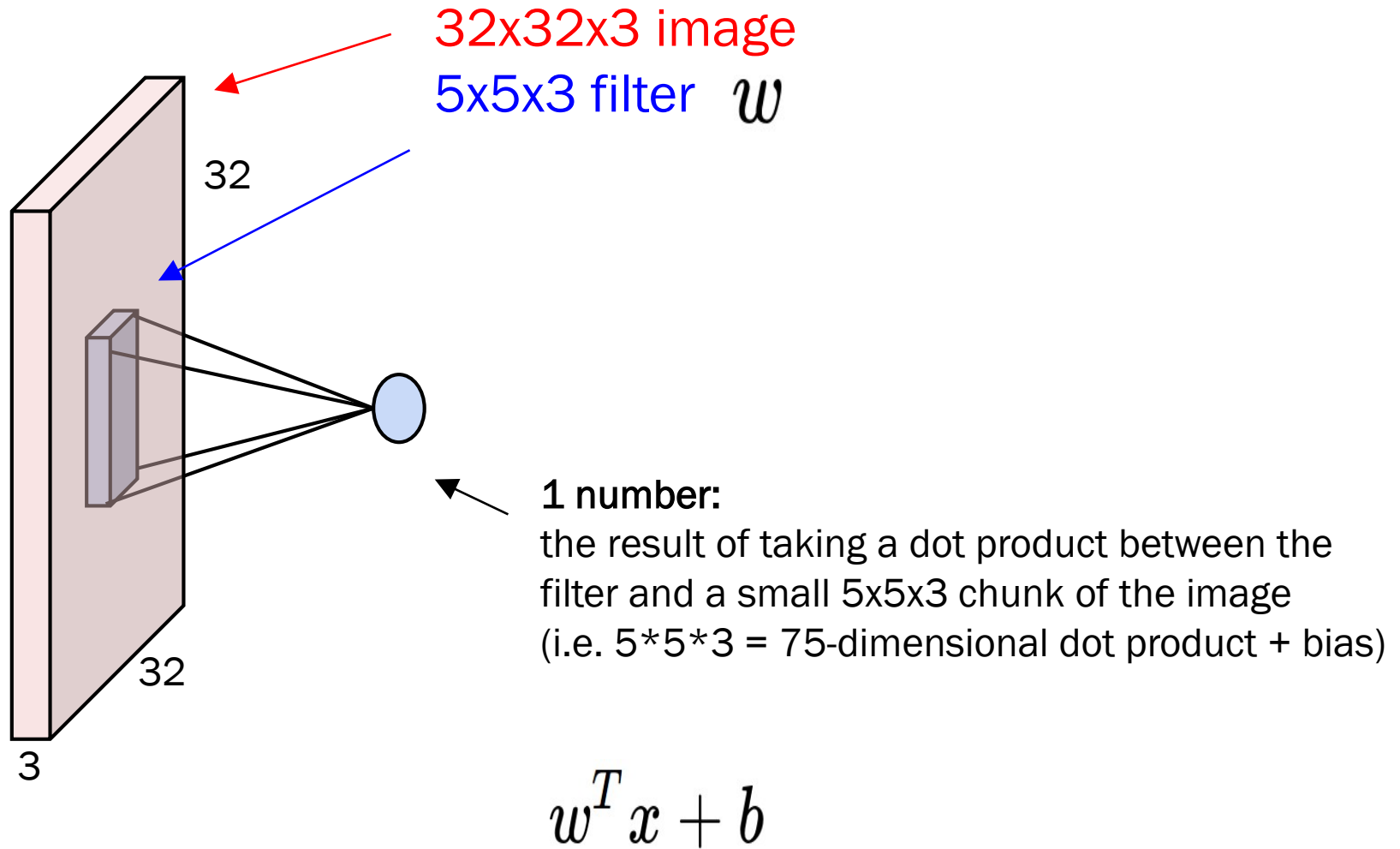
Filters always extend the full depth of the input volume

