

Template Matching and Introduction to CNNs

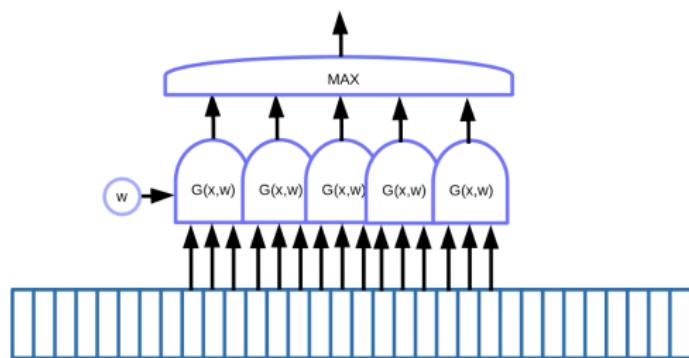
Tirtharaj Dash

Dept. of CS & IS and APPCAIR
BITS Pilani, Goa Campus

September 23, 2021

Motif detection in sequential data I

- Motif detection means to find some motifs in sequential data like keywords in speech or text.
- Idea: Use a sliding window on data, which moves the weight-sharing function to detect a particular motif (i.e. a particular sound in speech signal), and the outputs (i.e. a score) goes into a maximum function.



(source: <https://atcold.github.io/pytorch-Deep-Learning/>)

Motif detection in sequential data II

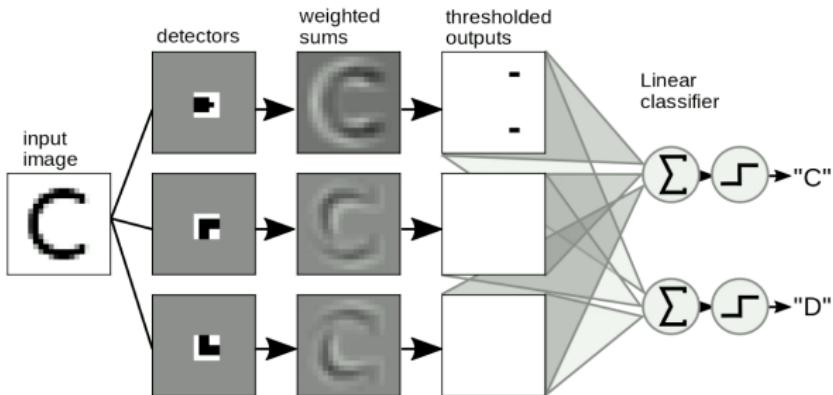
- There are 5 functions $G()$ s:
 - Each function behaves as the template to detect the motif in the input;
 - Each share the same parameters w
- For backprop, the five gradients are summed up to update w .

Motif detection in images I

- Slide a “templates” (or multiple templates) over an image to detect the shapes independent of position and distortion of the shapes.
- E.g. distinguish between two images: one of character ‘C’ and another ‘D’.
- ‘C’ has two endpoints; ‘D’ has two corners.
- So there are three templates: endpoint template, corner templates.

Motif detection in images II

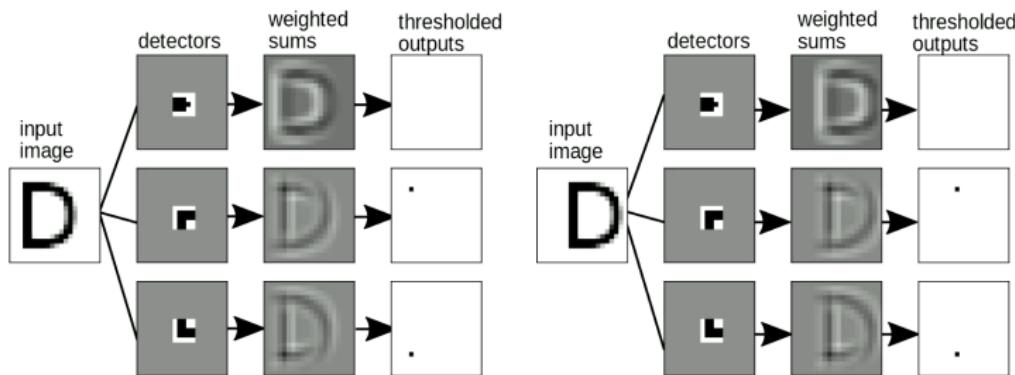
- 'C' (image) has no corners but has endpoints so that output for 'C' is activated.



(source: <https://atcold.github.io/pytorch-Deep-Learning/>)

Motif detection in images III

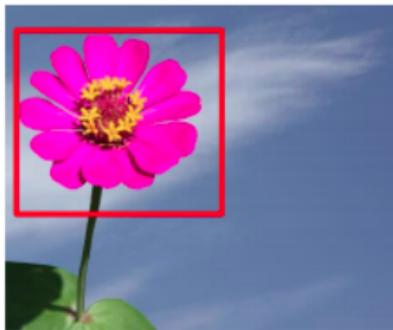
- It is also important that our “template matching” should be shift-invariant: When we shift the input, the output (i.e. the letter detected) shouldn't change. This can be solved with weight sharing transformation.



(source: <https://atcold.github.io/pytorch-Deep-Learning/>)

Motif detection in images IV

Another example:



Credit: The illustrations with the flower image are adapted from Bhiksha Raj's [DL course](#).

Motif detection in images V

The output must be correct regardless of the precise location of the target object (here, flower).

What weight sharing does is a kind of scanning over the original image with the weight template.

Motif detection in images VI

- This was the approach for long time: hand-crafted templates.
- Questions:
 - How to design ‘templates’ automatically?
 - How many templates to be designed?
 - What about rotations, noisy images etc?
 - Can we use neural nets to learn (or construct) these?
 - ...

Discrete Convolution I

Convolution:

- In 1-dimensional case:

$$y_i = \sum_j w_j x_{i-j}$$

- The i -th output is computed as the dot product between the reversed w and a window of the same size in x .
- To compute the full output, start the window at the beginning, shift this window by one entry each time and repeat until x is exhausted.

Discrete Convolution II

Cross-correlation:

- Most common deep learning libraries don't reverse w so it becomes cross-correlation:

$$y_i = \sum_j w_j x_{i+j}$$

- It is however a method of convention (with due respect to signal processing guys :-)) based on how you read w from memory.

Discrete Convolution III

Higher-dimensional convolution:

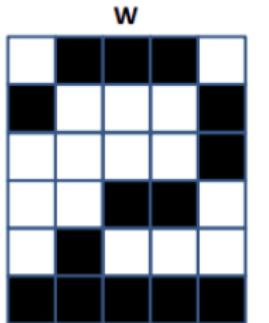
- For two dimensional inputs such as images, we make use of the two dimensional version of convolution:

$$y_{ij} = \sum_{kl} w_{kl} x_{i+k, j+l}$$

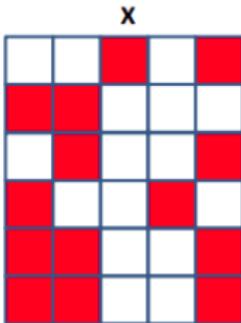
- Can be extended beyond two dimensions to three or four dimensions.
- Here w is called the convolution kernel.

Discrete Convolution IV

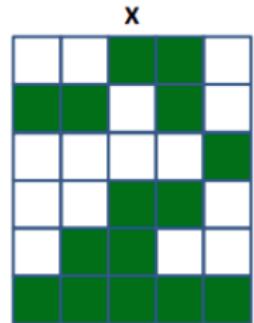
Example: Checking if an image is of the digit '2'



$$y = \begin{cases} 1 & \text{if } \sum_i w_i x_i \geq T \\ 0 & \text{else} \end{cases}$$



Correlation = 0.57



Correlation = 0.82

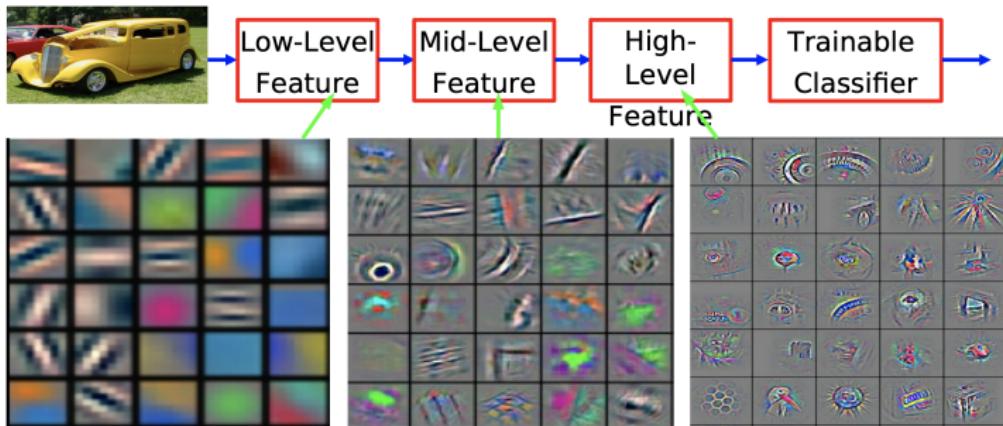


Deep Convolutional Neural Network I

- DNNs: `repeat(linear transofrmation, non-linearity)`
- CNNs: `repeat(convolution, non-linearity, [pooling])`
- Reason for stacking such layers: build a hierarchical representation of the data.