5x5x3 filter

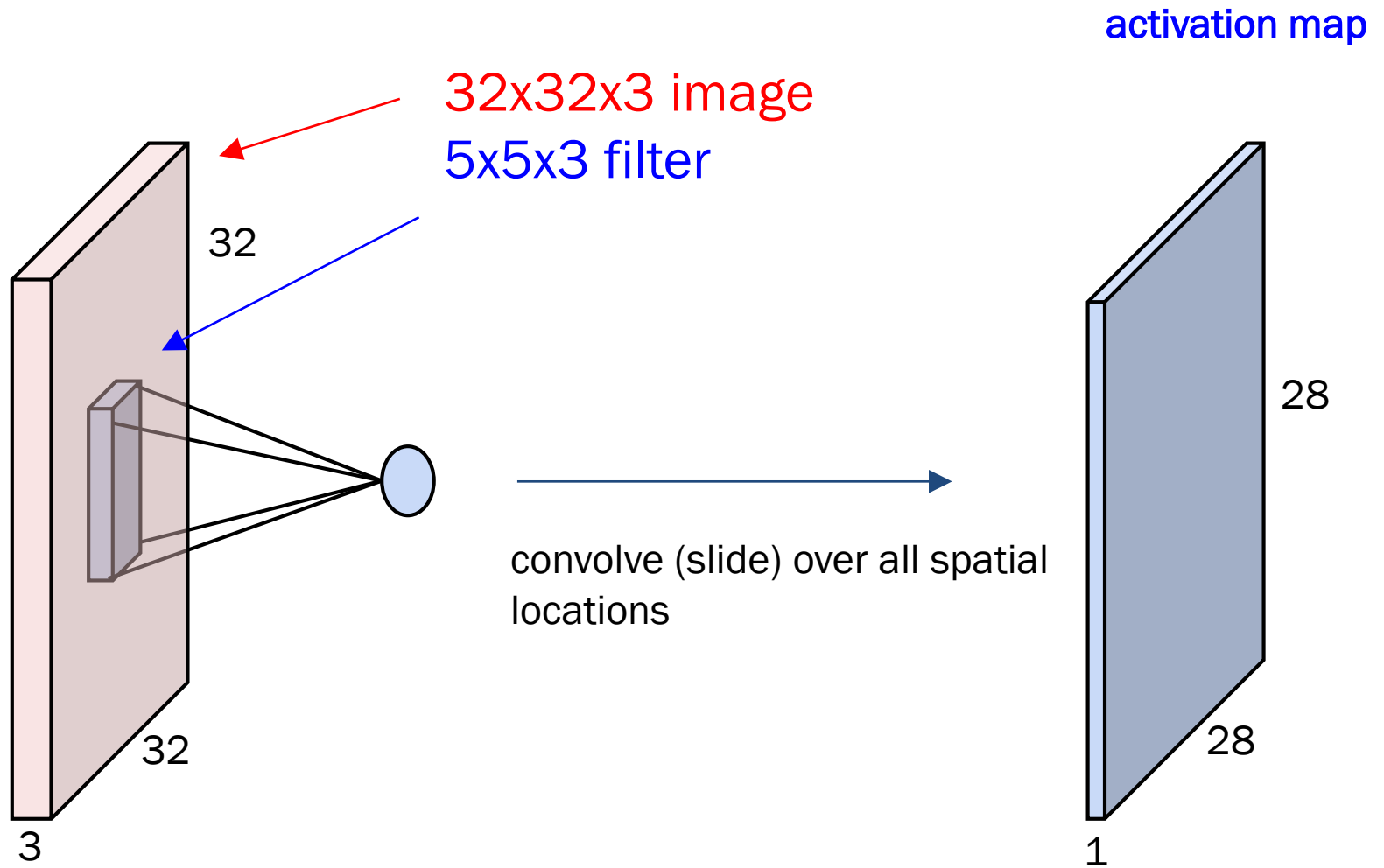


**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# Convolution Layer

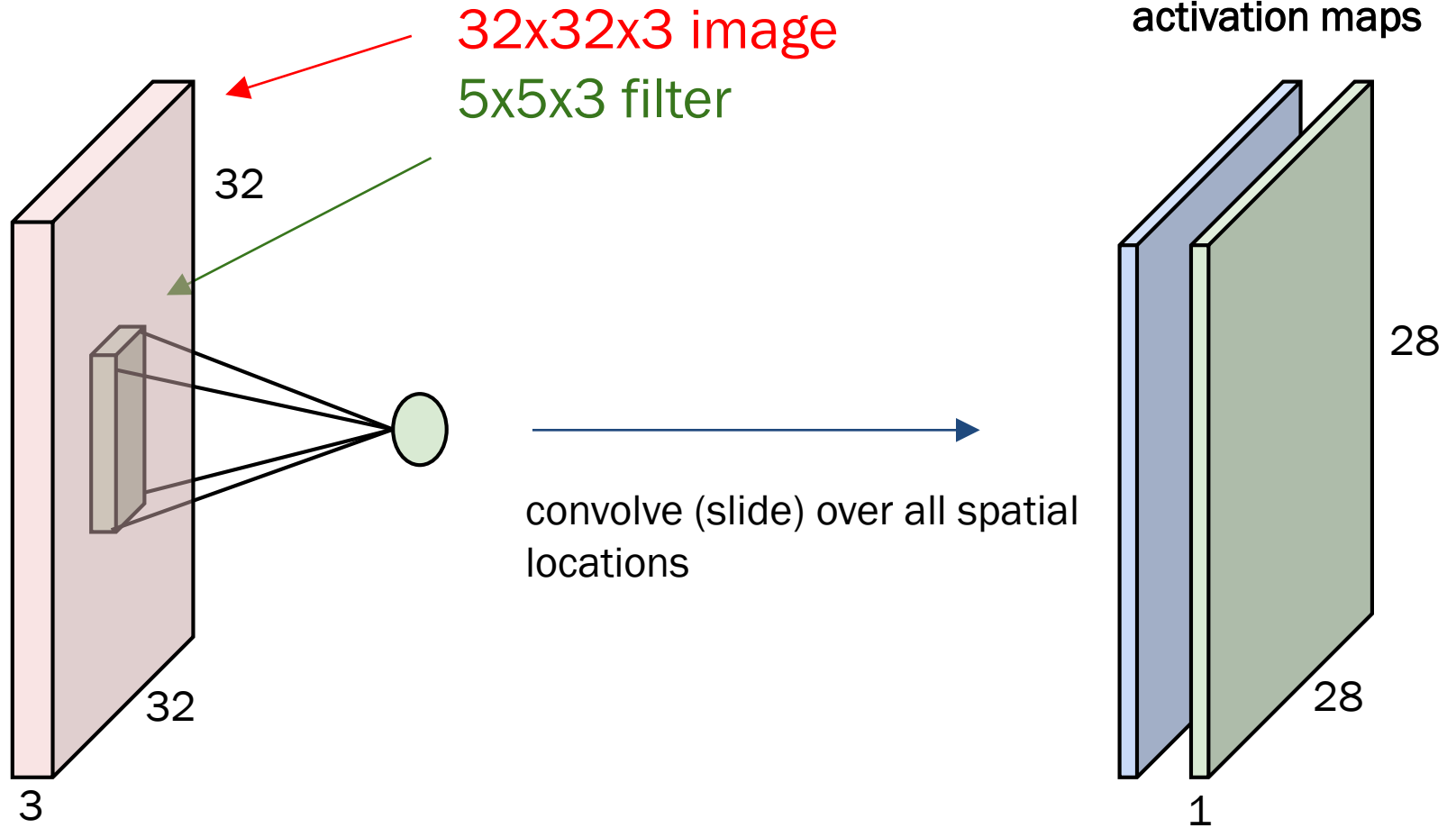


# Convolution Layer

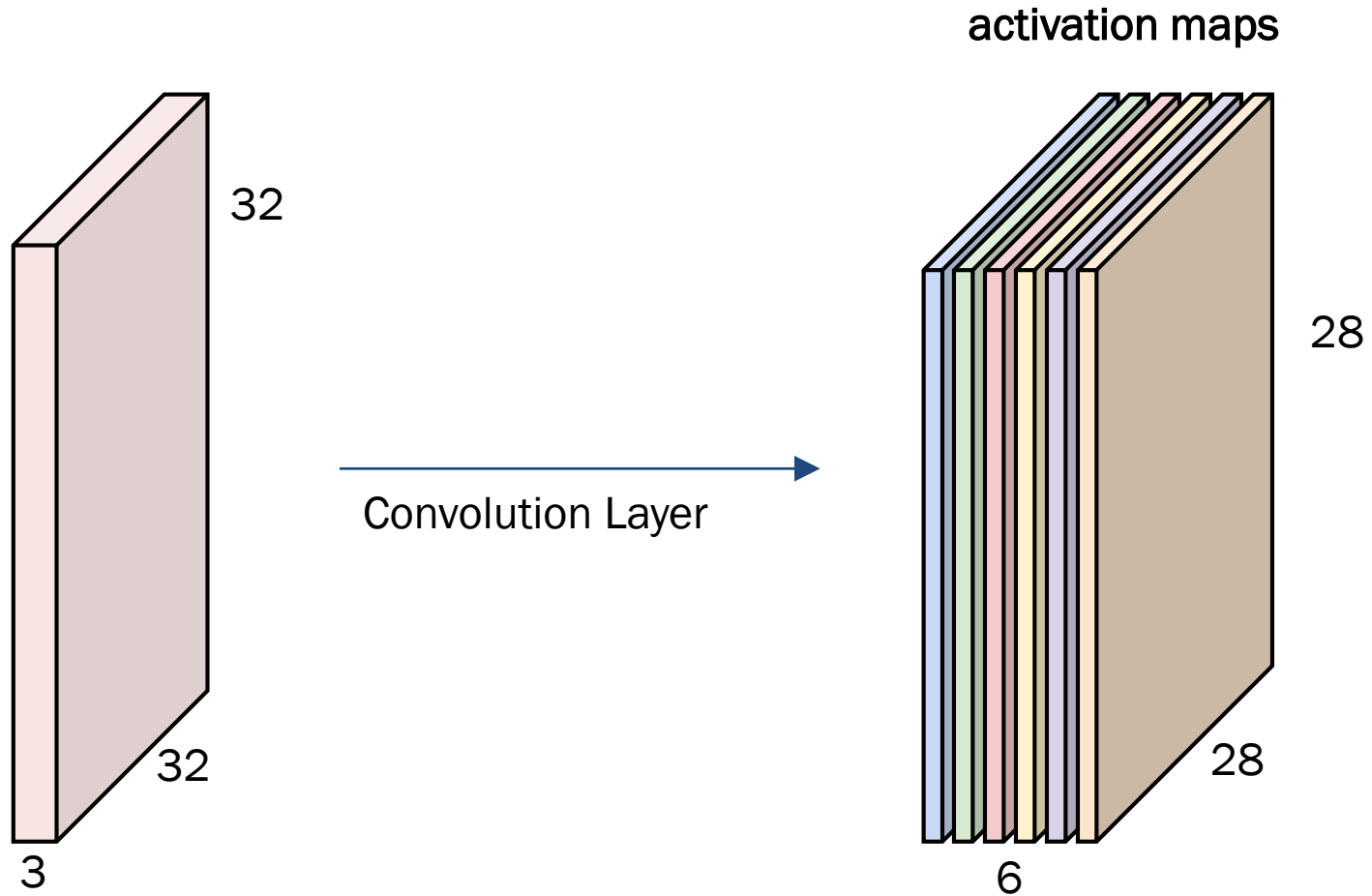


# Convolution Layer

consider a second, green filter

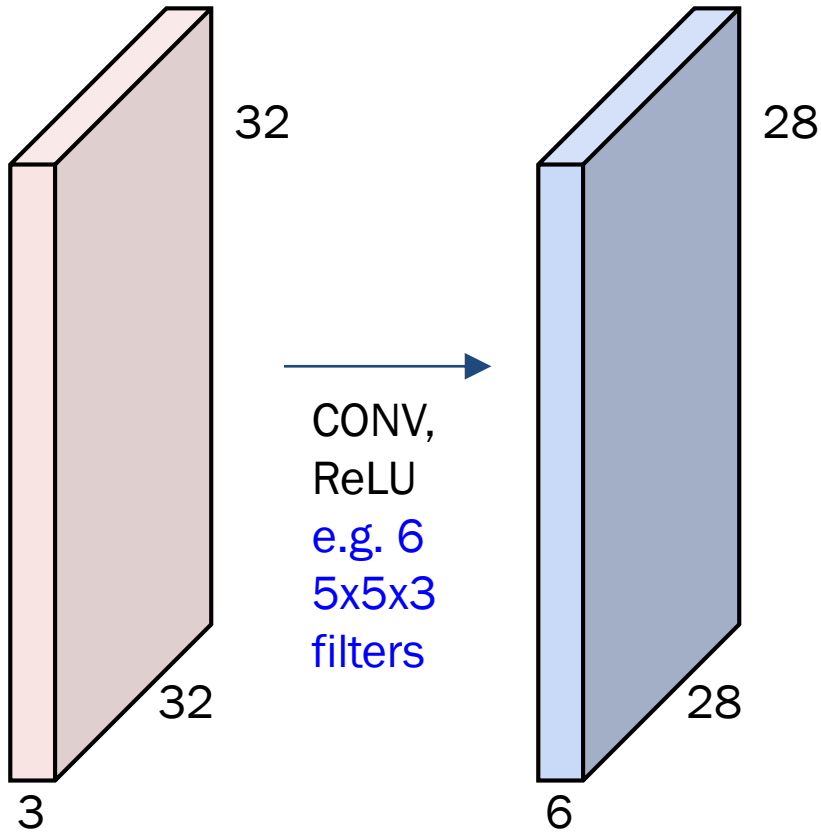


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

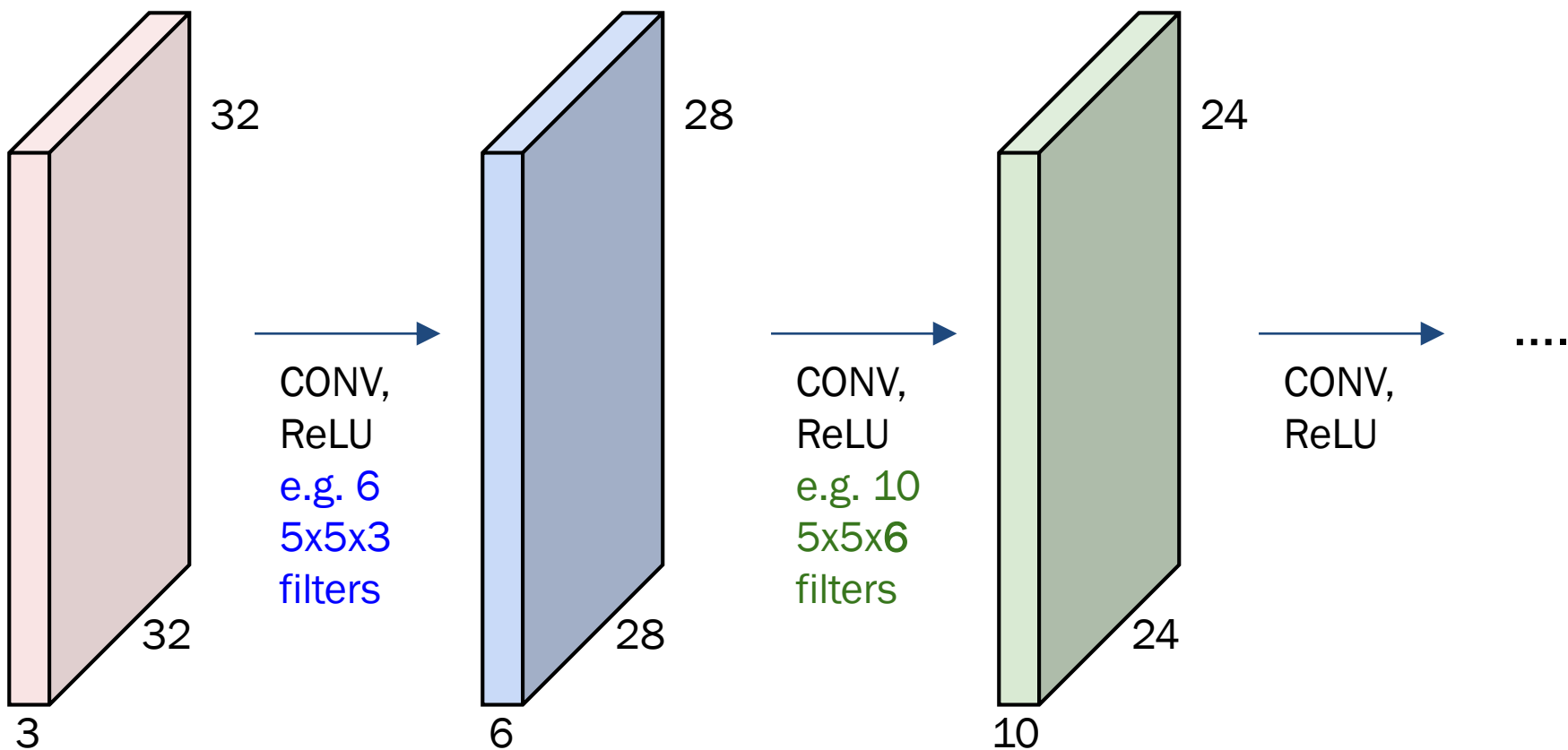


We stack these up to get a “new image” of size 28x28x6!

**Preview:** ConvNet is a sequence of Convolution Layers, interspersed with activation functions

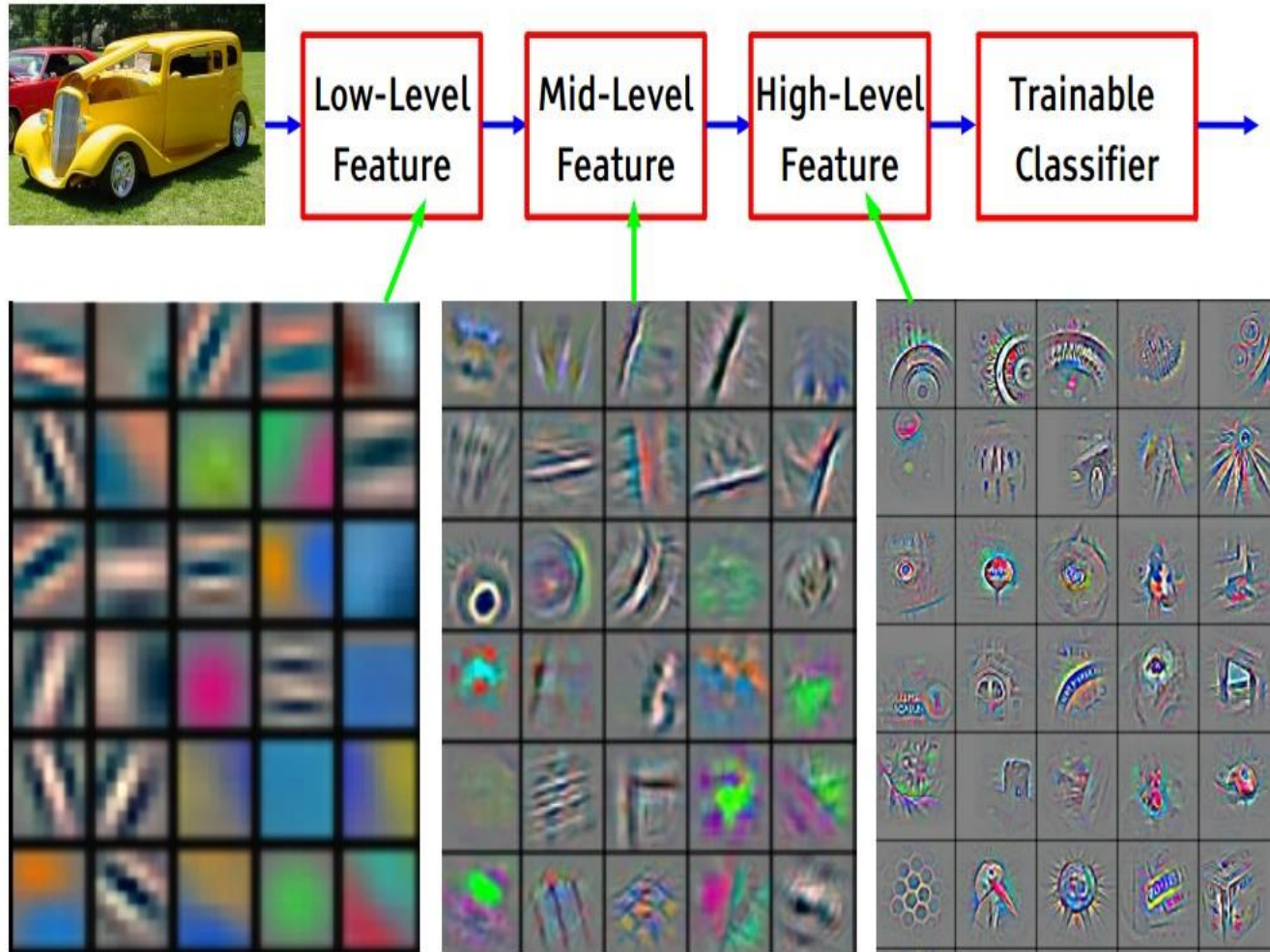


**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



# Preview

[From recent Yann LeCun slides]



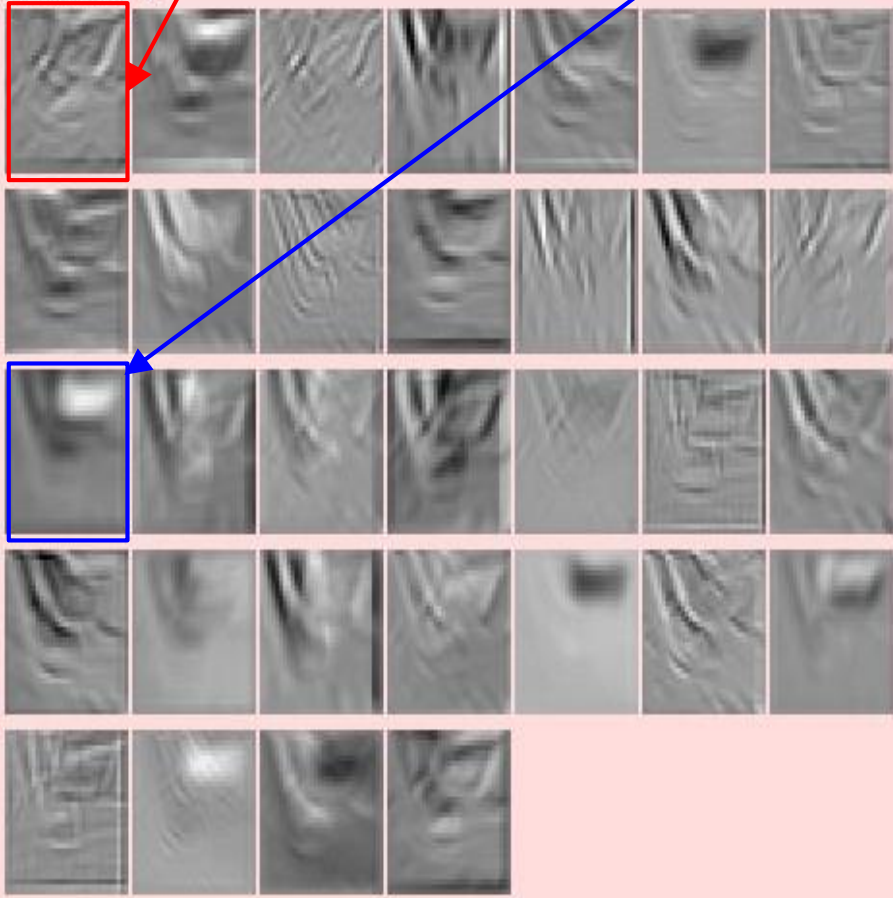
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]



one filter =>  
one activation map

example 5x5 filters  
(32 total)

Activations:

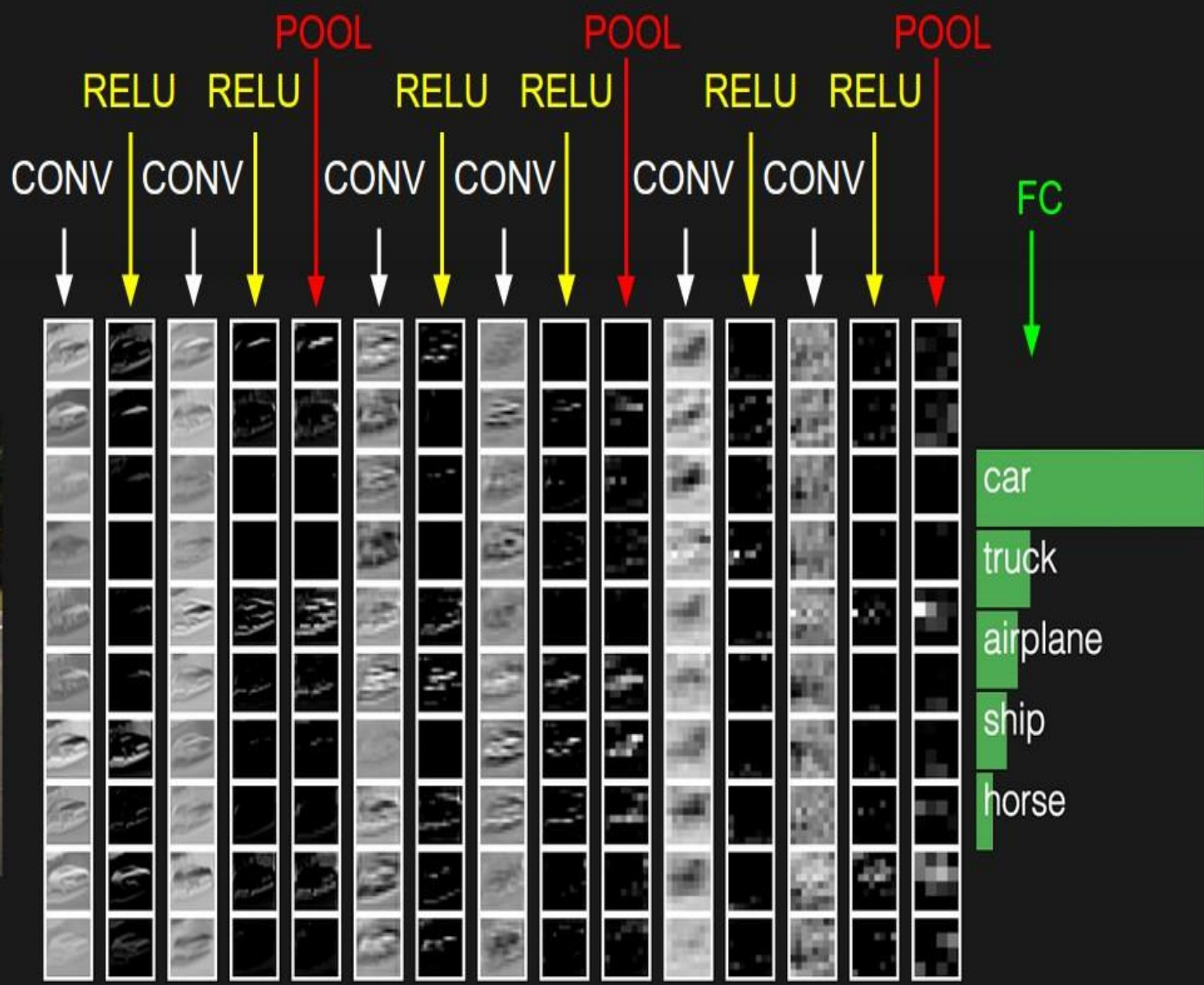


We call the layer convolutional because it is related to convolution of two signals:

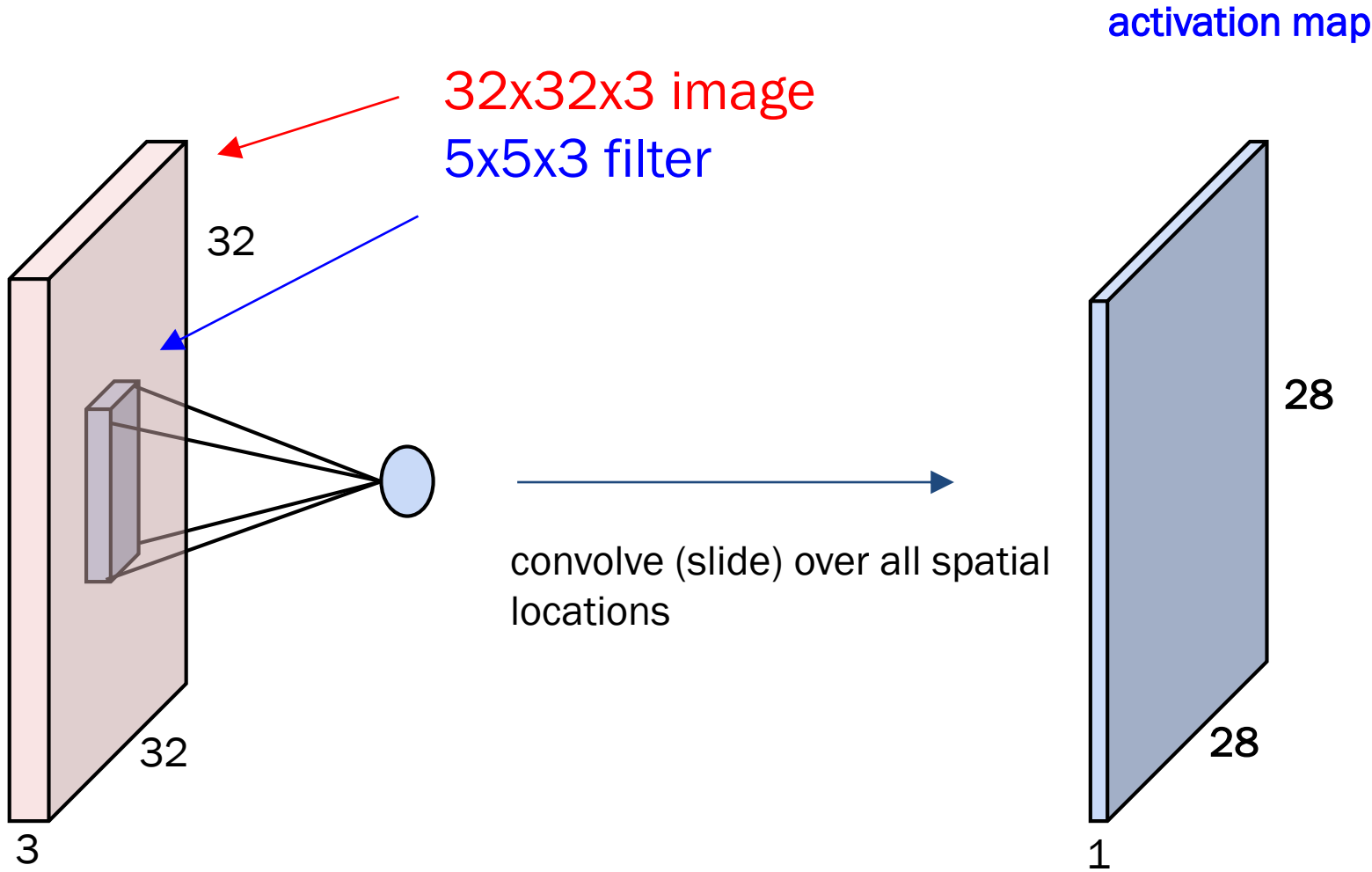
$$f[x,y] * g[x,y] = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1,n_2] \cdot g[x-n_1,y-n_2]$$

↑  
elementwise multiplication and  
sum of a filter and the signal  
(image)

preview:

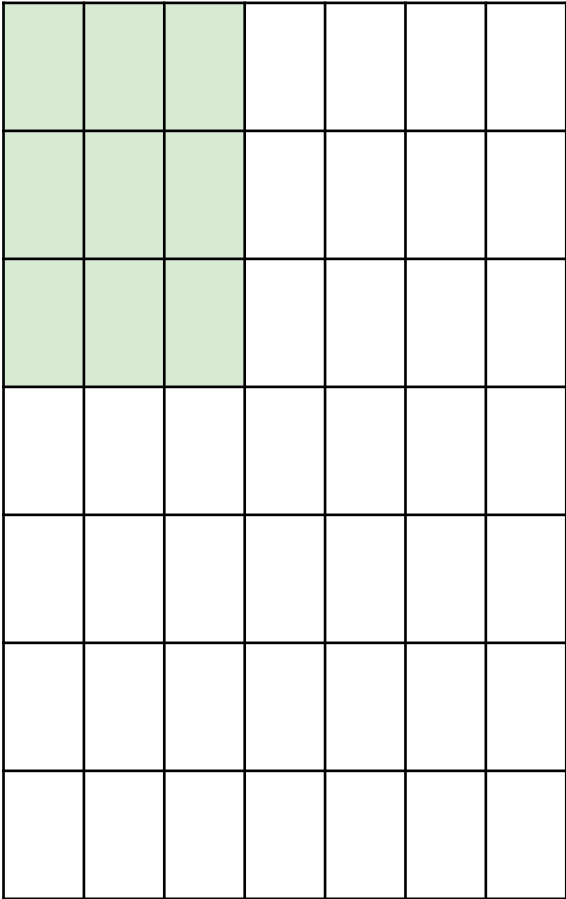


A closer look at spatial dimensions:



A closer look at spatial dimensions:

7

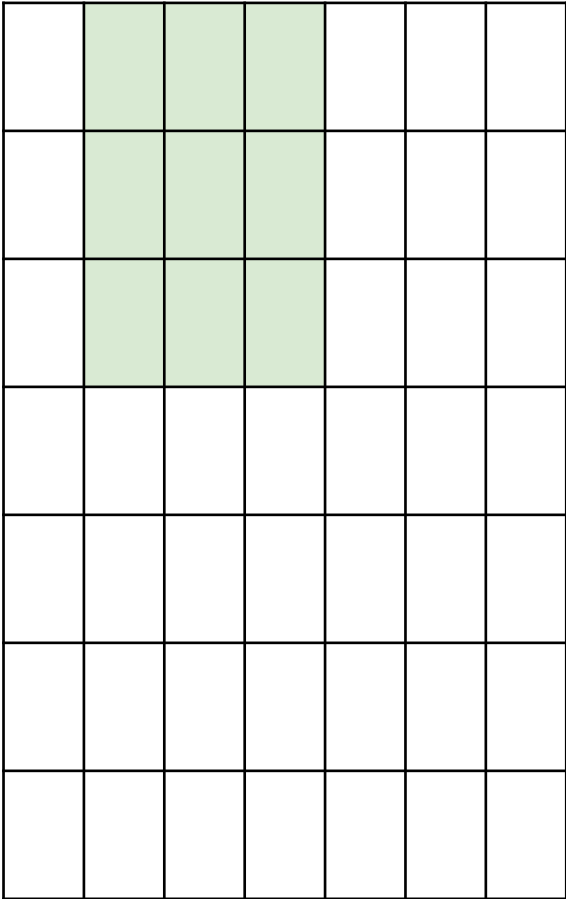


7x7 input (spatially)  
assume 3x3 filter

7

A closer look at spatial dimensions:

7

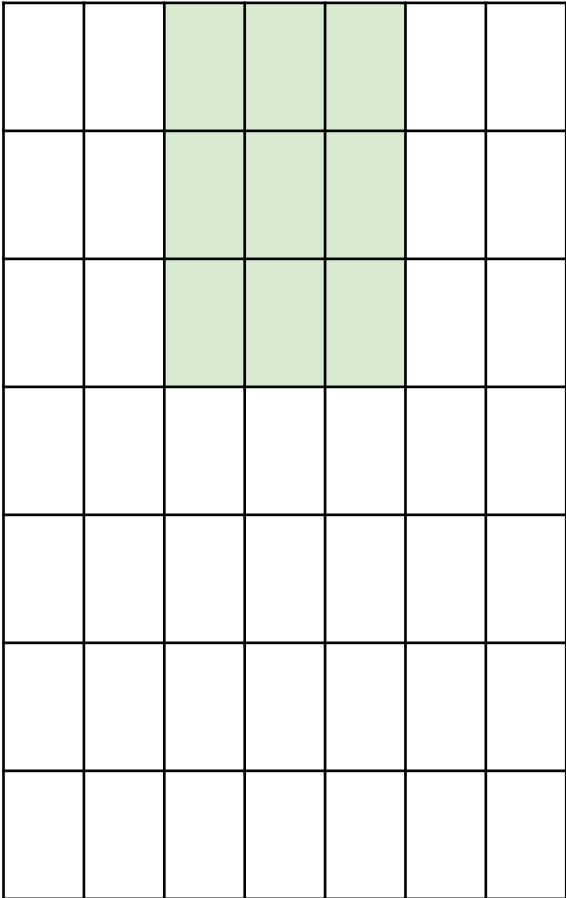


7x7 input (spatially)  
assume 3x3 filter

7

A closer look at spatial dimensions:

7

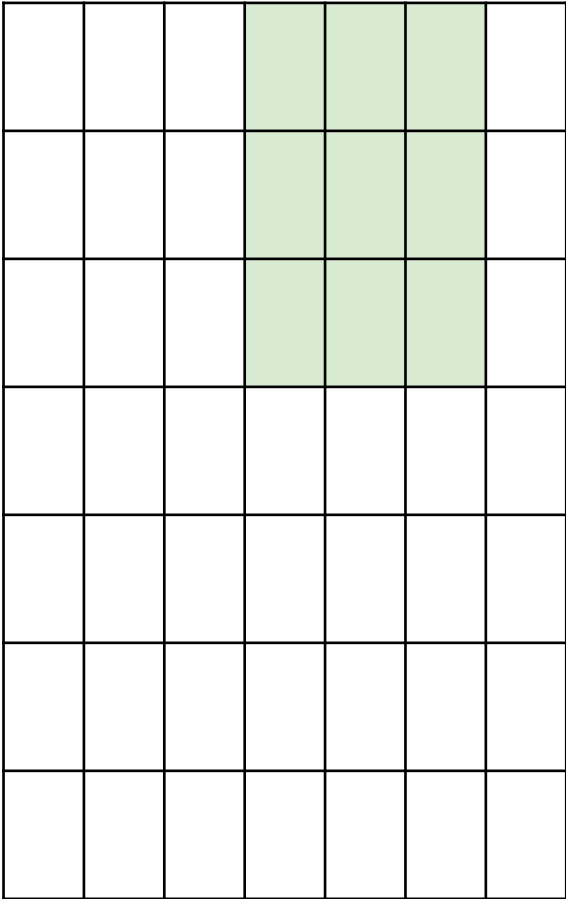


7x7 input (spatially)  
assume 3x3 filter

7

A closer look at spatial dimensions:

7

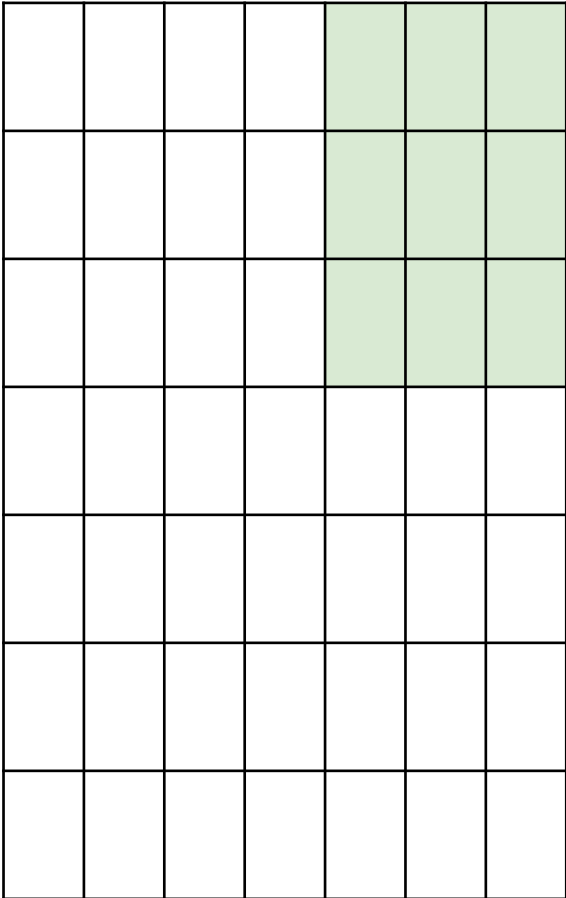


7x7 input (spatially)  
assume 3x3 filter

7

A closer look at spatial dimensions:

7



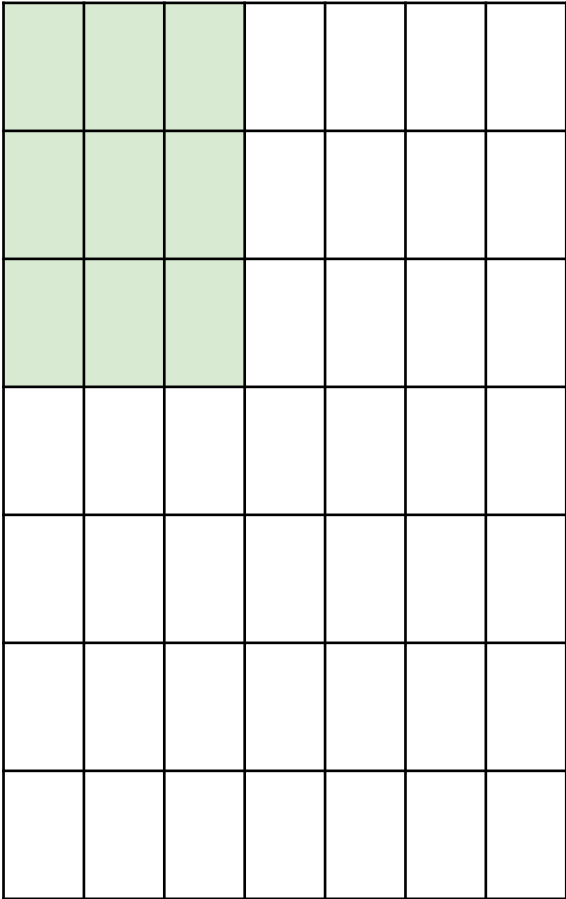
7x7 input (spatially)  
assume 3x3 filter

=> 5x5 output

7

A closer look at spatial dimensions:

7

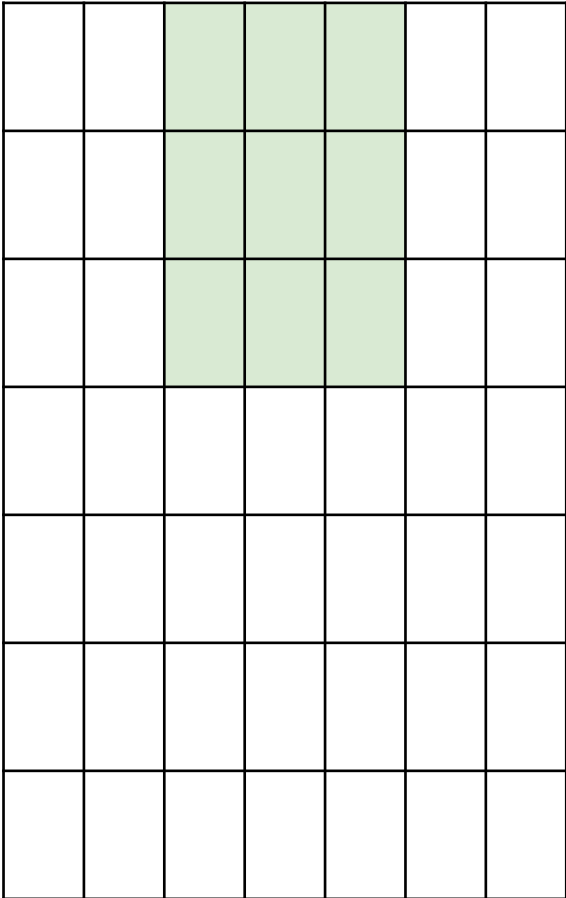


7

7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

A closer look at spatial dimensions:

7

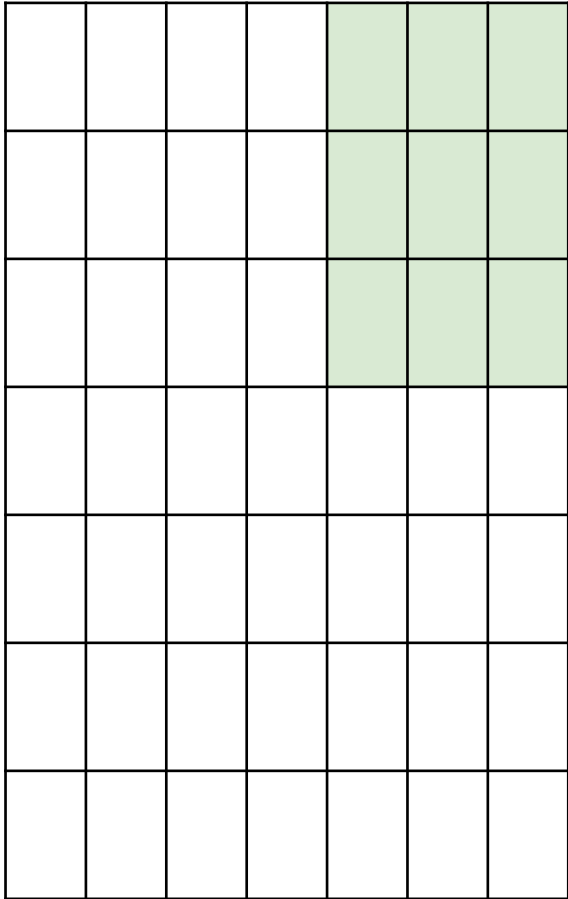


7

7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

A closer look at spatial dimensions:

7

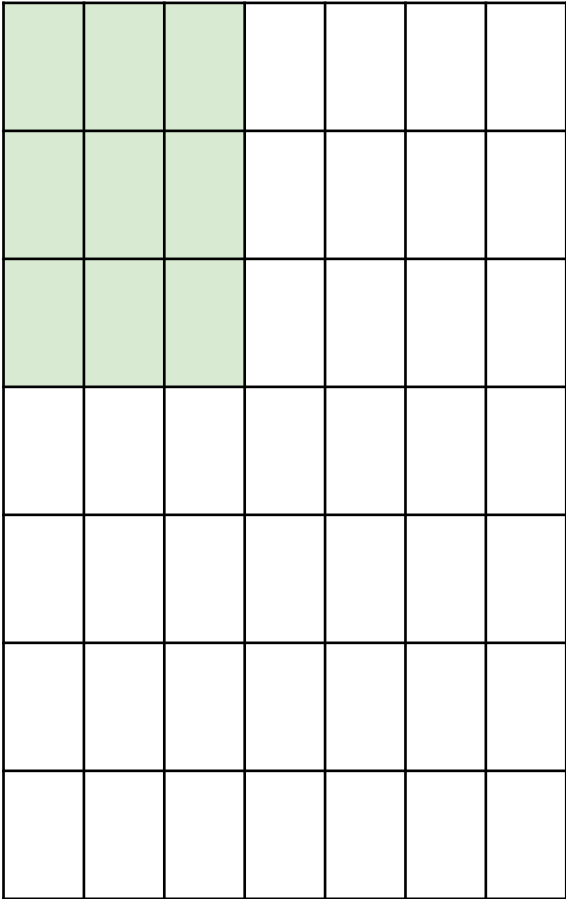


7

7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**  
**=> 3x3 output!**

A closer look at spatial dimensions:

7

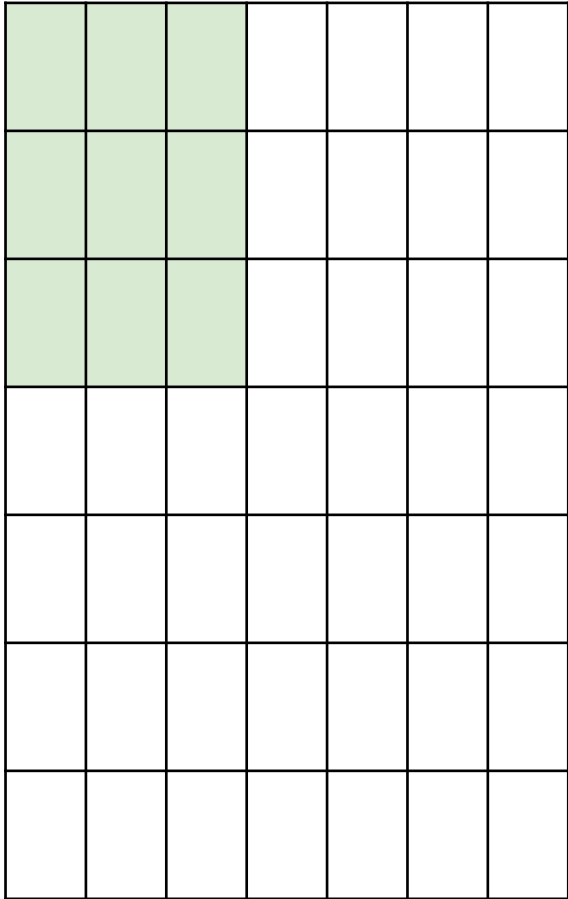


7

7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 3?**

A closer look at spatial dimensions:

7

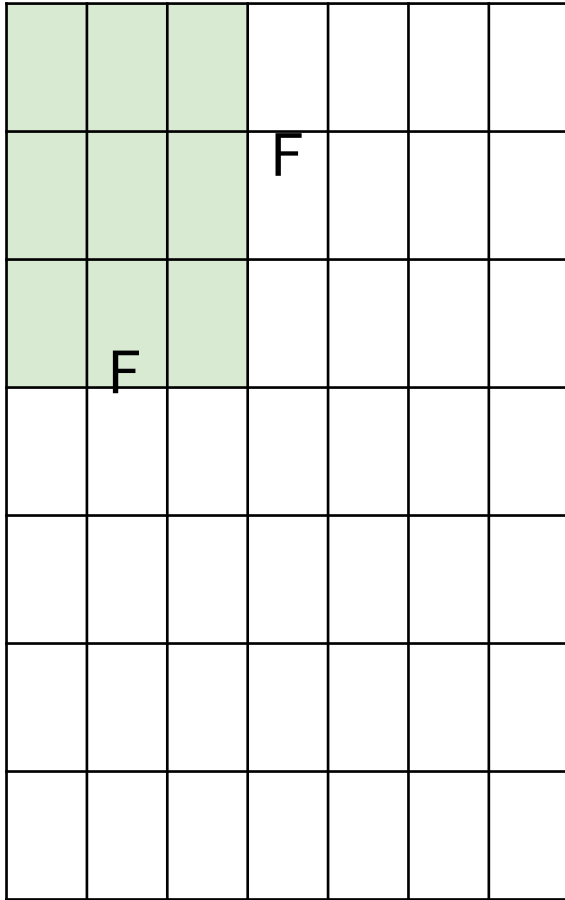


7

7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 3?**

**doesn't fit!**  
cannot apply 3x3 filter on  
7x7 input with stride 3.

N



Output size:

$$(N - F) / \text{stride} + 1$$

e.g.  $N = 7, F = 3$ :

$$\text{stride } 1 \Rightarrow (7 - 3) / 1 + 1 = 5$$

$$\text{stride } 2 \Rightarrow (7 - 3) / 2 + 1 = 3$$

$$\text{stride } 3 \Rightarrow (7 - 3) / 3 + 1 = 2.33 : \backslash$$

# In practice: Common to zero pad the border

|   |   |   |   |   |   |  |  |  |
|---|---|---|---|---|---|--|--|--|
| 0 | 0 | 0 | 0 | 0 | 0 |  |  |  |
| 0 |   |   |   |   |   |  |  |  |
| 0 |   |   |   |   |   |  |  |  |
| 0 |   |   |   |   |   |  |  |  |
| 0 |   |   |   |   |   |  |  |  |
|   |   |   |   |   |   |  |  |  |
|   |   |   |   |   |   |  |  |  |
|   |   |   |   |   |   |  |  |  |
|   |   |   |   |   |   |  |  |  |

e.g. input 7x7

**3x3** filter, applied with **stride 1**

**pad with 1 pixel** border => what is the output?

(recall:)

$$(N - F) / \text{stride} + 1$$

# In practice: Common to zero pad the border

|   |   |   |   |   |   |  |  |  |
|---|---|---|---|---|---|--|--|--|
| 0 | 0 | 0 | 0 | 0 | 0 |  |  |  |
| 0 |   |   |   |   |   |  |  |  |
| 0 |   |   |   |   |   |  |  |  |
| 0 |   |   |   |   |   |  |  |  |
| 0 |   |   |   |   |   |  |  |  |
|   |   |   |   |   |   |  |  |  |
|   |   |   |   |   |   |  |  |  |
|   |   |   |   |   |   |  |  |  |
|   |   |   |   |   |   |  |  |  |

e.g. input 7x7

**3x3** filter, applied with **stride 1**

**pad with 1 pixel** border => what is the output?

**7x7** output!

# In practice: Common to zero pad the border

|   |   |   |   |   |   |  |  |  |
|---|---|---|---|---|---|--|--|--|
| 0 | 0 | 0 | 0 | 0 | 0 |  |  |  |
| 0 |   |   |   |   |   |  |  |  |
| 0 |   |   |   |   |   |  |  |  |
| 0 |   |   |   |   |   |  |  |  |
| 0 |   |   |   |   |   |  |  |  |
|   |   |   |   |   |   |  |  |  |
|   |   |   |   |   |   |  |  |  |
|   |   |   |   |   |   |  |  |  |
|   |   |   |   |   |   |  |  |  |

e.g. input 7x7

**3x3** filter, applied with **stride 1**

**pad with 1 pixel** border => what is the output?

**7x7 output!**

in general, common to see CONV layers with stride 1, filters of size  $F \times F$ , and zero-padding with  $(F-1)/2$ . (will preserve size spatially)

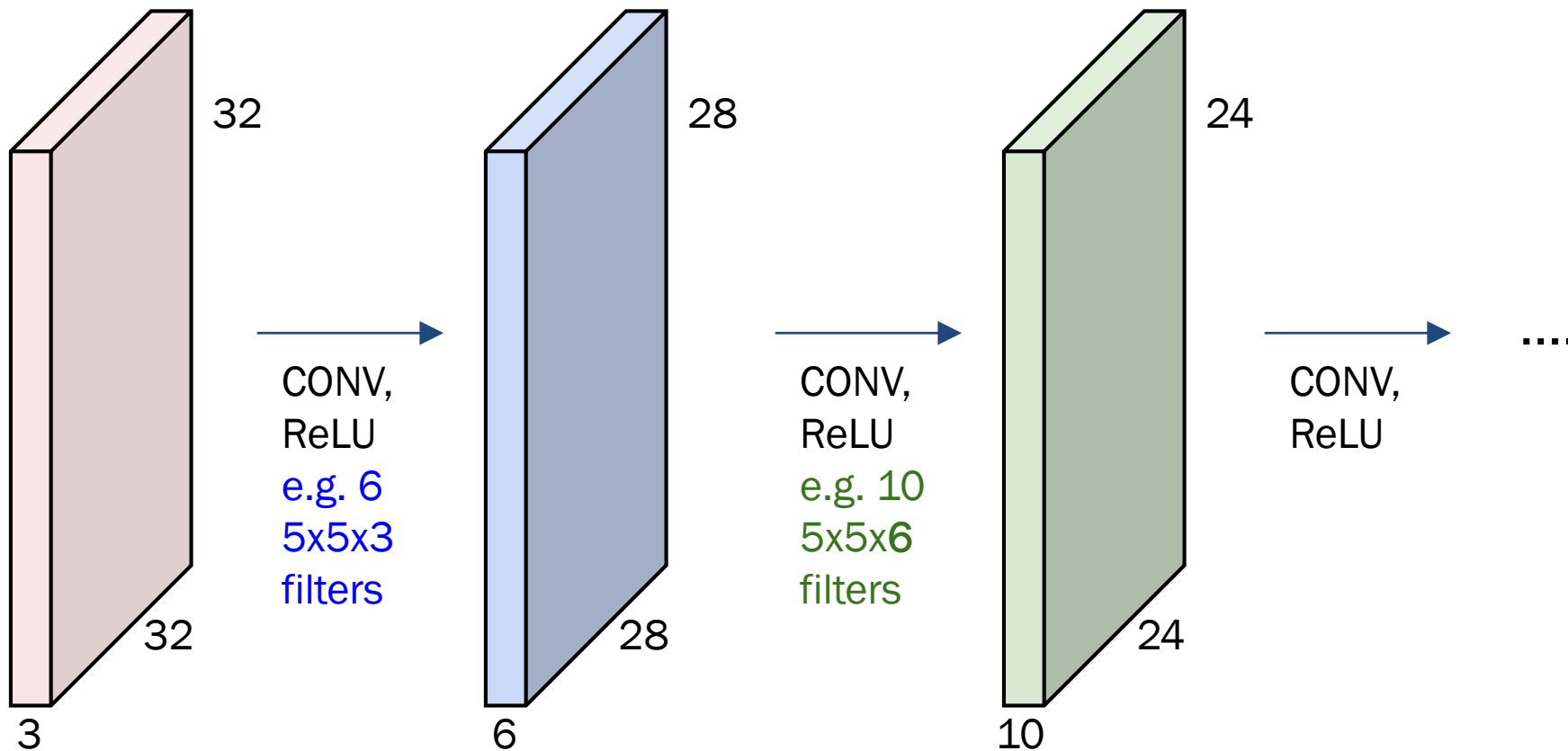
e.g.  $F = 3 \Rightarrow$  zero pad with 1

$F = 5 \Rightarrow$  zero pad with 2

$F = 7 \Rightarrow$  zero pad with 3

## Remember back to...

E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially!  
(32 -> 28 -> 24 ...). Shrinking too fast is not good, doesn't work well.

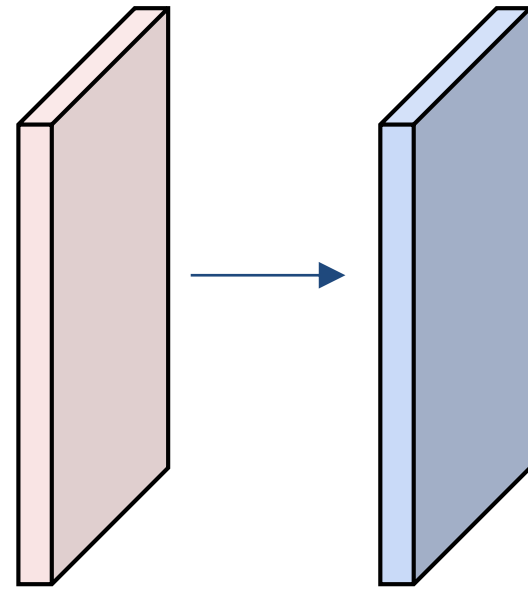


Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

Output volume size: ?



Examples time:

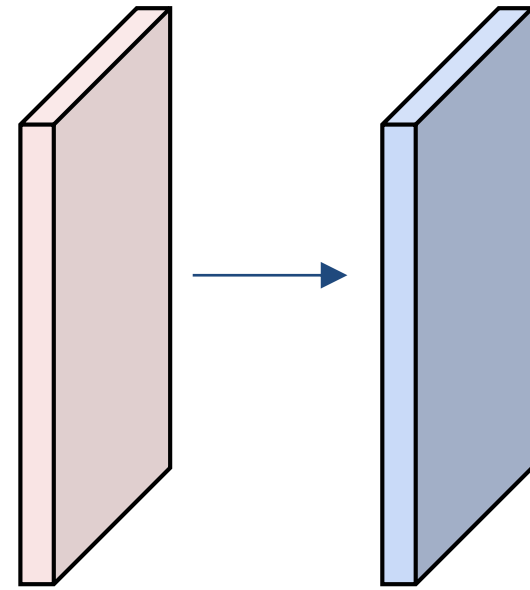
Input volume: **32x32x3**

**10** **5x5** filters with stride **1**, pad **2**

Output volume size:

$(32 + 2 * 2 - 5) / 1 + 1 = 32$  spatially, so

**32x32x10**

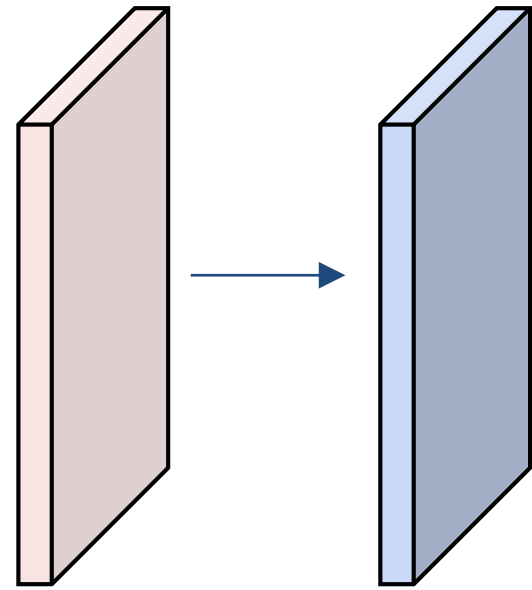


Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

Number of parameters in this layer?



Examples time:

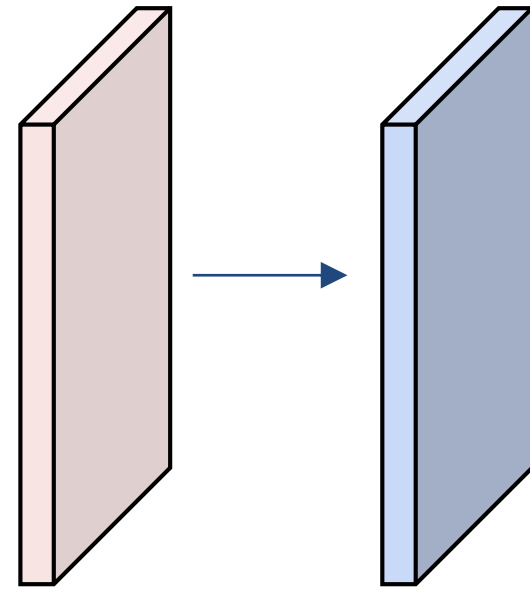
Input volume: **32x32x3**

**10** **5x5** filters with stride 1, pad 2

Number of parameters in this layer?

each filter has  $5 * 5 * 3 + 1 = 76$  params (+1 for bias)

=>  $76 * 10 = 760$



**Summary.** To summarize, the Conv Layer:

- 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.
- In the output volume, the  $d$ -th depth slice (of size  $W_2 \times H_2$ ) is the result of performing a valid convolution of the  $d$ -th filter over the input volume with a stride of  $S$ , and then offset by  $d$ -th bias.

## Common settings:

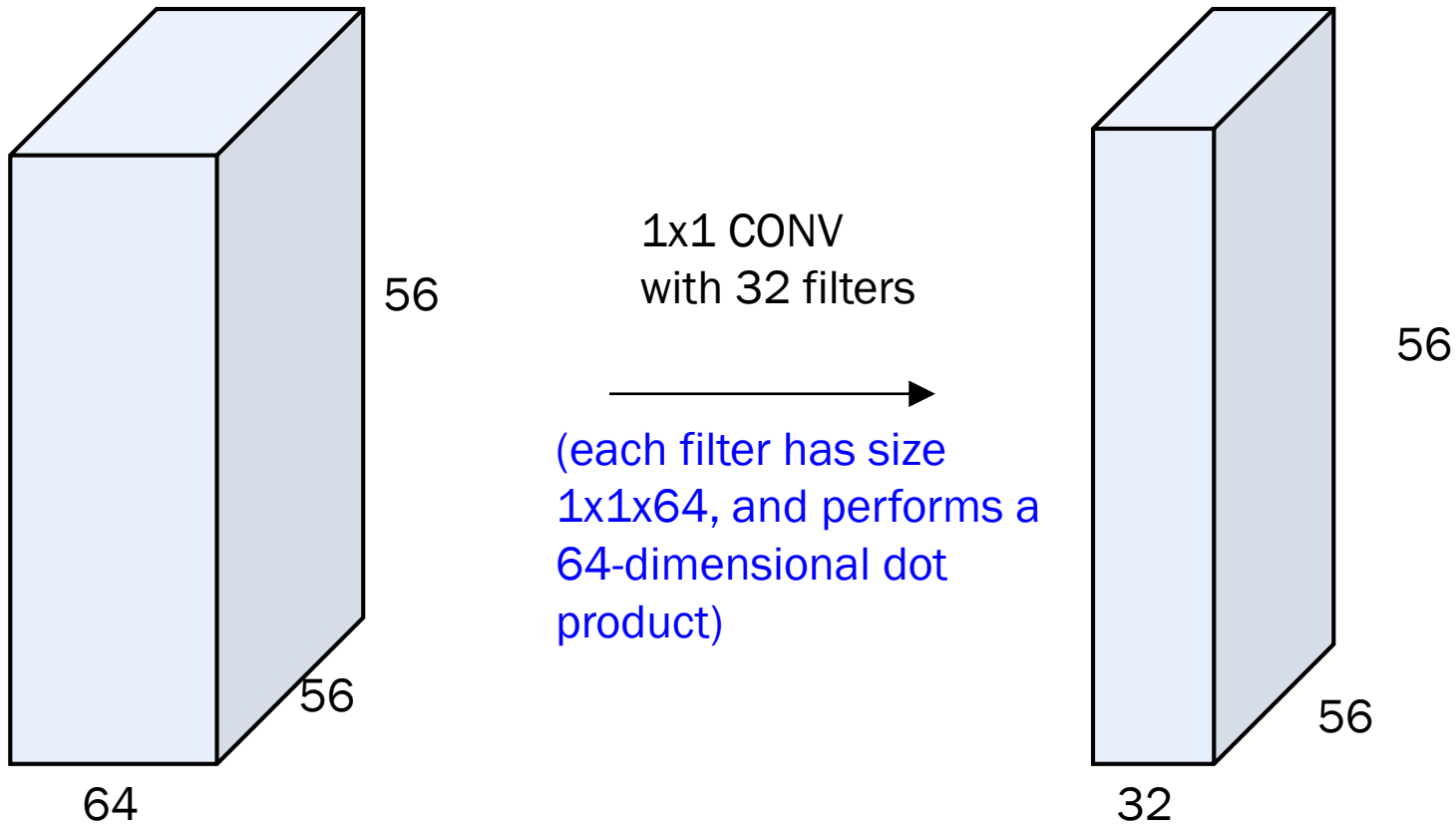
$K =$  (powers of 2, e.g. 32, 64, 128, 512)

- $F = 3, S = 1, P = 1$
- $F = 5, S = 1, P = 2$
- $F = 5, S = 2, P = ?$  (whatever fits)
- $F = 1, S = 1, P = 0$

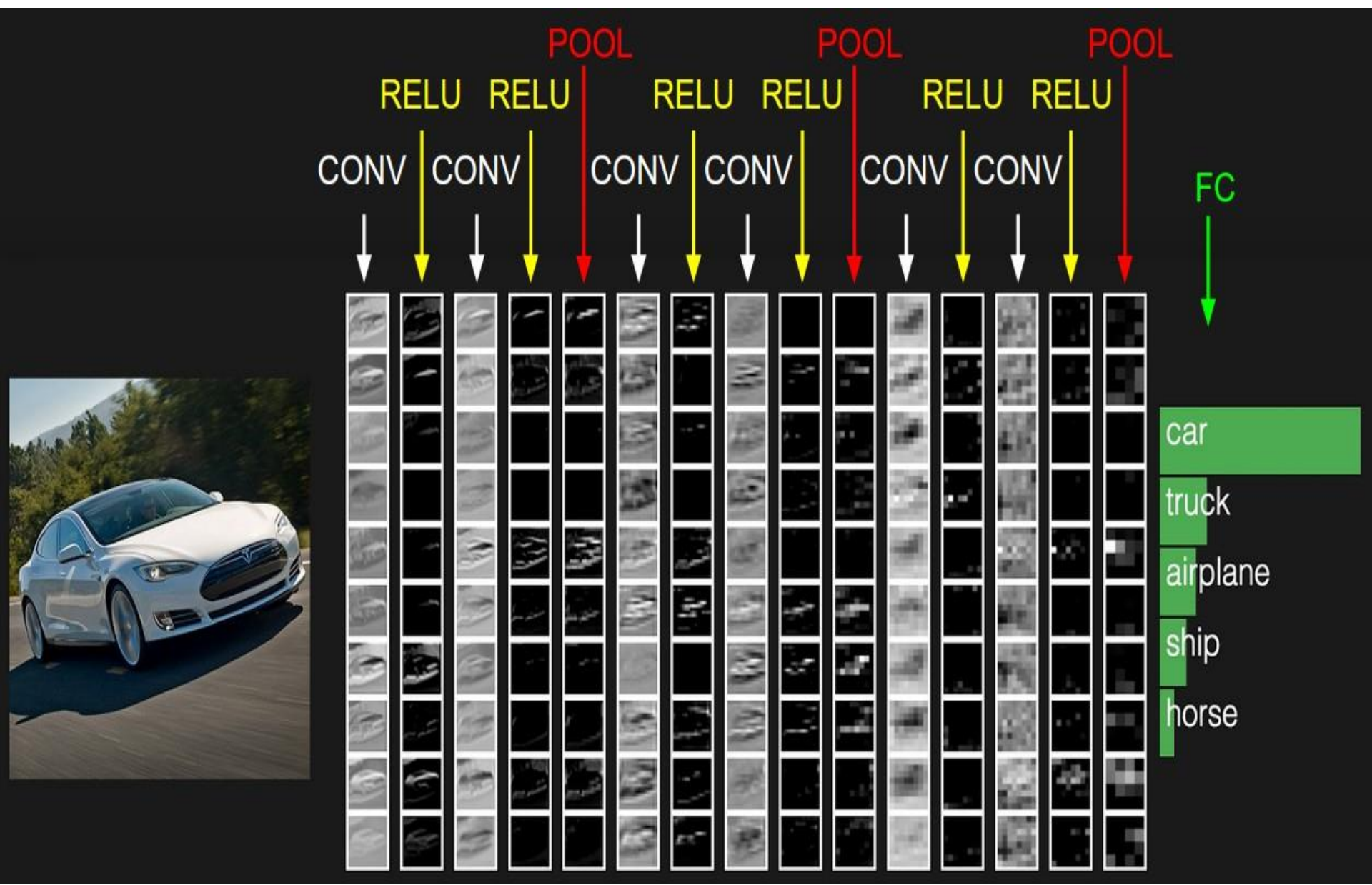
**Summary.** To summarize, the Conv Layer:

- 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.
- In the output volume, the  $d$ -th depth slice (of size  $W_2 \times H_2$ ) is the result of performing a valid convolution of the  $d$ -th filter over the input volume with a stride of  $S$ , and then offset by  $d$ -th bias.

(btw, 1x1 convolution layers make perfect sense)

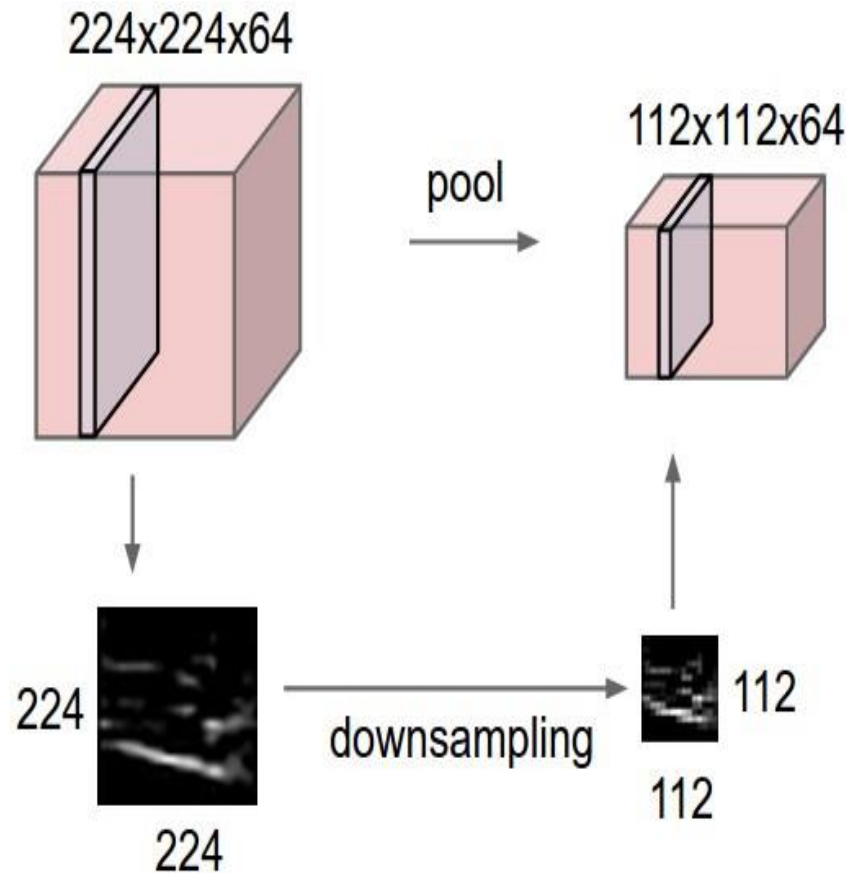


two more layers to go: POOL/FC



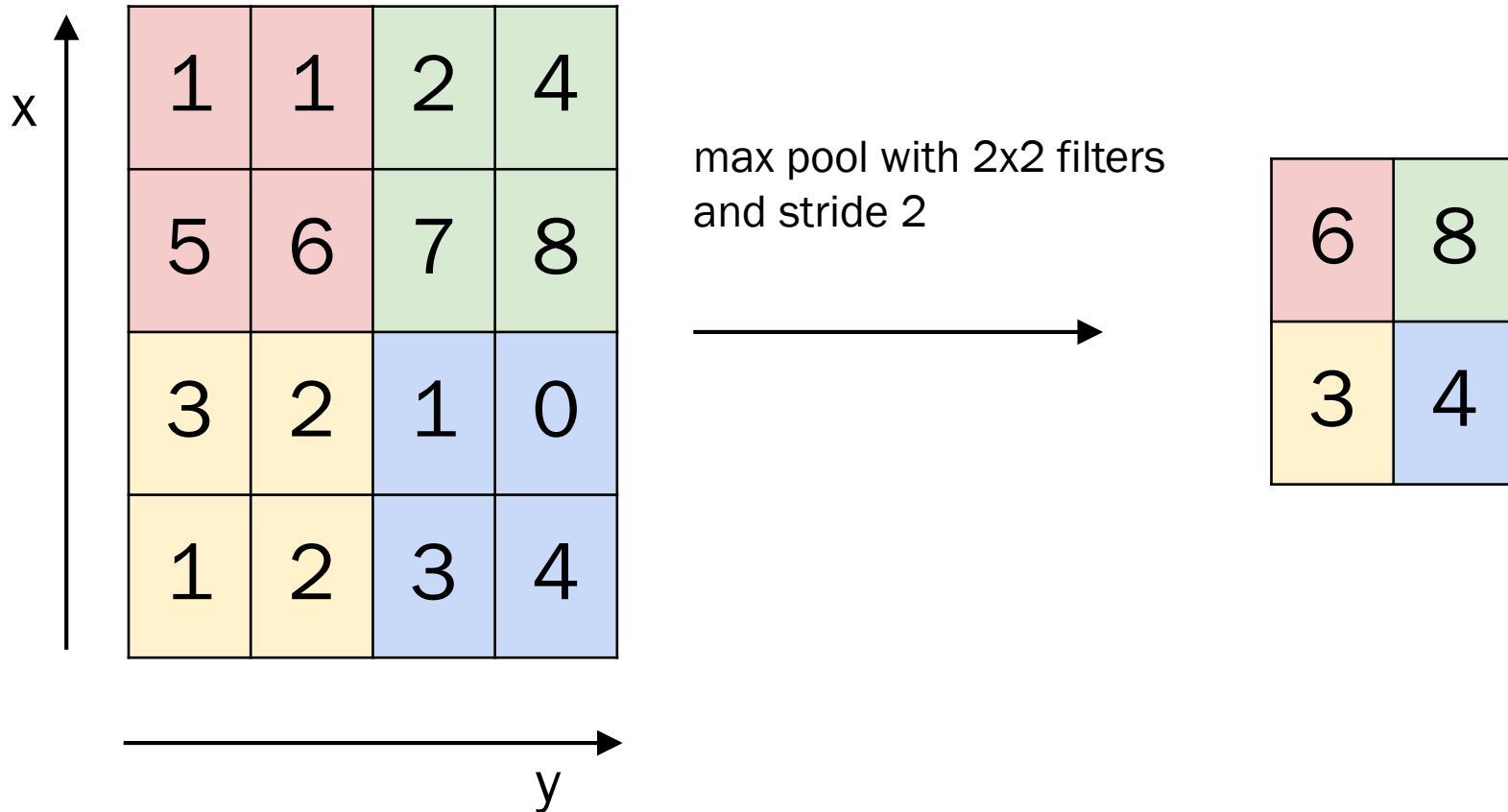
# Pooling layer

- makes the representations smaller and more manageable
- operates over each activation map independently:



# MAX POOLING

Single depth slice



- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
  - their spatial extent  $F$ ,
  - the stride  $S$ ,
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F) / S + 1$
  - $H_2 = (H_1 - F) / S + 1$
  - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

Common settings:

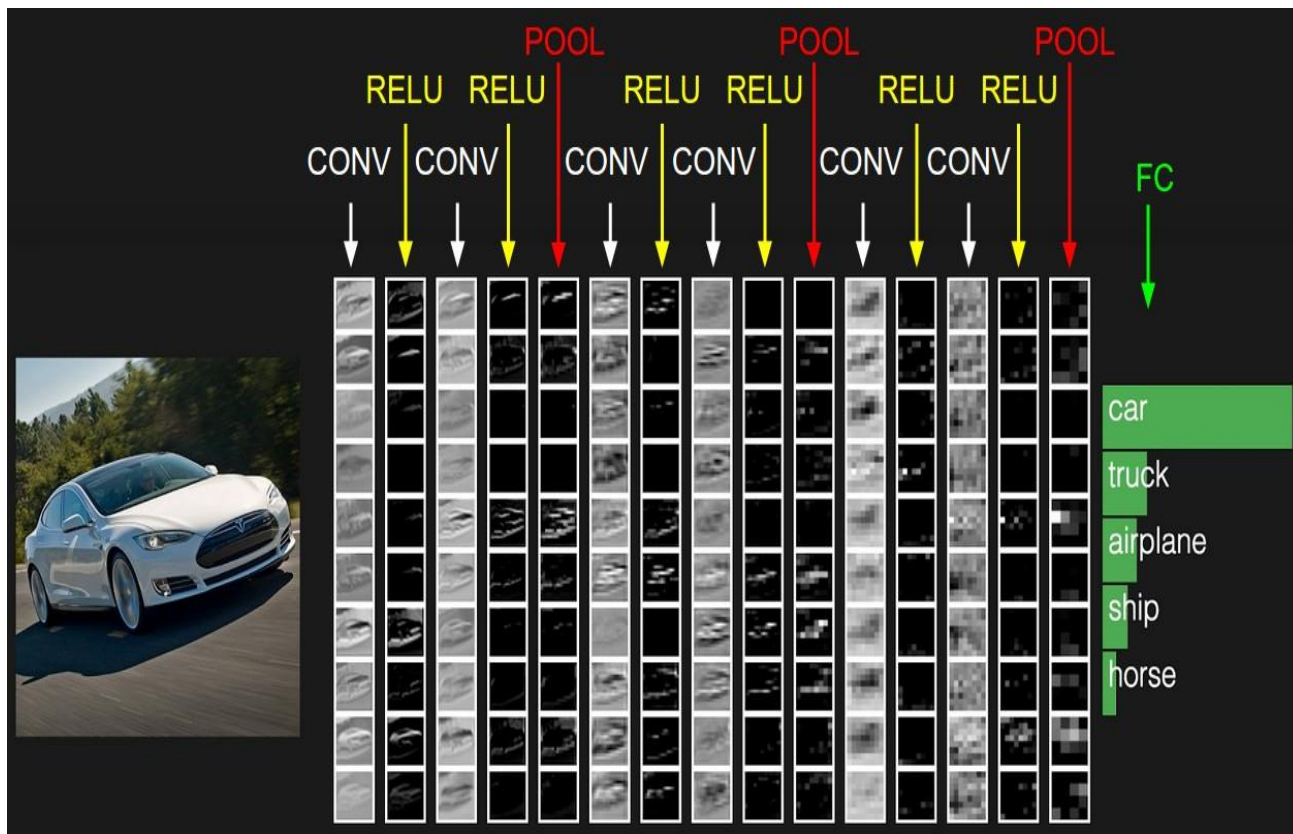
$$F = 2, S = 2$$

$$F = 3, S = 2$$

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
  - their spatial extent  $F$ ,
  - the stride  $S$ ,
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F)/S + 1$
  - $H_2 = (H_1 - F)/S + 1$
  - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