Deep Convolutional Neural Network II

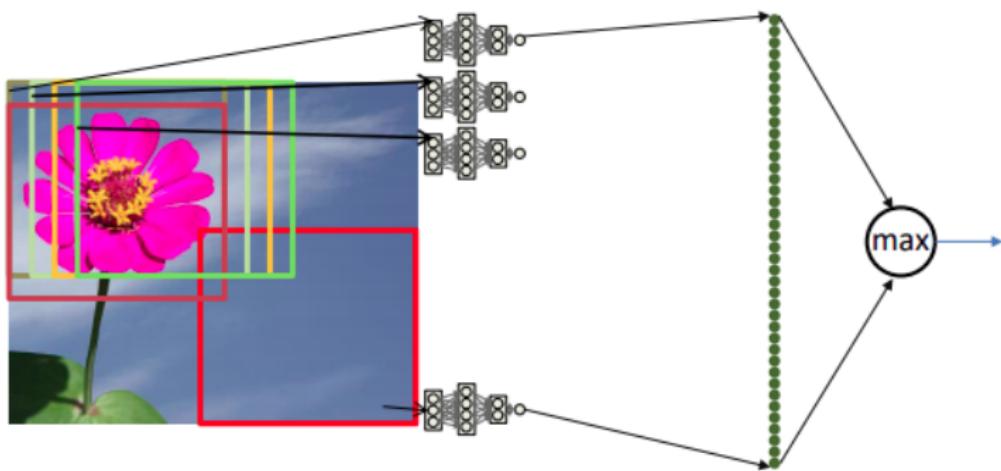
- Why would we want to capture the hierarchical representation of the world? — Because the world we live in is compositional.



Deep Convolutional Neural Network III

Scanning an image to recognise a flower:

- This is how it is usually done:
 - We get one classification output per scanned location
 - The score output by the MLP
 - Look at the maximum value
 - Or, pass it through an MLP

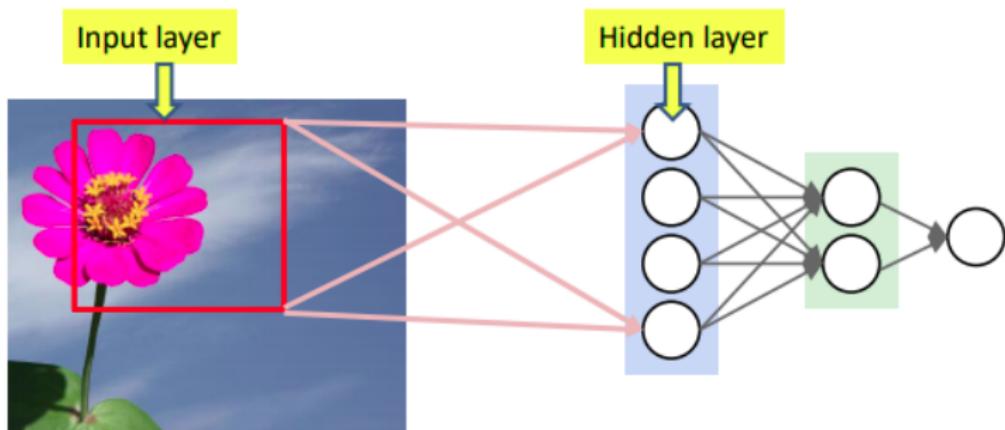


Deep Convolutional Neural Network IV

- The entire operation can be viewed as a single giant network
 - Composed of many “subnets” (one per window)
 - With one key feature: all subnets are identical
 - These are shared parameter networks: all are searching for the same pattern
 - Any update of the parameters of one copy of the subnet must equally update all copies

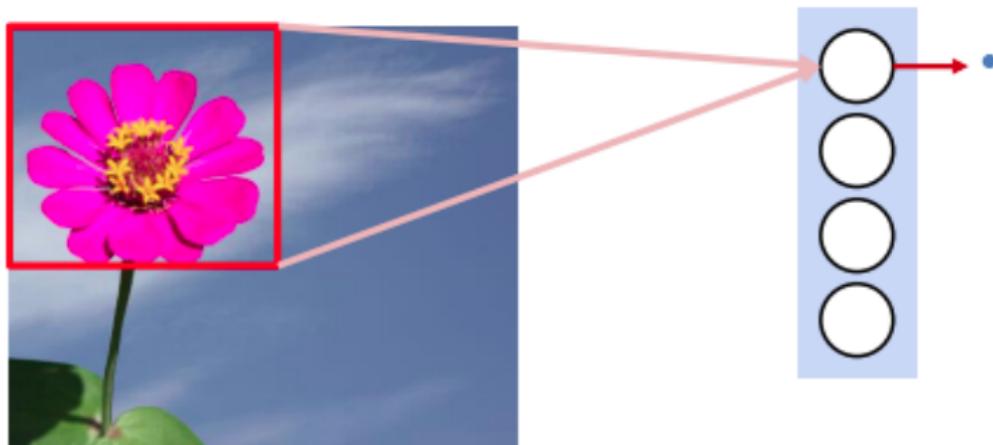
Deep Convolutional Neural Network V

- If we closely look at any small network:

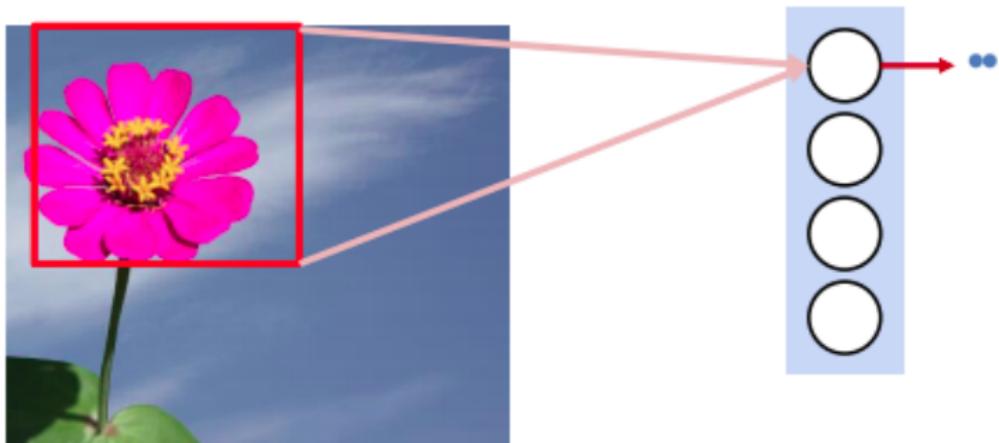


- At each position of the box, the perceptron is evaluating the part of the picture in the box as part of the classification for that region.
- Outputs are then arranged as we will see next.

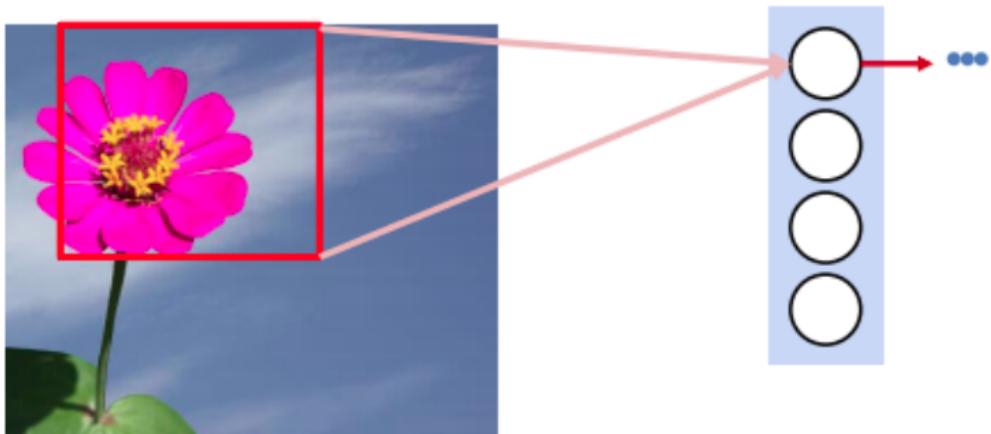
Deep Convolutional Neural Network VI



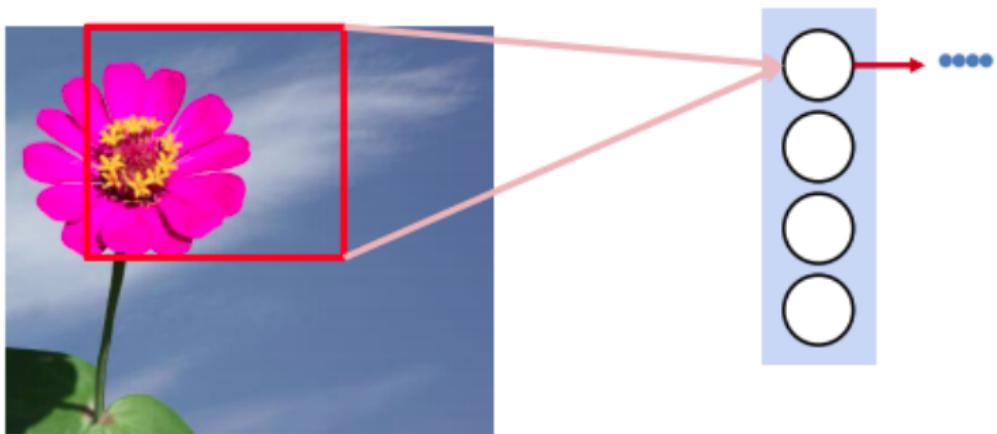
Deep Convolutional Neural Network VII



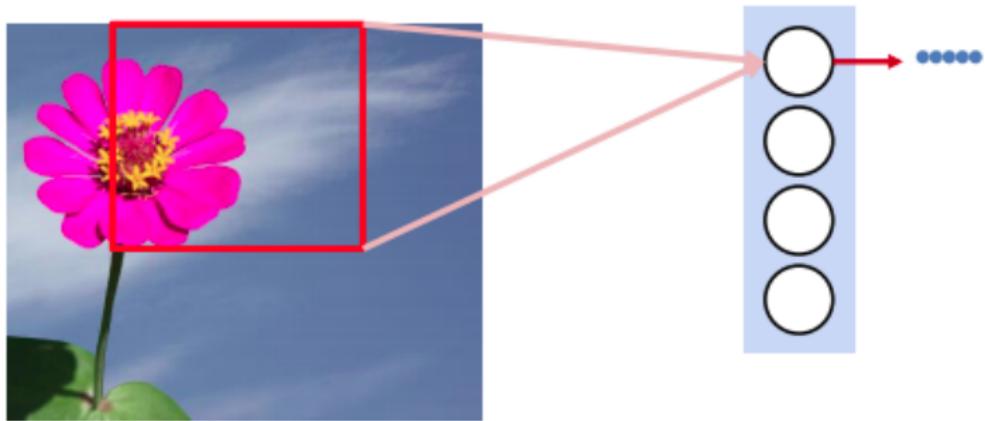
Deep Convolutional Neural Network VIII



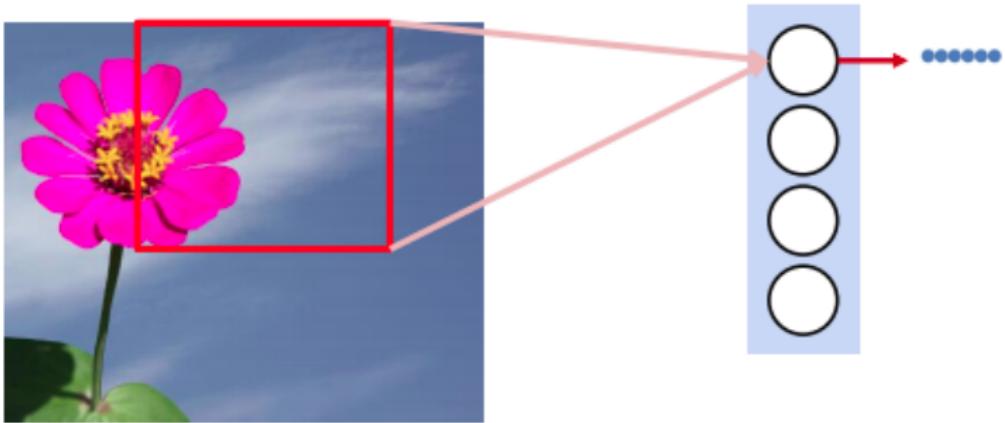
Deep Convolutional Neural Network IX



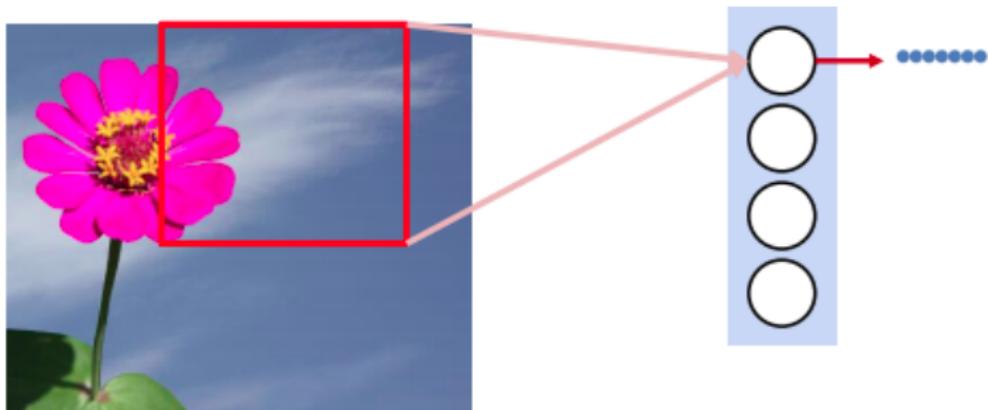
Deep Convolutional Neural Network X



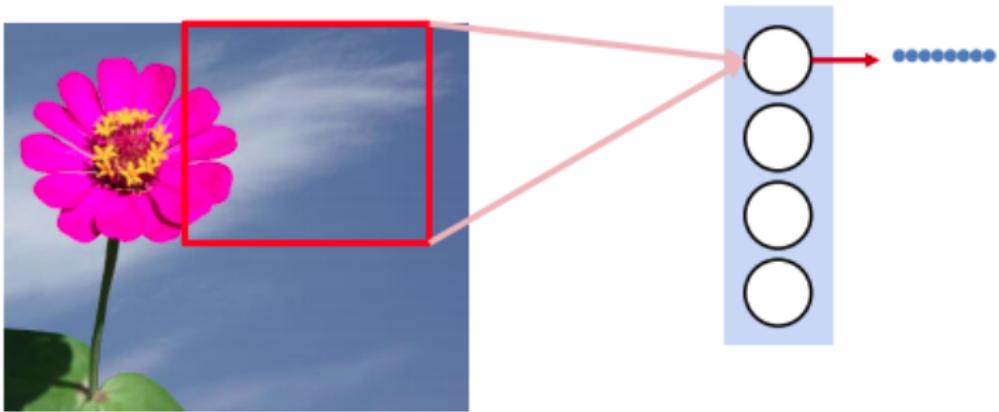
Deep Convolutional Neural Network XI



Deep Convolutional Neural Network XII



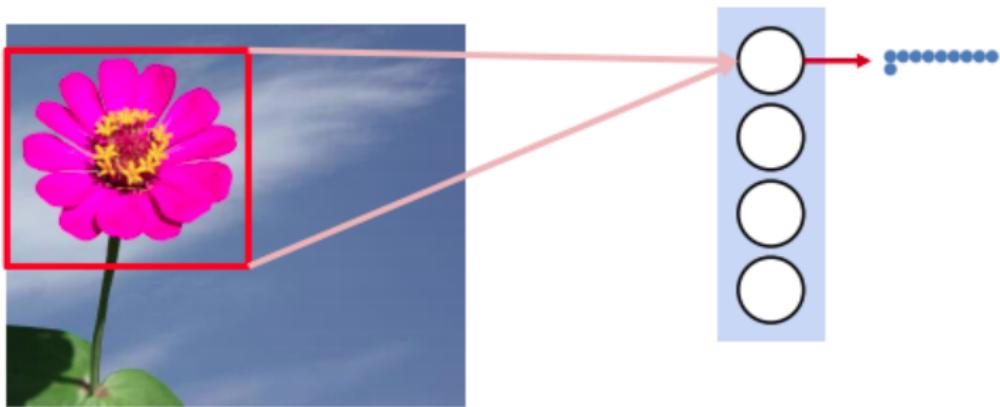
Deep Convolutional Neural Network XIII



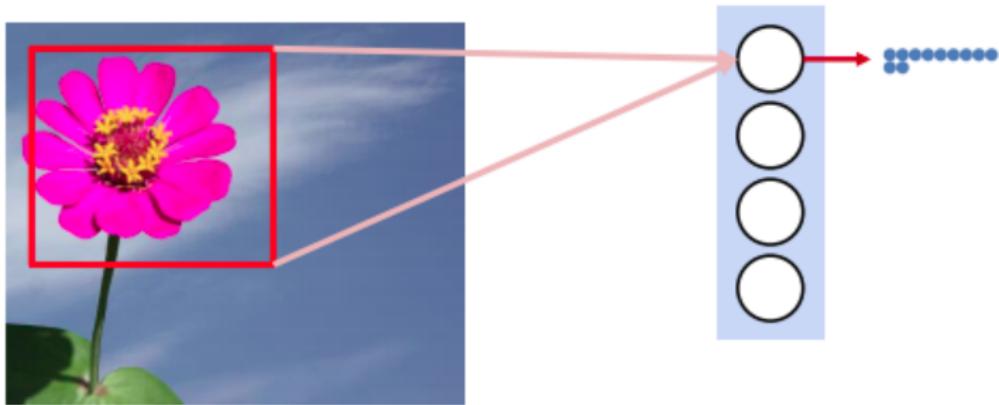
Deep Convolutional Neural Network XIV



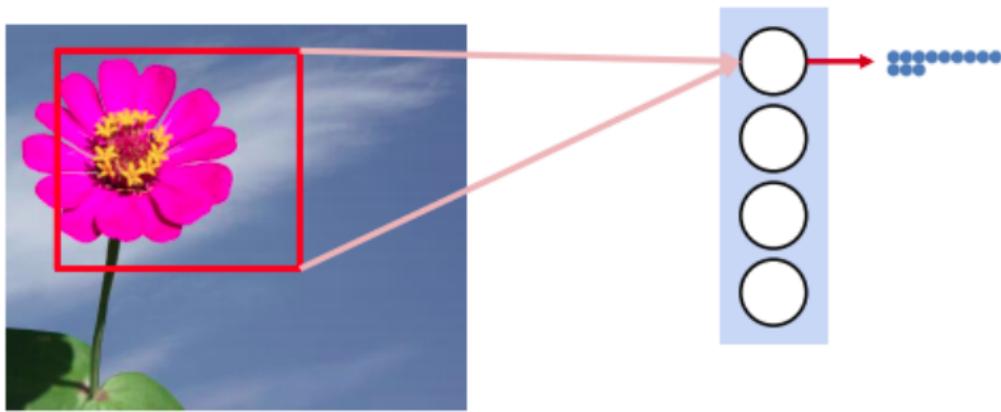