# Fully Connected Layer (FC layer)

- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks



# Demo

## Example predictions on Test set

test accuracy based on last 200 test images: 0.45



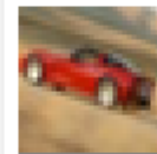
car  
ship  
airplane



car  
truck  
ship



ship  
airplane  
car



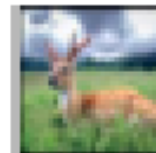
deer  
airplane  
bird



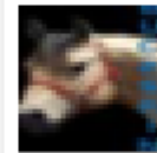
deer  
cat  
horse



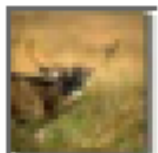
ship  
airplane  
horse



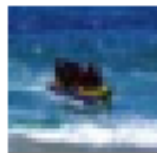
horse  
deer  
bird



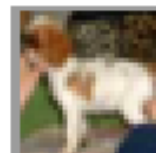
cat  
truck  
horse



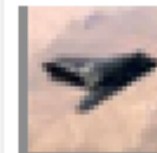
deer  
bird  
frog



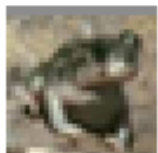
airplane  
deer  
bird



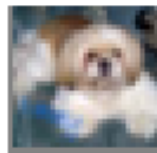
deer  
dog  
cat



airplane  
bird  
deer



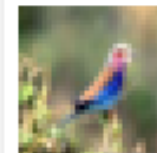
frog  
deer  
cat



dog  
cat  
ship



bird  
deer  
frog

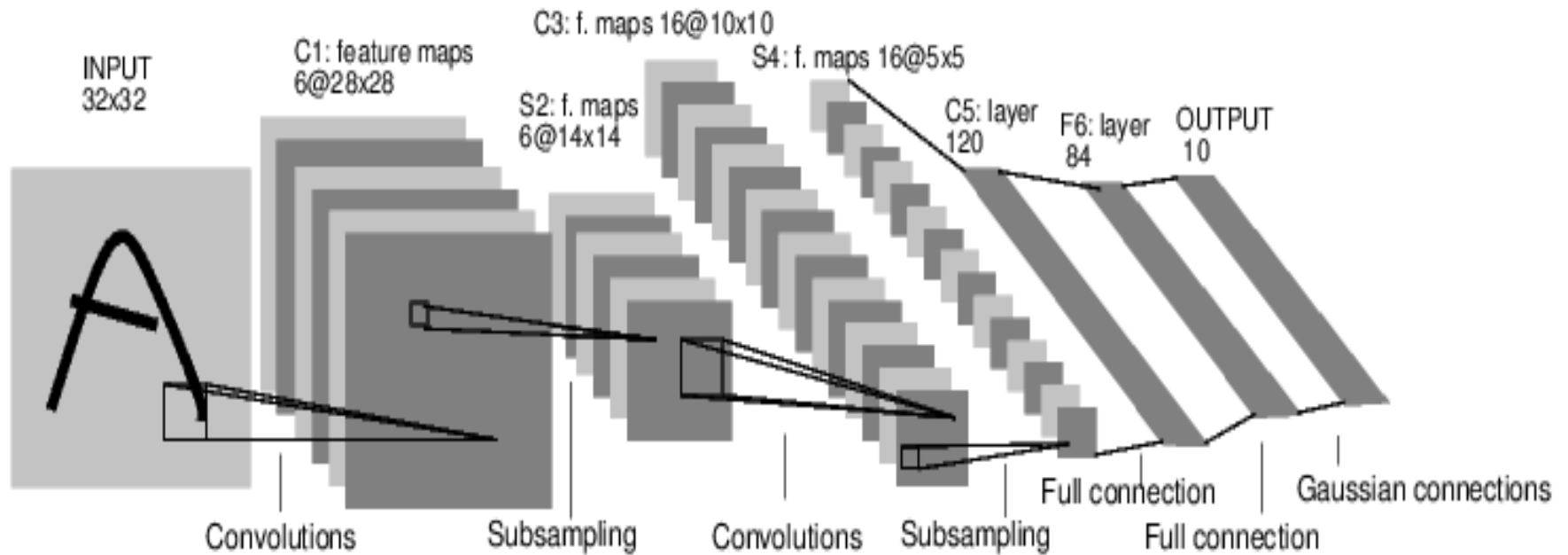


frog  
bird  
deer

<http://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

# Case Study: LeNet-5

[LeCun et al., 1998]



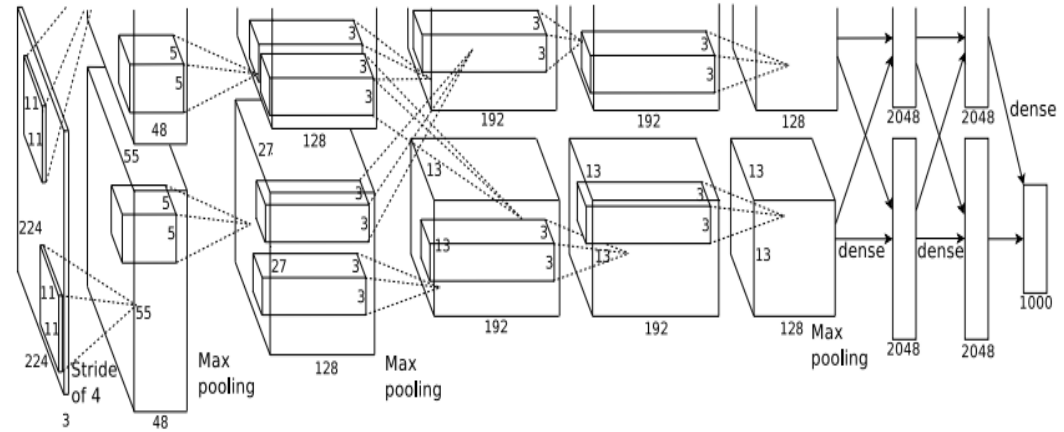
Conv filters were 5x5, applied at stride 1

Subsampling (Pooling) layers were 2x2 applied at stride 2

i.e. architecture is [CONV-POOL-CONV-POOL-CONV-FC]

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

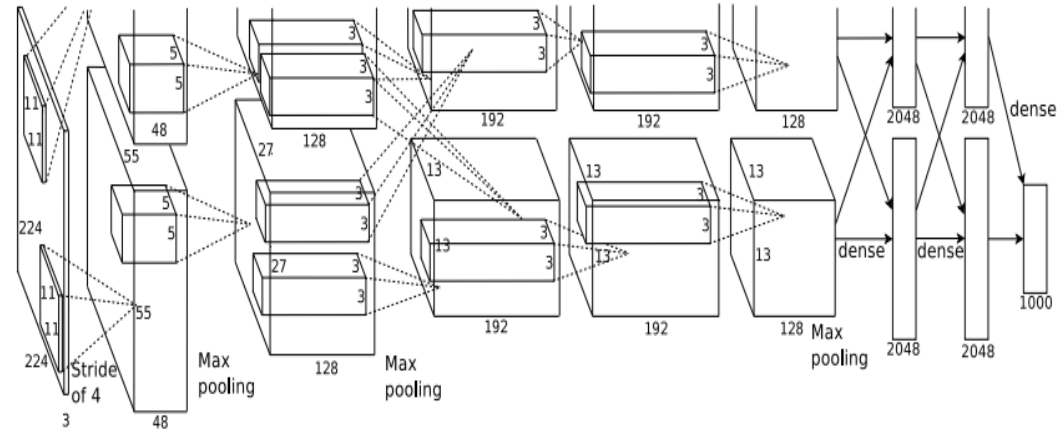
First layer (CONV1): 96 11x11 filters applied at stride 4

=>

Q: what is the output volume size? Hint:  $(227-11)/4+1 = 55$

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

**First layer (CONV1):** 96 11x11 filters applied at stride 4

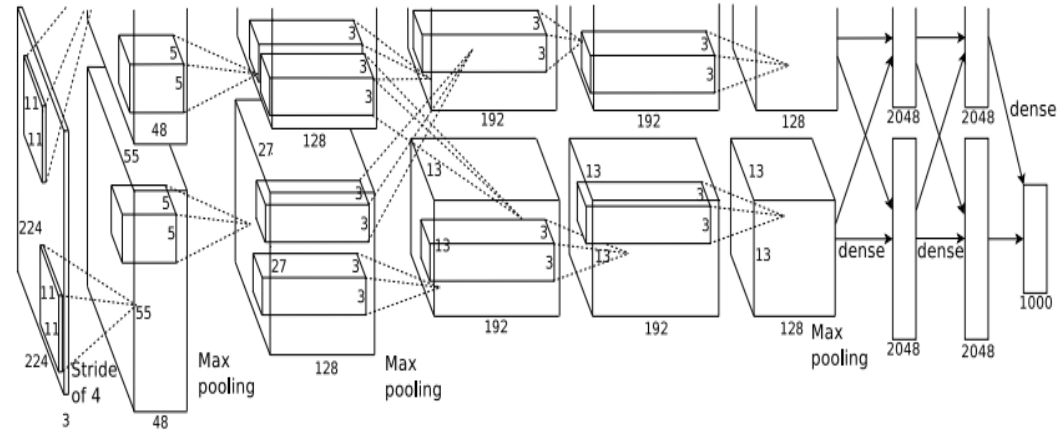
=>

Output volume [55x55x96]

Q: What is the total number of parameters in this layer?

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

**First layer (CONV1): 96 11x11 filters applied at stride 4**

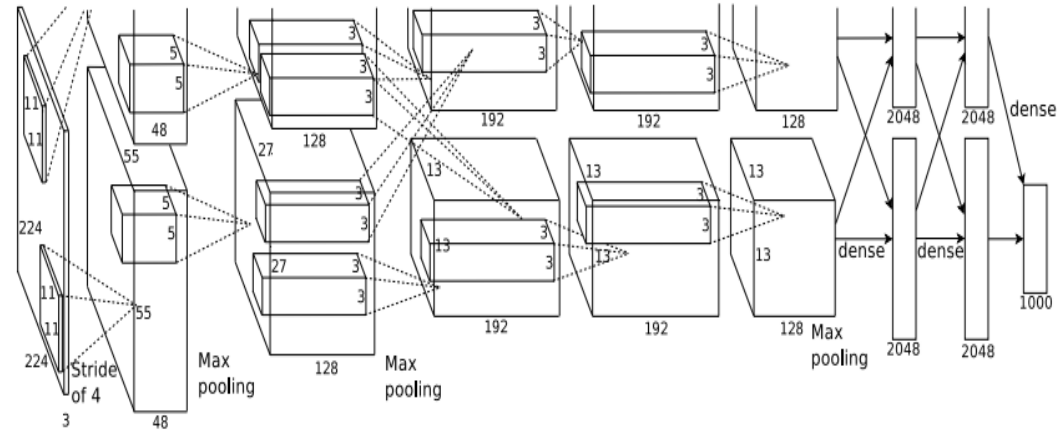
=>

Output volume [55x55x96]

Parameters:  $(11 * 11 * 3) * 96 = 35\text{K}$

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

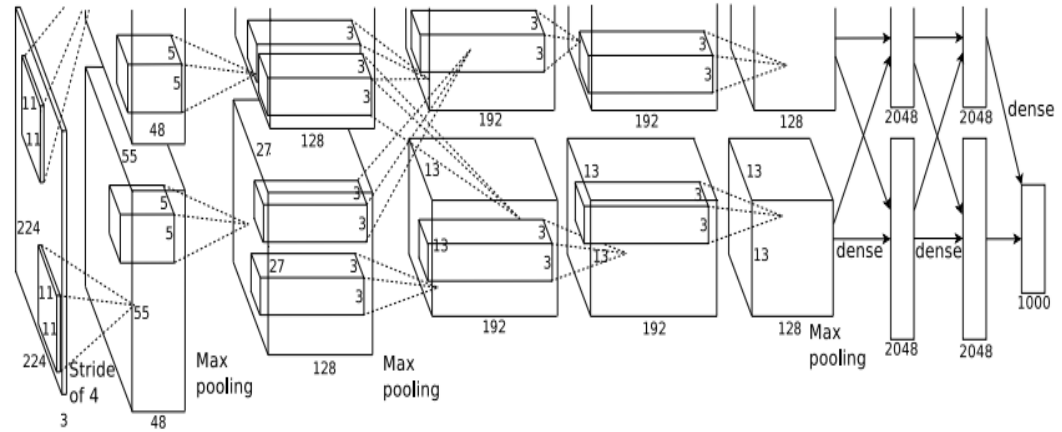
After CONV1: 55x55x96

Second layer (POOL1): 3x3 filters applied at stride 2

Q: what is the output volume size? Hint:  $(55-3)/2+1 = 27$

# Case Study: AlexNet

[Krizhevsky et al. 2012]



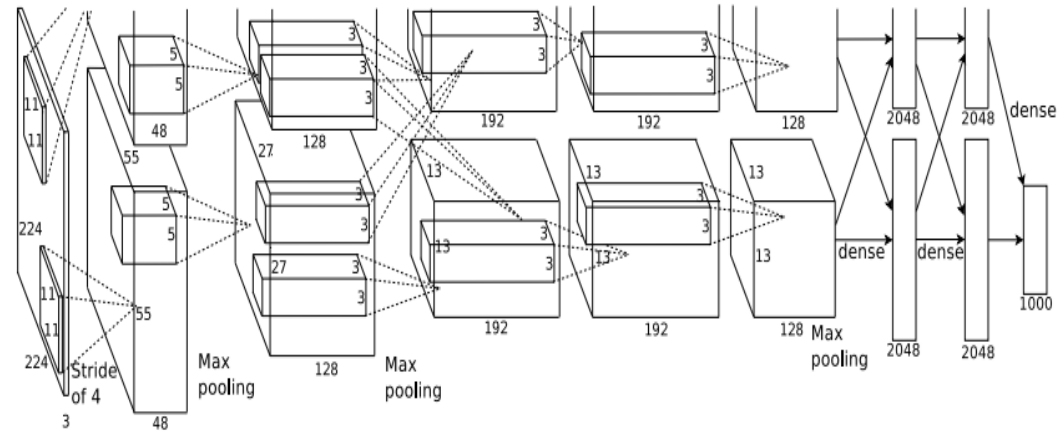
Input: 227x227x3 images  
After CONV1: 55x55x96

**Second layer (POOL1):** 3x3 filters applied at stride 2  
Output volume: 27x27x96

Q: what is the number of parameters in this layer?