Deep Convolutional Neural Network XV



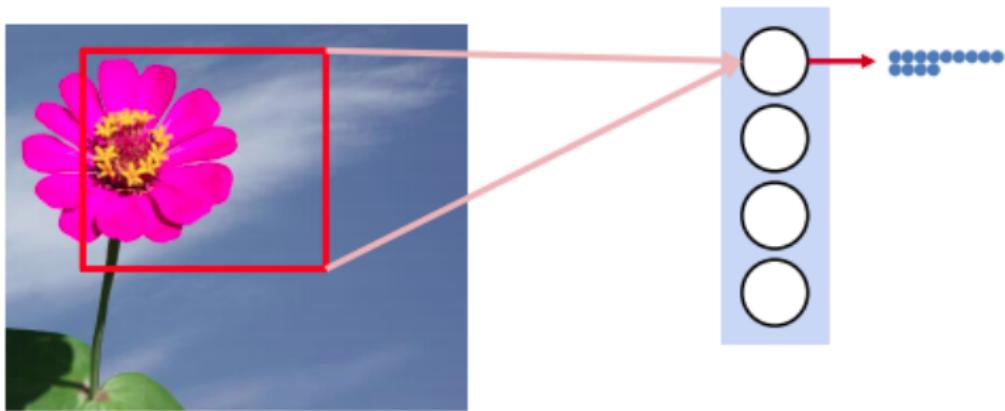
Deep Convolutional Neural Network XVI



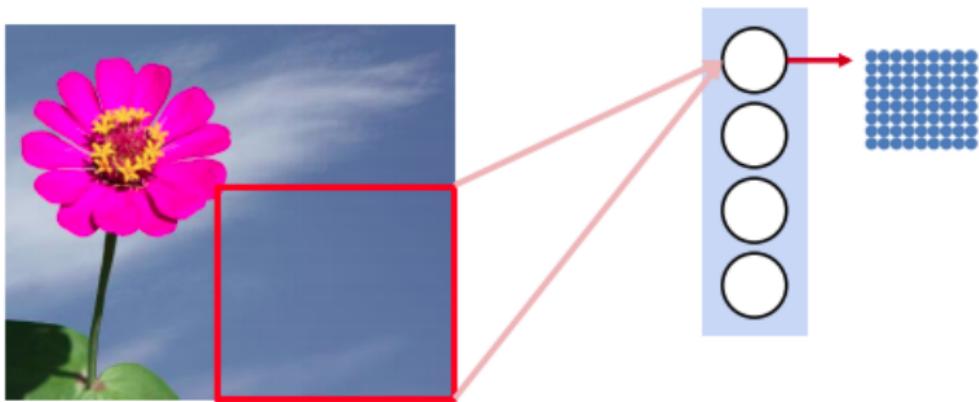
Deep Convolutional Neural Network XVII



Deep Convolutional Neural Network XVIII



Deep Convolutional Neural Network XIX



The same idea can be recursively applied:

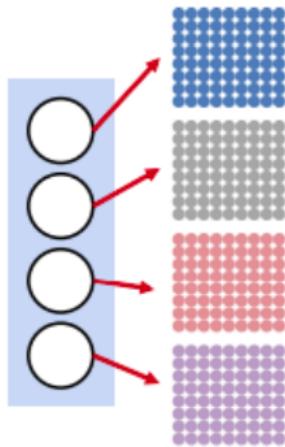
- The second level neurons too are “scanning” the rectangular outputs of the first-level neurons
- But, they are jointly scanning **multiple** “pictures”

Deep Convolutional Neural Network XXI

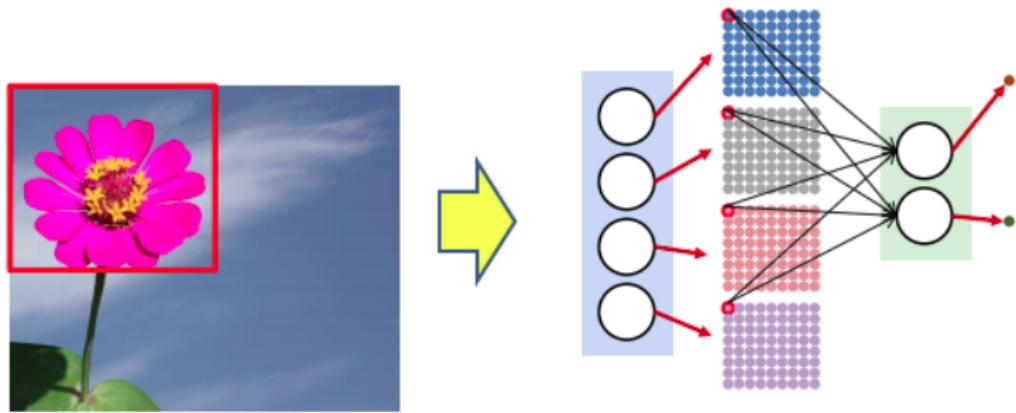


Each location in the output of the second level neuron considers the corresponding locations from the outputs of all the first-level neurons

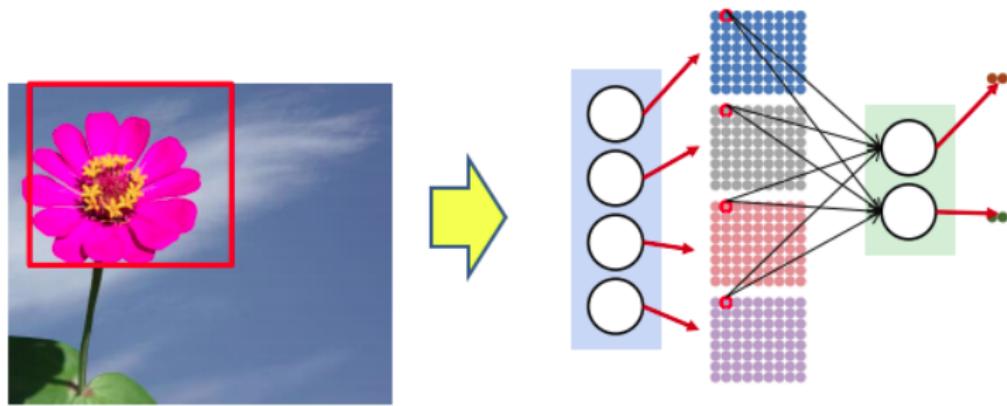
Deep Convolutional Neural Network XXII



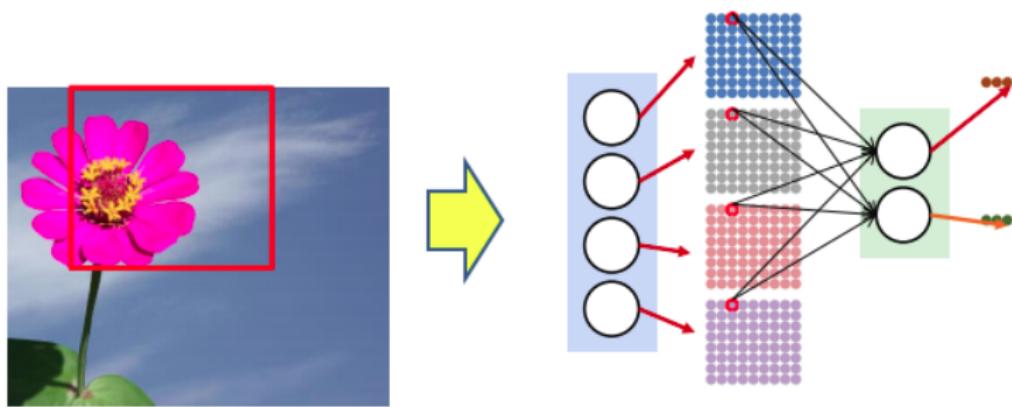
Deep Convolutional Neural Network XXIII



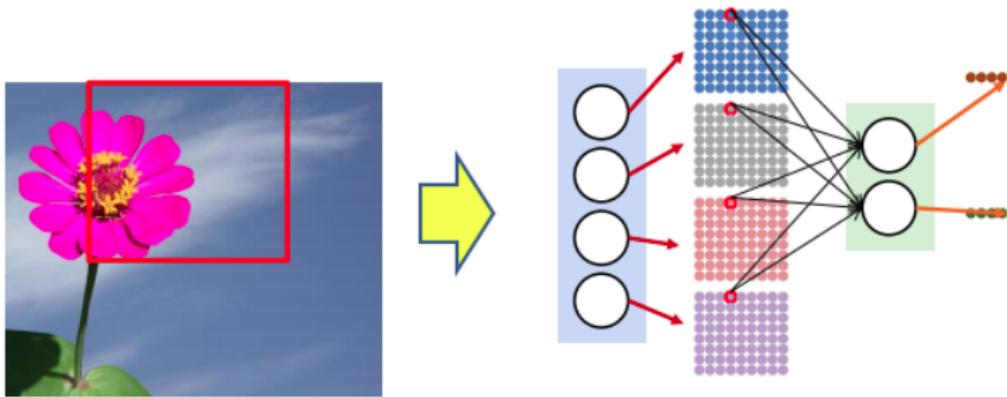
Deep Convolutional Neural Network XXIV



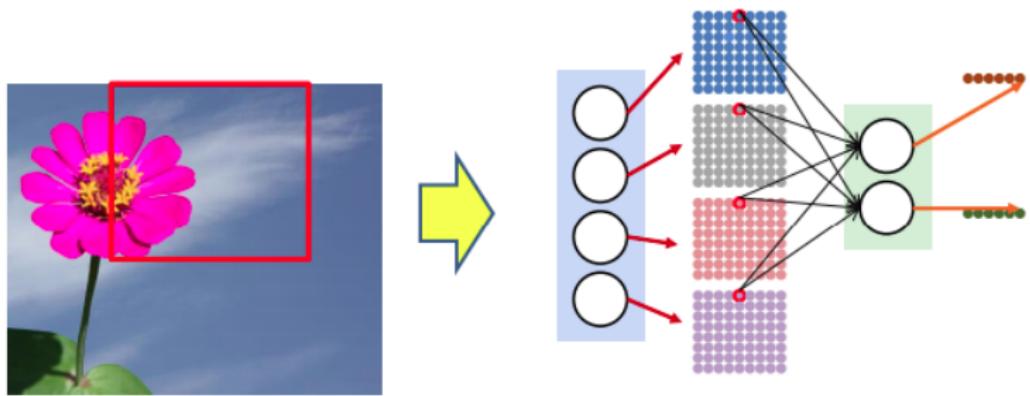
Deep Convolutional Neural Network XXV



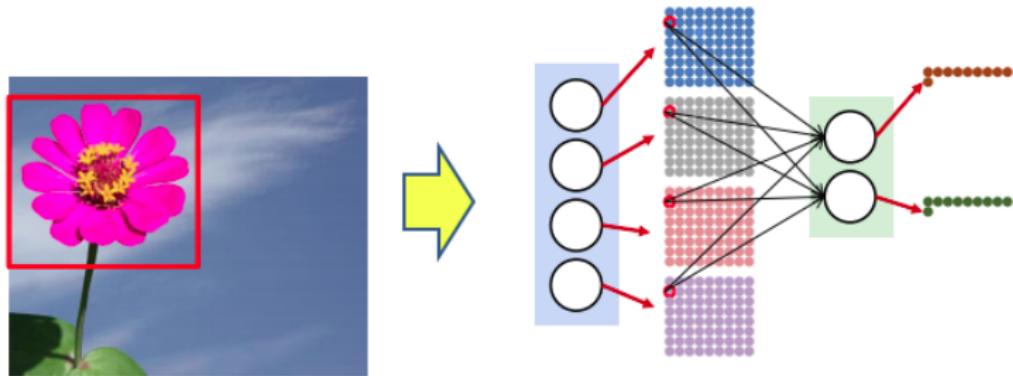
Deep Convolutional Neural Network XXVI



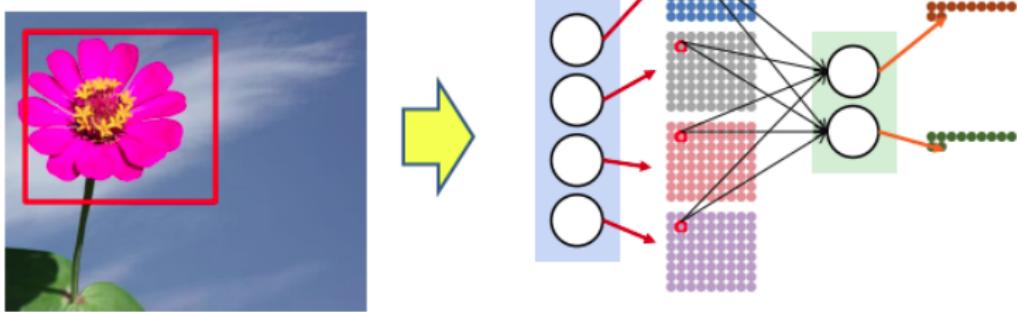
Deep Convolutional Neural Network XXVII



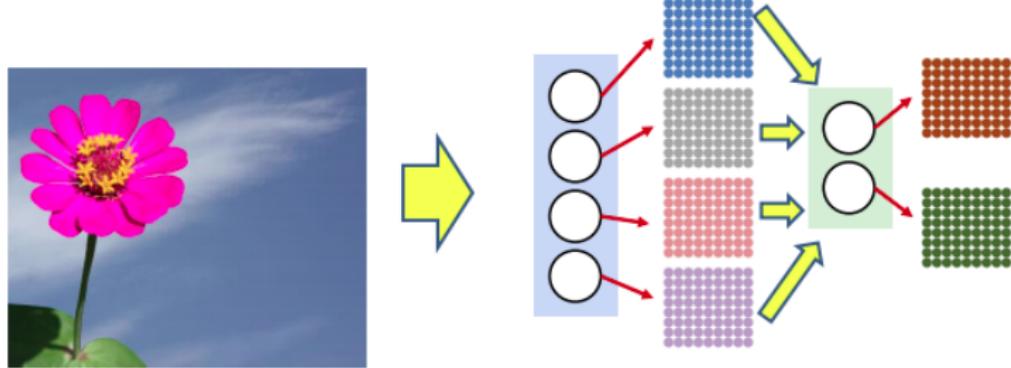
Deep Convolutional Neural Network XXVIII



Deep Convolutional Neural Network XXIX



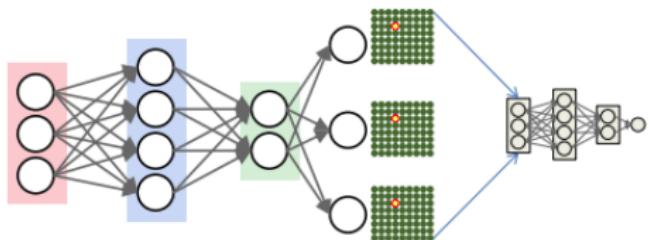
Deep Convolutional Neural Network XXX



- Is there a flower in the picture?
 - The output of the almost-last layer is also a grid/picture
 - The entire grid can be sent into a final neuron that performs a logical “OR” to detect a picture
 - Finds the max output from all positions

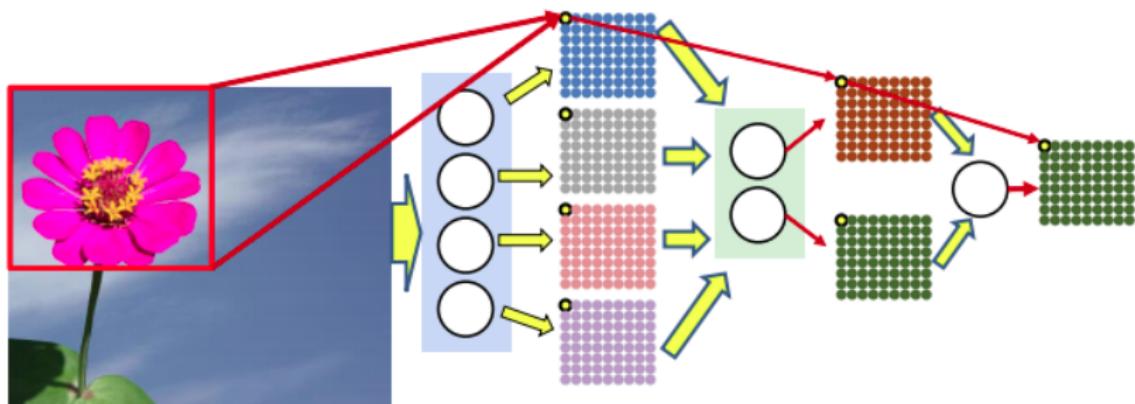
Deep Convolutional Neural Network XXXII

- In more general sense: At each location, the net searches for a flower
 - The entire map of outputs is sent through a follow-up perceptron (or MLP) to determine if there really is a flower in the picture