# Case Study: AlexNet

[Krizhevsky et al. 2012]

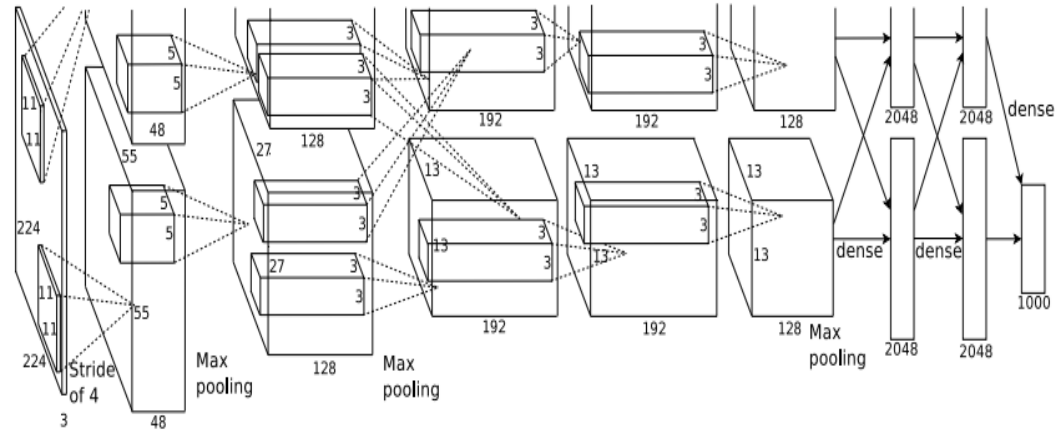


Input: 227x227x3 images  
After CONV1: 55x55x96

**Second layer (POOL1):** 3x3 filters applied at stride 2  
Output volume: 27x27x96  
Parameters: 0!

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

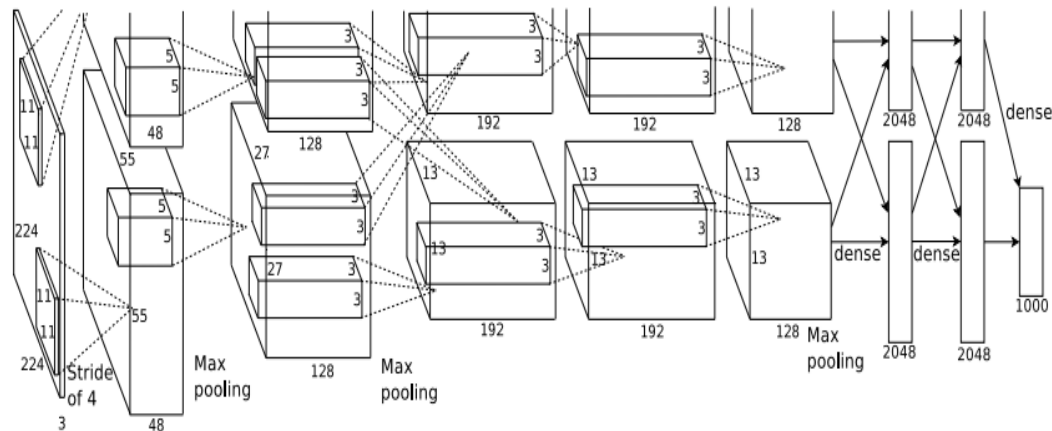
After CONV1: 55x55x96

After POOL1: 27x27x96

...

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1**: 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1**: 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2**: 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3**: 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4**: 384 3x3 filters at stride 1, pad 1

[13x13x256] **CONV5**: 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3**: 3x3 filters at stride 2

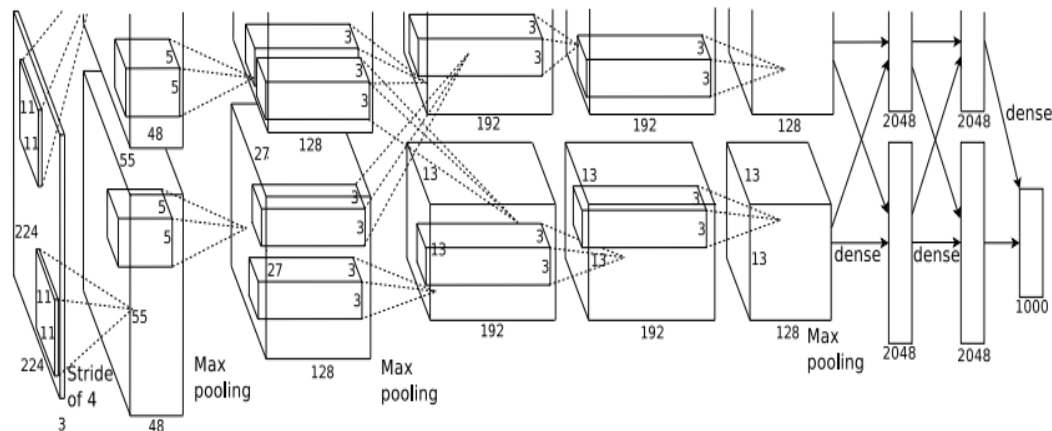
[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1**: 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1**: 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2**: 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3**: 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4**: 384 3x3 filters at stride 1, pad 1

[13x13x256] **CONV5**: 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3**: 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

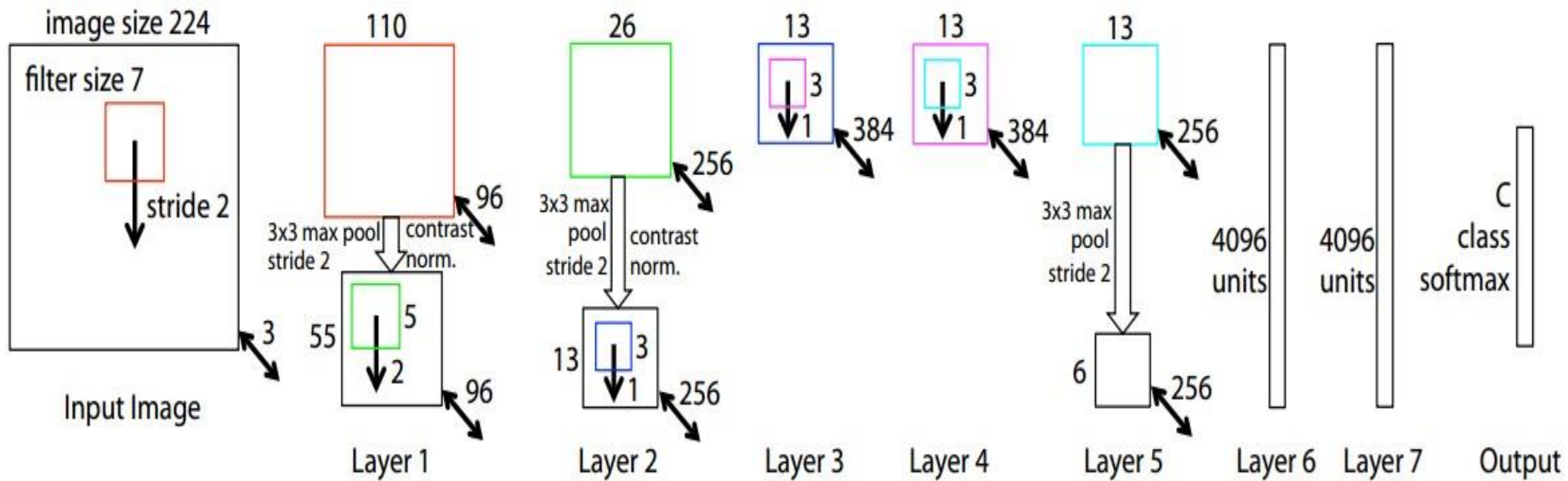
[1000] **FC8**: 1000 neurons (class scores)

## Details/Retrospectives:

- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% -> 15.4%

# Case Study: ZFNet

[Zeiler and Fergus, 2013]



AlexNet but:

CONV1: change from (11x11 stride 4) to (7x7 stride 2)

CONV3,4,5: instead of 384, 384, 256 filters use 512, 1024, 512

ImageNet top 5 error: 15.4% -> 14.8%

# Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Only 3x3 CONV stride 1, pad 1 and 2x2 MAX POOL stride 2

best model

11.2% top 5 error in ILSVRC 2013

->

7.3% top 5 error

| ConvNet Configuration       |                        |                               |  |  |   |
|-----------------------------|------------------------|-------------------------------|--|--|---|
| A                           | A-LRN                  | B                             | C  | D  | E   |
| 11 weight layers            | 11 weight layers       | 13 weight layers              | 16 weight layers                           | 16 weight layers                           | 19 weight layers  |
| input (224 × 224 RGB image) |                        |                               |  |  |   |
| conv3-64                    | conv3-64<br><b>LRN</b> | conv3-64<br><b>conv3-64</b>   | conv3-64<br>conv3-64                       | conv3-64<br>conv3-64                       | conv3-64<br>conv3-64                                    |
| maxpool                     |                        |                               |  |  |   |
| conv3-128                   | conv3-128              | conv3-128<br><b>conv3-128</b> | conv3-128<br>conv3-128                     | conv3-128<br>conv3-128                     | conv3-128<br>conv3-128                                  |
| maxpool                     |                        |                               |  |  |   |
| conv3-256<br>conv3-256      | conv3-256<br>conv3-256 | conv3-256<br>conv3-256        | conv3-256<br>conv3-256<br><b>conv1-256</b> | conv3-256<br>conv3-256<br><b>conv3-256</b> | conv3-256<br>conv3-256<br>conv3-256<br><b>conv3-256</b> |
| maxpool                     |                        |                               |  |  |   |
| conv3-512<br>conv3-512      | conv3-512<br>conv3-512 | conv3-512<br>conv3-512        | conv3-512<br>conv3-512<br><b>conv1-512</b> | conv3-512<br>conv3-512<br><b>conv3-512</b> | conv3-512<br>conv3-512<br>conv3-512<br><b>conv3-512</b> |
| maxpool                     |                        |                               |  |  |   |
| conv3-512<br>conv3-512      | conv3-512<br>conv3-512 | conv3-512<br>conv3-512        | conv3-512<br>conv3-512<br><b>conv1-512</b> | conv3-512<br>conv3-512<br><b>conv3-512</b> | conv3-512<br>conv3-512<br>conv3-512<br><b>conv3-512</b> |
| maxpool                     |                        |                               |  |  |   |
| FC-4096                     |                        |                               |  |  |   |
| FC-4096                     |                        |                               |  |  |   |
| FC-1000                     |                        |                               |  |  |   |
| soft-max                    |                        |                               |  |  |   |

Table 2: Number of parameters (in millions).

| Network              | A,A-LRN | B   | C   | D   | E   |
|----------------------|---------|-----|-----|-----|-----|
| Number of parameters | 133     | 133 | 134 | 138 | 144 |

(not counting biases)

INPUT: [224x224x3] memory:  $224*224*3=150K$  params: 0

CONV3-64: [224x224x64] memory:  $224*224*64=3.2M$  params:  $(3*3*3)*64 = 1,728$

CONV3-64: [224x224x64] memory:  $224*224*64=3.2M$  params:  $(3*3*64)*64 = 36,864$

POOL2: [112x112x64] memory:  $112*112*64=800K$  params: 0

CONV3-128: [112x112x128] memory:  $112*112*128=1.6M$  params:  $(3*3*64)*128 = 73,728$

CONV3-128: [112x112x128] memory:  $112*112*128=1.6M$  params:  $(3*3*128)*128 = 147,456$

POOL2: [56x56x128] memory:  $56*56*128=400K$  params: 0

CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*128)*256 = 294,912$

CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*256)*256 = 589,824$

CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*256)*256 = 589,824$

POOL2: [28x28x256] memory:  $28*28*256=200K$  params: 0

CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*256)*512 = 1,179,648$

CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*512)*512 = 2,359,296$

POOL2: [14x14x512] memory:  $14*14*512=100K$  params: 0

CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$

POOL2: [7x7x512] memory:  $7*7*512=25K$  params: 0

FC: [1x1x4096] memory: 4096 params:  $7*7*512*4096 = 102,760,448$

FC: [1x1x4096] memory: 4096 params:  $4096*4096 = 16,777,216$

FC: [1x1x1000] memory: 1000 params:  $4096*1000 = 4,096,000$

| ConvNet Configuration       |                  |                  |    |
|-----------------------------|------------------|------------------|----|
| B                           | C                | D                |    |
| 13 weight layers            | 16 weight layers | 16 weight layers | 19 |
| input (224 x 224 RGB image) |                  |                  |    |
| conv3-64                    | conv3-64         | conv3-64         | cc |
| <b>conv3-64</b>             | conv3-64         | conv3-64         | cc |
| maxpool                     |                  |                  |    |
| conv3-128                   | conv3-128        | conv3-128        | co |
| <b>conv3-128</b>            | conv3-128        | conv3-128        | co |
| maxpool                     |                  |                  |    |
| conv3-256                   | conv3-256        | conv3-256        | co |
| conv3-256                   | conv3-256        | conv3-256        | co |
|                             | <b>conv1-256</b> | <b>conv3-256</b> | co |
|                             |                  |                  | co |
| maxpool                     |                  |                  |    |
| conv3-512                   | conv3-512        | conv3-512        | co |
| conv3-512                   | conv3-512        | conv3-512        | co |
|                             | <b>conv1-512</b> | <b>conv3-512</b> | co |
|                             |                  |                  | co |
| maxpool                     |                  |                  |    |
| conv3-512                   | conv3-512        | conv3-512        | co |
| conv3-512                   | conv3-512        | conv3-512        | co |
|                             | <b>conv1-512</b> | <b>conv3-512</b> | co |
|                             |                  |                  | co |
| maxpool                     |                  |                  |    |
| FC-4096                     |                  |                  |    |
| FC-4096                     |                  |                  |    |
| FC-1000                     |                  |                  |    |
| soft-max                    |                  |                  |    |

(not counting biases)

INPUT: [224x224x3] memory:  $224*224*3=150K$  params: 0

CONV3-64: [224x224x64] memory:  $224*224*64=3.2M$  params:  $(3*3*3)*64 = 1,728$

CONV3-64: [224x224x64] memory:  $224*224*64=3.2M$  params:  $(3*3*64)*64 = 36,864$

POOL2: [112x112x64] memory:  $112*112*64=800K$  params: 0

CONV3-128: [112x112x128] memory:  $112*112*128=1.6M$  params:  $(3*3*64)*128 = 73,728$

CONV3-128: [112x112x128] memory:  $112*112*128=1.6M$  params:  $(3*3*128)*128 = 147,456$

POOL2: [56x56x128] memory:  $56*56*128=400K$  params: 0

CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*128)*256 = 294,912$

CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*256)*256 = 589,824$

CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*256)*256 = 589,824$

POOL2: [28x28x256] memory:  $28*28*256=200K$  params: 0

CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*256)*512 = 1,179,648$

CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*512)*512 = 2,359,296$

POOL2: [14x14x512] memory:  $14*14*512=100K$  params: 0

CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$

POOL2: [7x7x512] memory:  $7*7*512=25K$  params: 0

FC: [1x1x4096] memory: 4096 params:  $7*7*512*4096 = 102,760,448$

FC: [1x1x4096] memory: 4096 params:  $4096*4096 = 16,777,216$

FC: [1x1x1000] memory: 1000 params:  $4096*1000 = 4,096,000$

| ConvNet Configuration       |                  |                  |    |
|-----------------------------|------------------|------------------|----|
| B                           | C                | D                |    |
| 13 weight layers            | 16 weight layers | 16 weight layers | 19 |
| input (224 x 224 RGB image) |                  |                  |    |
| conv3-64                    | conv3-64         | conv3-64         | cc |
| <b>conv3-64</b>             | conv3-64         | conv3-64         | cc |
| maxpool                     |                  |                  |    |
| conv3-128                   | conv3-128        | conv3-128        | co |
| <b>conv3-128</b>            | conv3-128        | conv3-128        | co |
| maxpool                     |                  |                  |    |
| conv3-256                   | conv3-256        | conv3-256        | co |
| conv3-256                   | conv3-256        | conv3-256        | co |
|                             | <b>conv1-256</b> | <b>conv3-256</b> | co |
|                             |                  |                  | co |
| maxpool                     |                  |                  |    |
| conv3-512                   | conv3-512        | conv3-512        | co |
| conv3-512                   | conv3-512        | conv3-512        | co |
|                             | <b>conv1-512</b> | <b>conv3-512</b> | co |
|                             |                  |                  | co |
| maxpool                     |                  |                  |    |
| conv3-512                   | conv3-512        | conv3-512        | co |
| conv3-512                   | conv3-512        | conv3-512        | co |
|                             | <b>conv1-512</b> | <b>conv3-512</b> | co |
|                             |                  |                  | co |
| maxpool                     |                  |                  |    |
| FC-4096                     |                  |                  |    |
| FC-4096                     |                  |                  |    |
| FC-1000                     |                  |                  |    |
| soft-max                    |                  |                  |    |

TOTAL memory: 24M \* 4 bytes ~ = 93MB / image (only forward! ~\*2 for bwd)

TOTAL params: 138M parameters

(not counting biases)

Note:

Most memory is in early CONV

Most params are in late FC

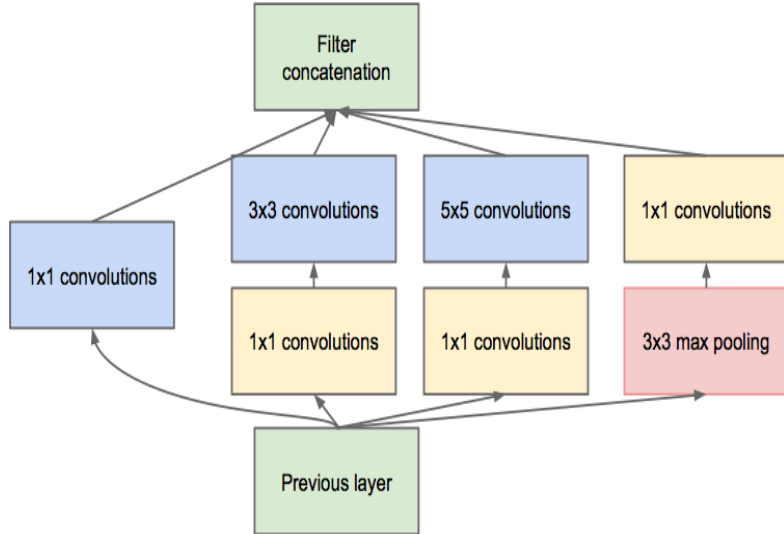
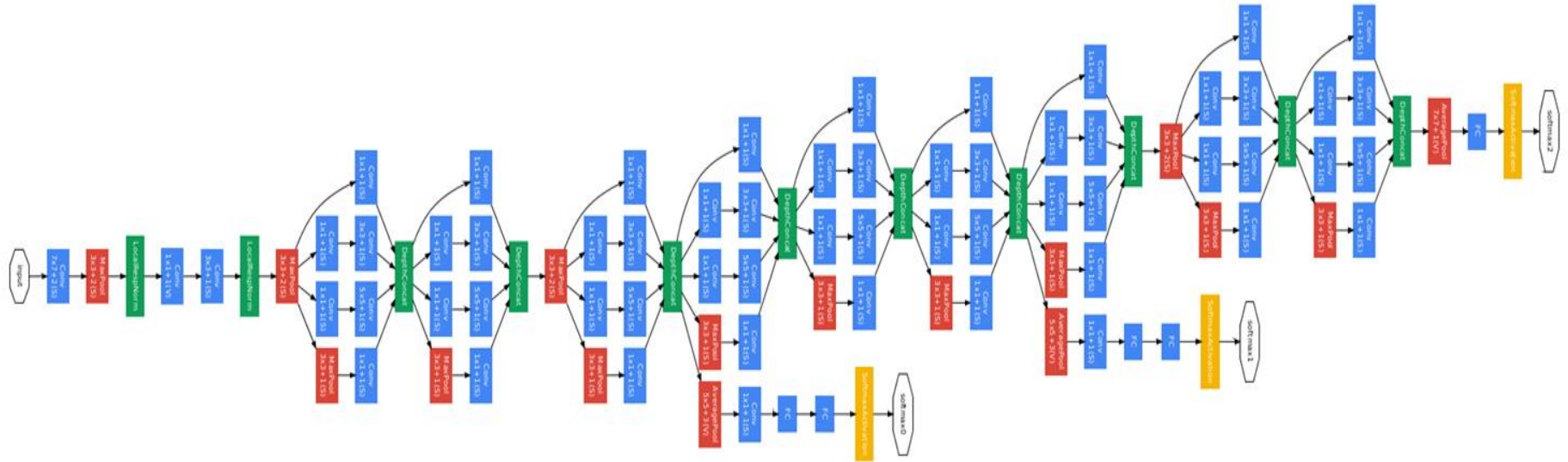
INPUT: [224x224x3] memory:  $224*224*3=150K$  params: 0  
 CONV3-64: [224x224x64] memory:  $224*224*64=3.2M$  params:  $(3*3*3)*64 = 1,728$   
 CONV3-64: [224x224x64] memory:  $224*224*64=3.2M$  params:  $(3*3*64)*64 = 36,864$   
 POOL2: [112x112x64] memory:  $112*112*64=800K$  params: 0  
 CONV3-128: [112x112x128] memory:  $112*112*128=1.6M$  params:  $(3*3*64)*128 = 73,728$   
 CONV3-128: [112x112x128] memory:  $112*112*128=1.6M$  params:  $(3*3*128)*128 = 147,456$   
 POOL2: [56x56x128] memory:  $56*56*128=400K$  params: 0  
 CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*128)*256 = 294,912$   
 CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*256)*256 = 589,824$   
 CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*256)*256 = 589,824$   
 POOL2: [28x28x256] memory:  $28*28*256=200K$  params: 0  
 CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*256)*512 = 1,179,648$   
 CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*512)*512 = 2,359,296$   
 CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*512)*512 = 2,359,296$   
 POOL2: [14x14x512] memory:  $14*14*512=100K$  params: 0  
 CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$   
 CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$   
 CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$   
 POOL2: [7x7x512] memory:  $7*7*512=25K$  params: 0  
 FC: [1x1x4096] memory: 4096 params:  $7*7*512*4096 = 102,760,448$   
 FC: [1x1x4096] memory: 4096 params:  $4096*4096 = 16,777,216$   
 FC: [1x1x1000] memory: 1000 params:  $4096*1000 = 4,096,000$

TOTAL memory: 24M \* 4 bytes  $\approx$  93MB / image (only forward!  $\sim$ \*2 for bwd)

TOTAL params: 138M parameters

# Case Study: GoogLeNet

[Szegedy et al., 2014]



## Inception module

ILSVRC 2014 winner (6.7% top 5 error)

# Case Study: GoogLeNet

| type           | patch size/<br>stride | output<br>size | depth | #1×1 | #3×3<br>reduce | #3×3 | #5×5<br>reduce | #5×5 | pool<br>proj | params | ops  |
|----------------|-----------------------|----------------|-------|------|----------------|------|----------------|------|--------------|--------|------|
| convolution    | 7×7/2                 | 112×112×64     | 1     |      |                |      |                |      |              | 2.7K   | 34M  |
| max pool       | 3×3/2                 | 56×56×64       | 0     |      |                |      |                |      |              |        |      |
| convolution    | 3×3/1                 | 56×56×192      | 2     |      | 64             | 192  |                |      |              | 112K   | 360M |
| max pool       | 3×3/2                 | 28×28×192      | 0     |      |                |      |                |      |              |        |      |
| inception (3a) |                       | 28×28×256      | 2     | 64   | 96             | 128  | 16             | 32   | 32           | 159K   | 128M |
| inception (3b) |                       | 28×28×480      | 2     | 128  | 128            | 192  | 32             | 96   | 64           | 380K   | 304M |
| max pool       | 3×3/2                 | 14×14×480      | 0     |      |                |      |                |      |              |        |      |
| inception (4a) |                       | 14×14×512      | 2     | 192  | 96             | 208  | 16             | 48   | 64           | 364K   | 73M  |
| inception (4b) |                       | 14×14×512      | 2     | 160  | 112            | 224  | 24             | 64   | 64           | 437K   | 88M  |
| inception (4c) |                       | 14×14×512      | 2     | 128  | 128            | 256  | 24             | 64   | 64           | 463K   | 100M |
| inception (4d) |                       | 14×14×528      | 2     | 112  | 144            | 288  | 32             | 64   | 64           | 580K   | 119M |
| inception (4e) |                       | 14×14×832      | 2     | 256  | 160            | 320  | 32             | 128  | 128          | 840K   | 170M |
| max pool       | 3×3/2                 | 7×7×832        | 0     |      |                |      |                |      |              |        |      |
| inception (5a) |                       | 7×7×832        | 2     | 256  | 160            | 320  | 32             | 128  | 128          | 1072K  | 54M  |
| inception (5b) |                       | 7×7×1024       | 2     | 384  | 192            | 384  | 48             | 128  | 128          | 1388K  | 71M  |
| avg pool       | 7×7/1                 | 1×1×1024       | 0     |      |                |      |                |      |              |        |      |
| dropout (40%)  |                       | 1×1×1024       | 0     |      |                |      |                |      |              |        |      |
| linear         |                       | 1×1×1000       | 1     |      |                |      |                |      |              | 1000K  | 1M   |
| softmax        |                       | 1×1×1000       | 0     |      |                |      |                |      |              |        |      |

Fun features:

- Only 5 million params!  
(Removes FC layers completely)

Compared to AlexNet:

- 12X less params
- 2x more compute
- 6.67% (vs. 16.4%)

# Case Study: ResNet

[He et al., 2015]


ILSVRC 2015 winner (3.6% top 5 error)

Microsoft  
Research

## MSRA @ ILSVRC & COCO 2015 Competitions

- **1st places in all five main tracks**
  - ImageNet Classification: *“Ultra-deep”* (quote Yann) **152-layer** nets
  - ImageNet Detection: **16%** better than 2nd
  - ImageNet Localization: **27%** better than 2nd
  - COCO Detection: **11%** better than 2nd
  - COCO Segmentation: **12%** better than 2nd

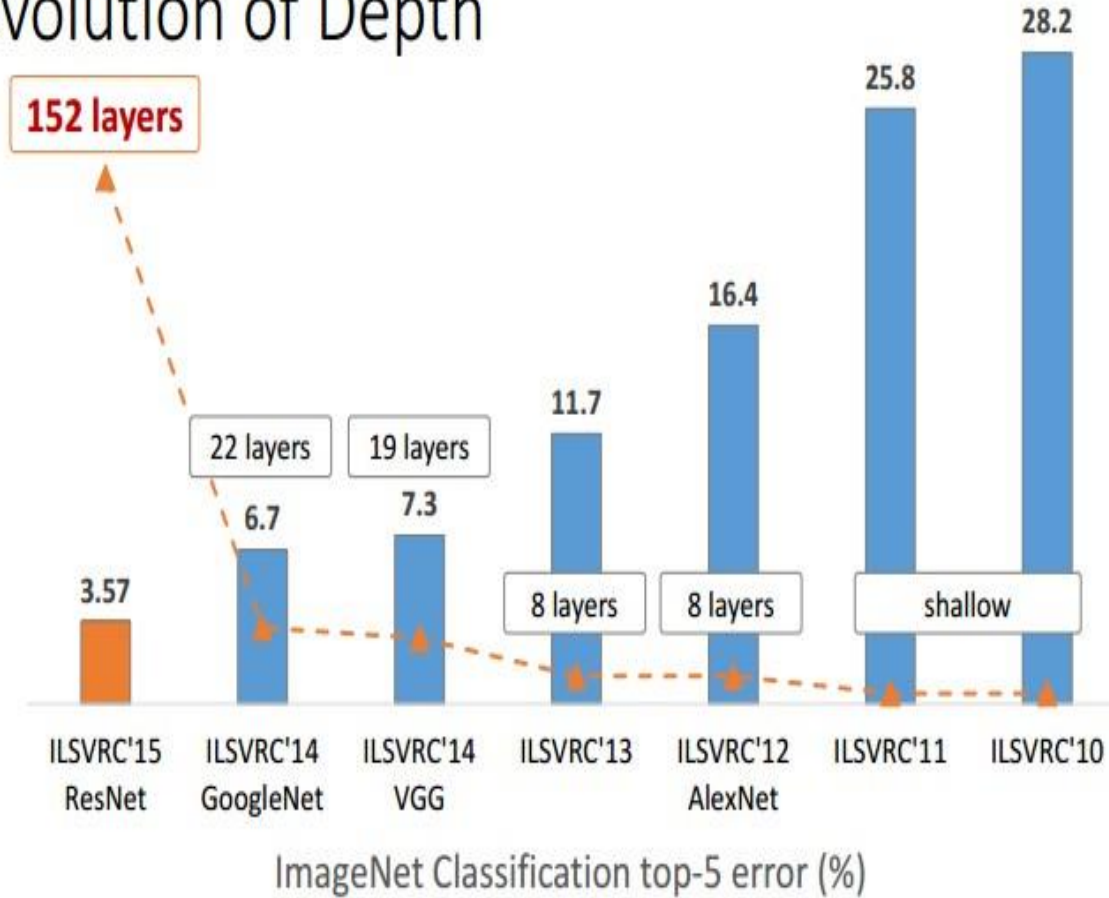
\*improvements are relative numbers

 Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. “Deep Residual Learning for Image Recognition”. arXiv 2015.

Slide from Kaiming He’s recent presentation

<https://www.youtube.com/watch?v=1PGLj-uKT1w>

# Revolution of Depth



ImageNet Classification top-5 error (%)



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

(slide from Kaiming He's recent presentation)

# Summary

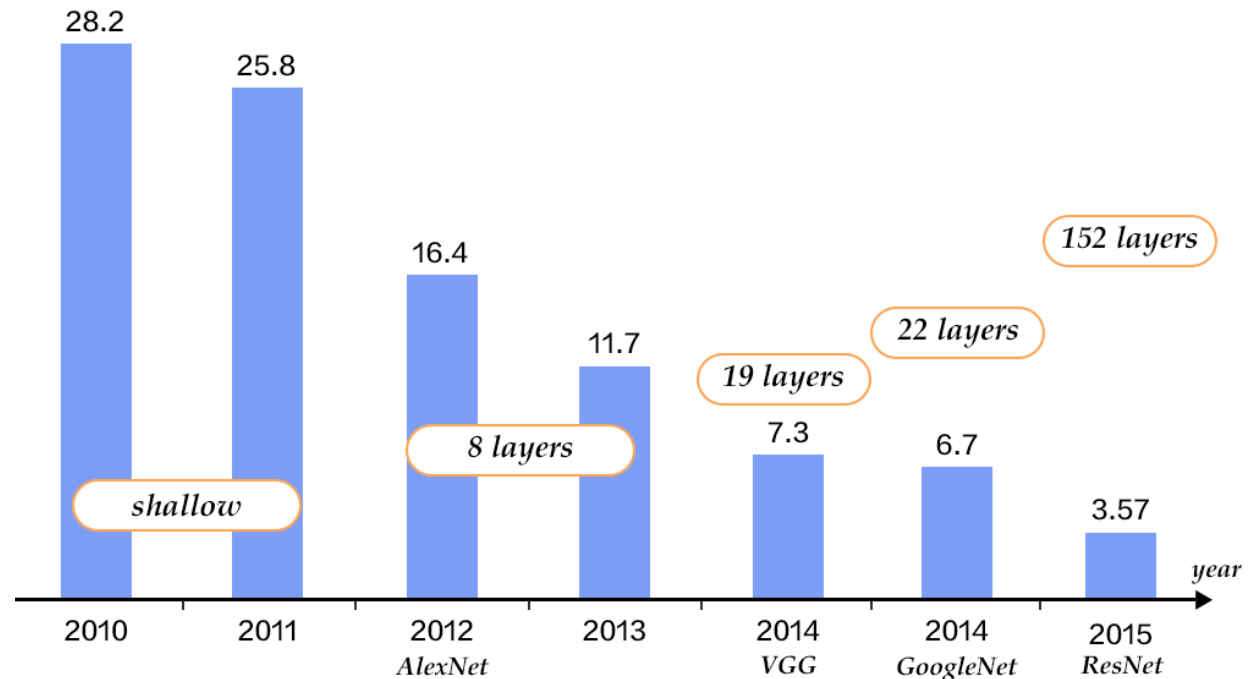
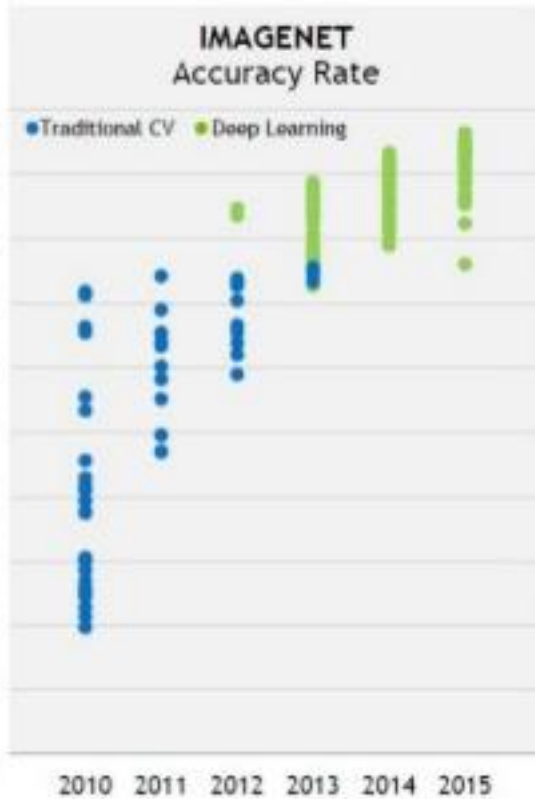
- ConvNets stack CONV, POOL, FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Typical architectures look like

**[(CONV-RELU)\*N-POOL?]\*M-(FC-RELU)\*K,SOFTMAX**

where N is usually up to ~5, M is large,  $0 \leq K \leq 2$ .

- but recent advances such as ResNet/GoogLeNet challenge this paradigm

# ILVRC over the years



<https://image.slidesharecdn.com/goto-ai-pastpresentandfuture-160801162251/95/artificial-intelligence-past-present-and-future-43-638.jpg?cb=1470068690>

[http://book.paddlepaddle.org/03.image\\_classification/image/ilsvrc.png](http://book.paddlepaddle.org/03.image_classification/image/ilsvrc.png)