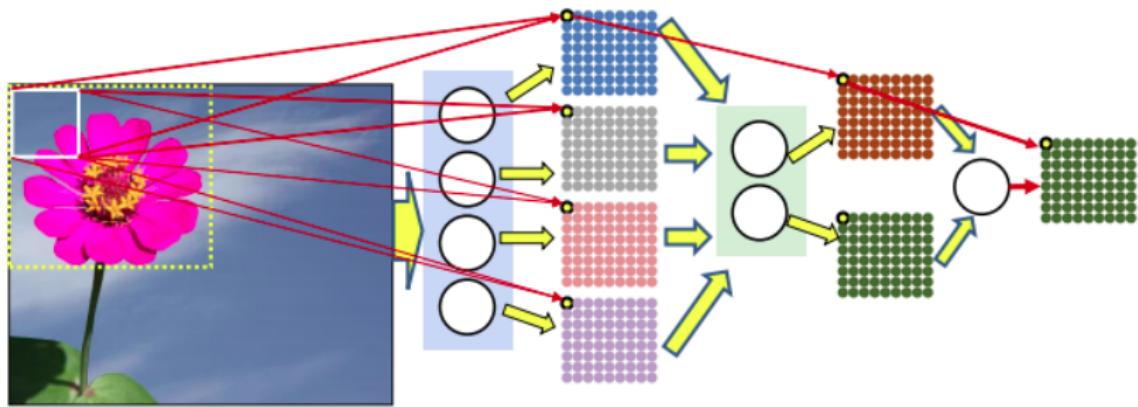


Deep Convolutional Neural Network XXXIII

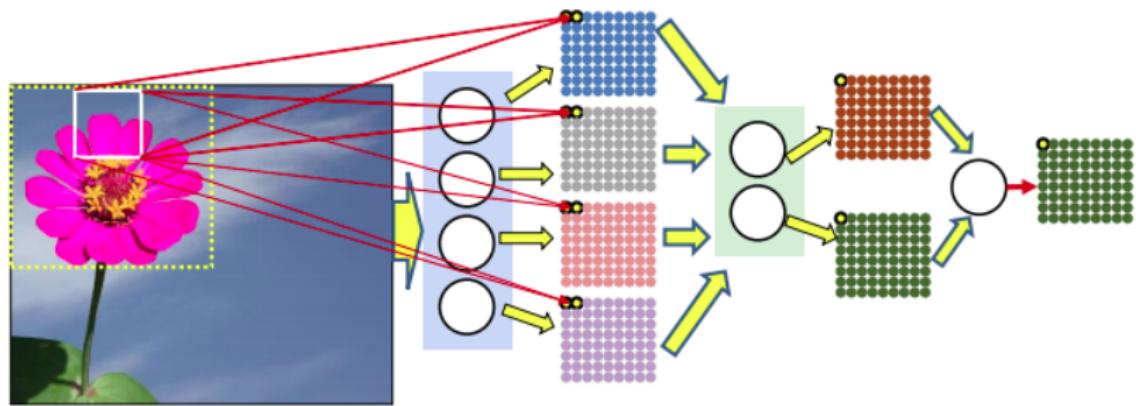
- Scan was actually distributed:



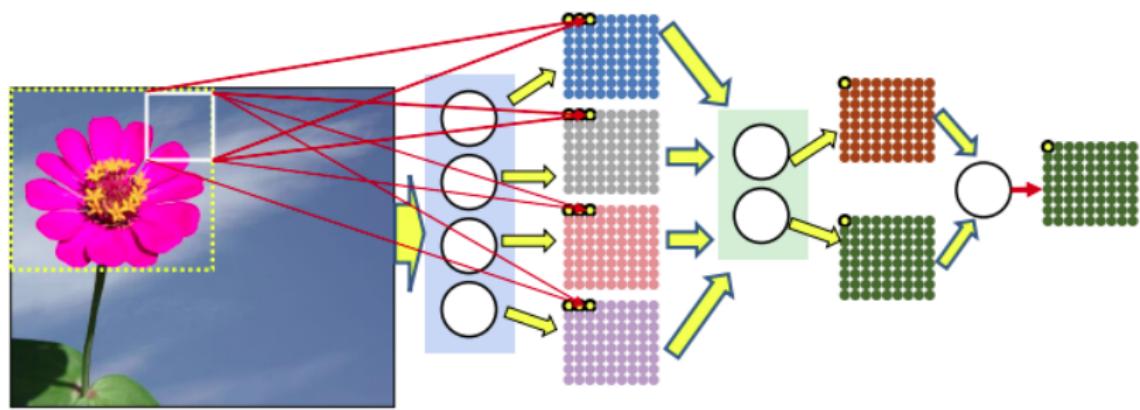
Deep Convolutional Neural Network XXXIV



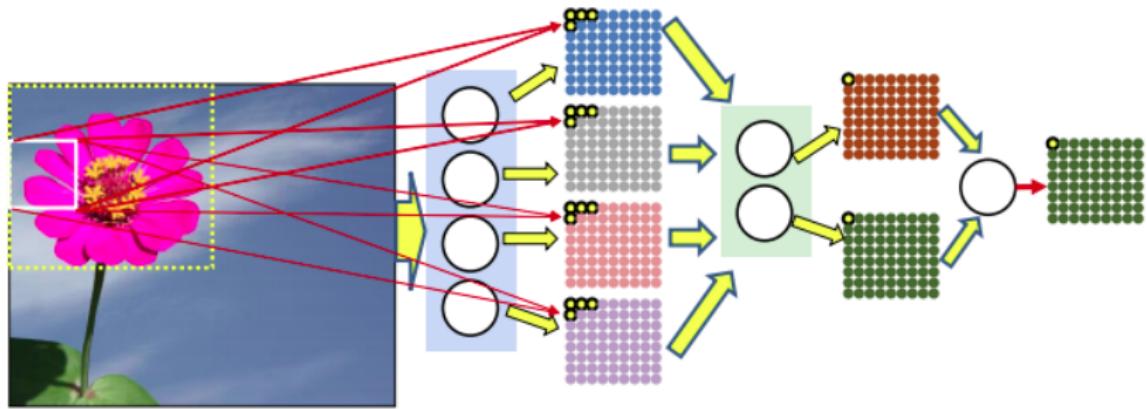
Deep Convolutional Neural Network XXXV



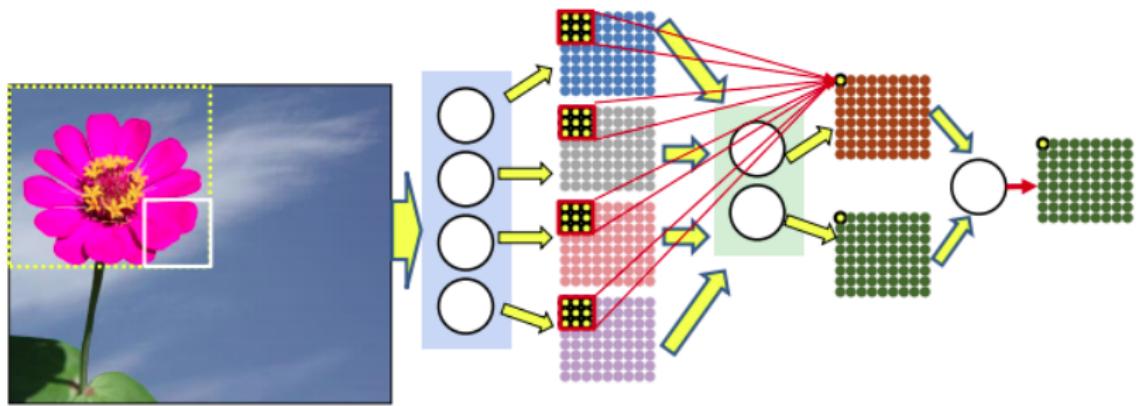
Deep Convolutional Neural Network XXXVI



Deep Convolutional Neural Network XXXVII



Deep Convolutional Neural Network XXXVIII



Generic CNN architecture I

- Normalisation
 - Adjusting whitening (optional)
 - Subtractive methods e.g. average removal, high pass filtering
 - Divisive: local contrast normalisation, variance normalisation
- Filter banks:
 - Increase dimensionality
 - Projection on overcomplete basis
 - Edge detections
- Non-linearities
 - Sparsification
 - Typically Rectified Linear Unit (ReLU)
- Handling jitter: Pooling
 - Aggregating over a feature map
 - Max pooling
 - L_p -norm pooling
 - ...

What CNNs are good for? I

- Locality:
 - There is a strong local correlation between values.
 - Nearby pixels same color; similarity decreases as they move distant
 - The local correlations can help us detect local features (CNNs do this).
 - If we feed the CNN with permuted pixels, it will not perform well at recognizing the input images, while FC will not be affected.
 - The local correlation justifies local connections.

What CNNs are good for? II

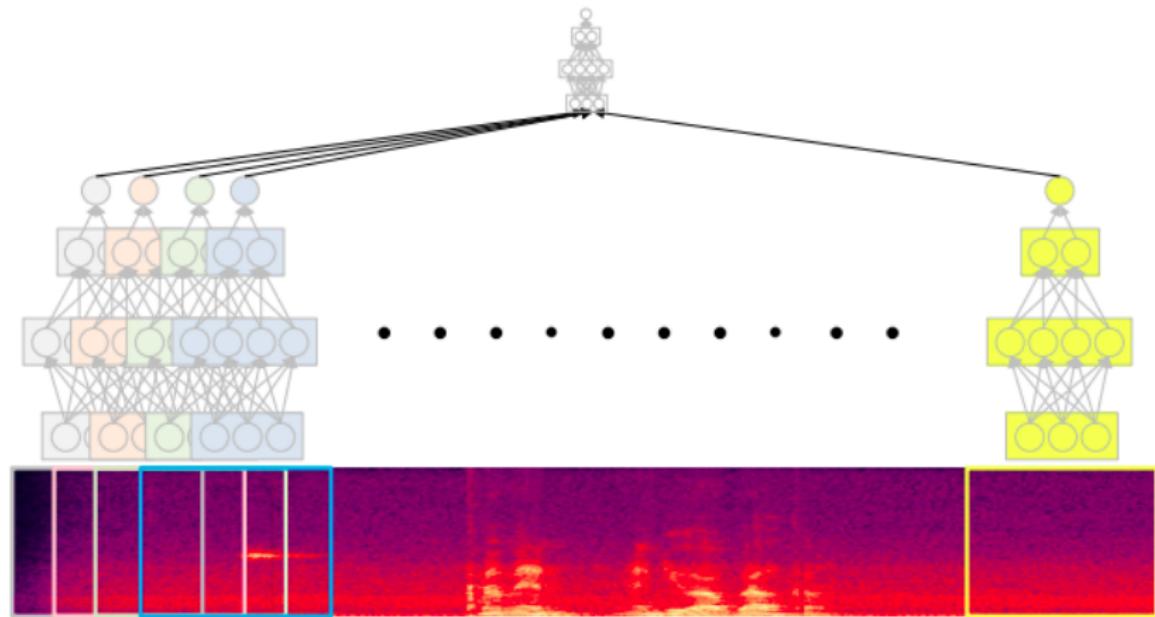
- Stationarity:
 - The features are essential and can appear anywhere on the image (shared weight and pooling)
 - Moreover, statistical signals are uniformly distributed, which means we need to repeat the feature detection for every location on the input image.

What CNNs are good for? III

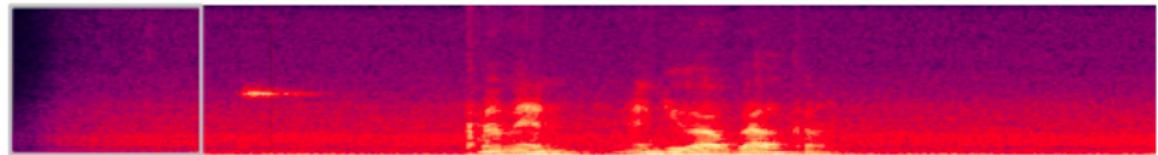
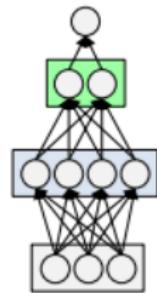
- Compositionality:
 - The natural images are compositional, meaning the features compose an image in a hierarchical manner.
 - Justifies use of multiple layers
- Success in various application domains: videos, images, texts, and speech recognition.

A simple case-study: Speech Recognition I

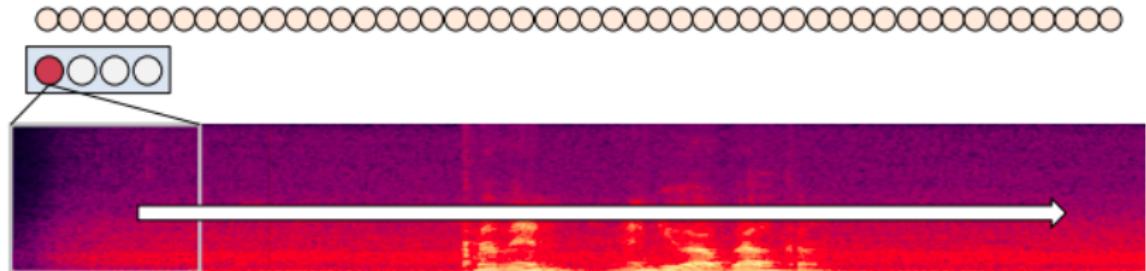
The spectrographic time-frequency components are fed as input to the 1D CNN.



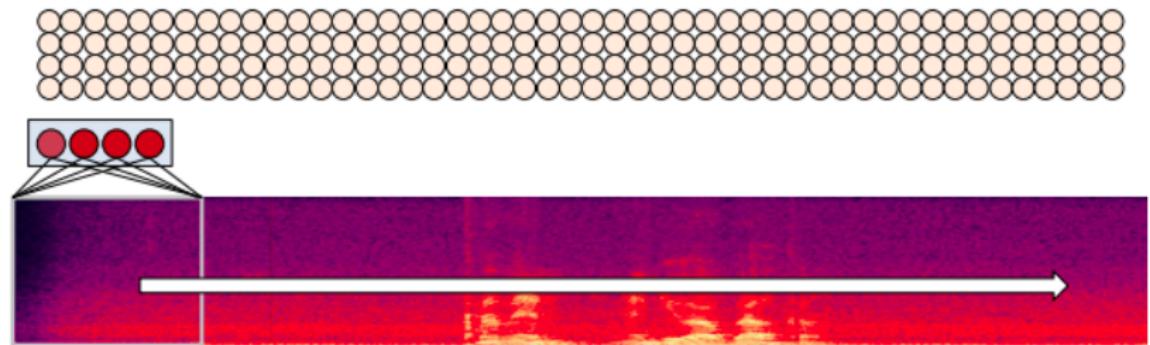
A simple case-study: Speech Recognition II



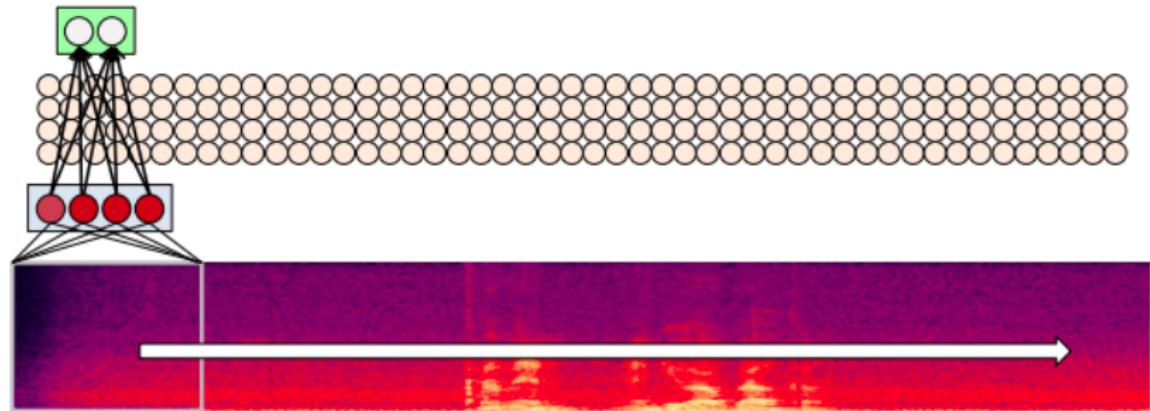
A simple case-study: Speech Recognition III



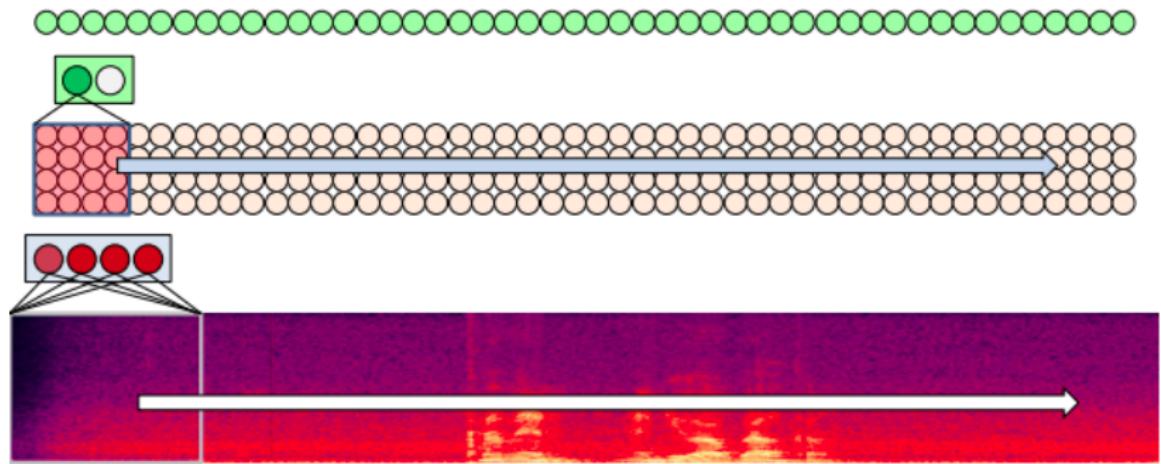
A simple case-study: Speech Recognition IV



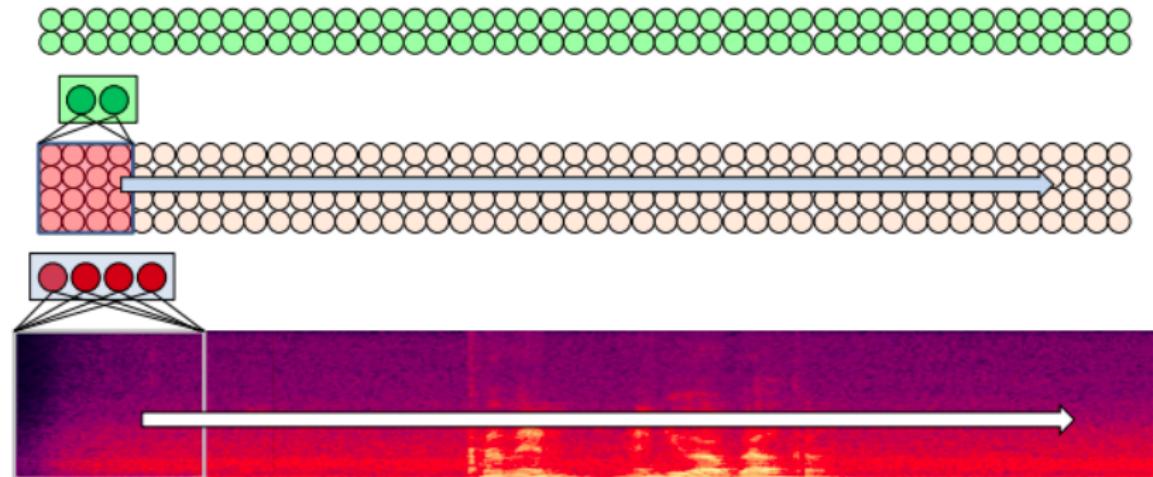
A simple case-study: Speech Recognition V



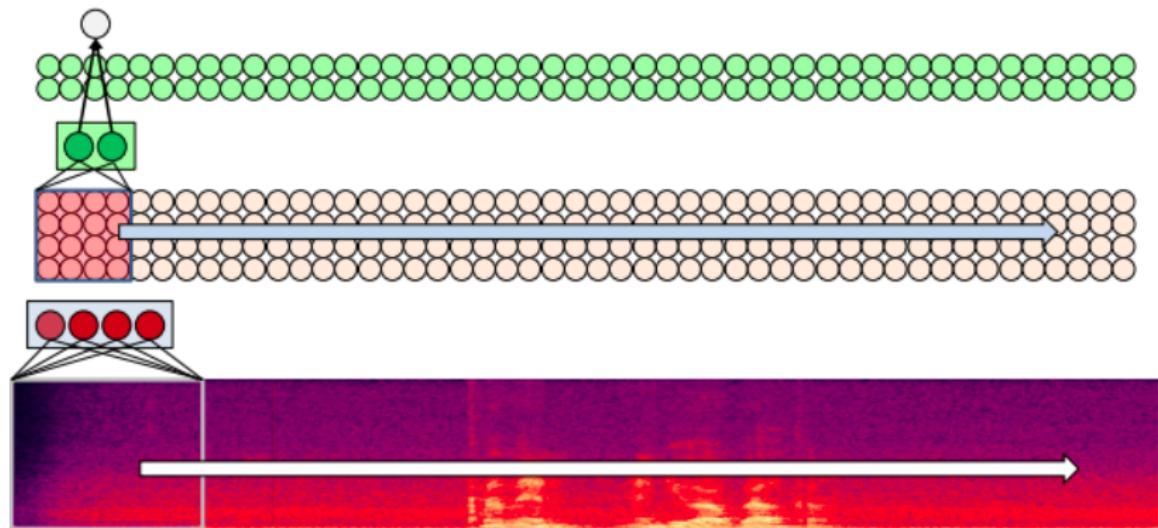
A simple case-study: Speech Recognition VI



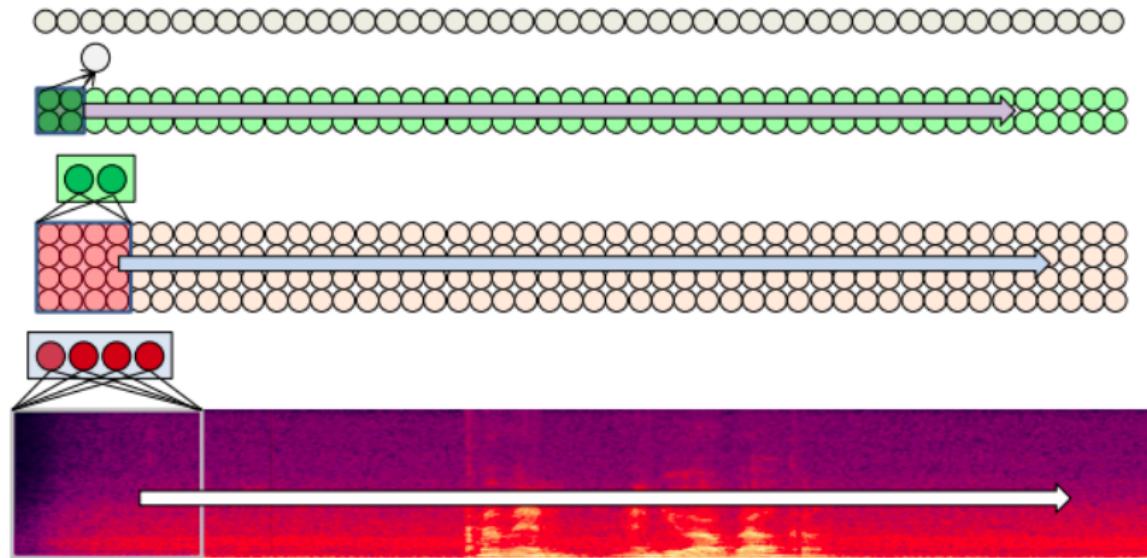
A simple case-study: Speech Recognition VII



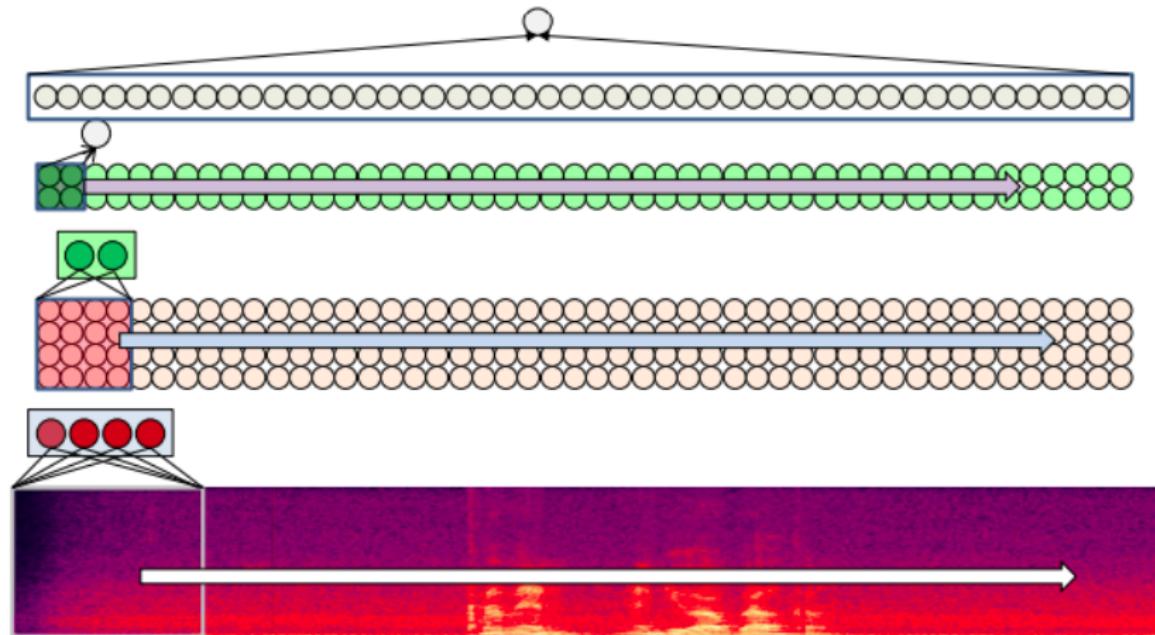
A simple case-study: Speech Recognition VIII



A simple case-study: Speech Recognition IX



A simple case-study: Speech Recognition X



A simple case-study: Speech Recognition XI

E.g. A final perceptron (or MLP) to aggregate evidence: “Does this recording have the target word?”