

Convolutional Neural Network

Hung-yi Lee

Can the network be simplified by
considering the properties of images?

What does machine learn?

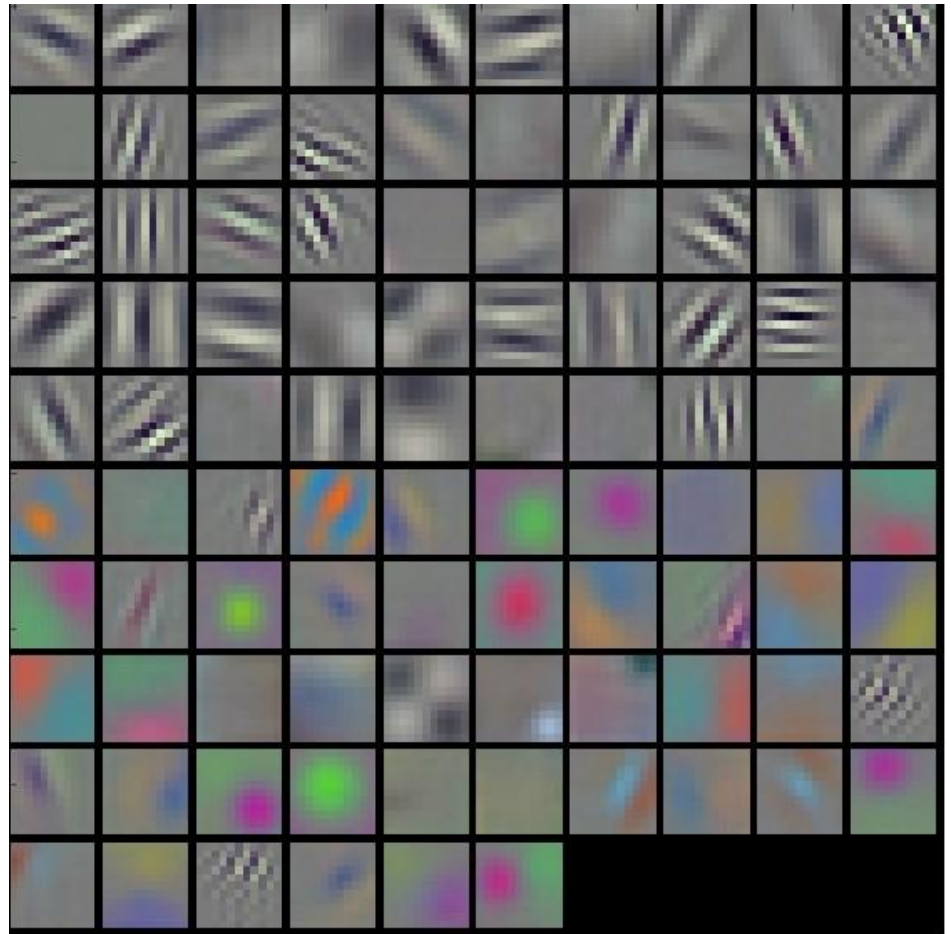


<http://newsneakernews.wpengine.netdna-cdn.com/wp-content/uploads/2016/11/rihanna-puma-creeper-velvet-release-date-02.jpg>

First Convolution Layer

- Typical-looking filters on the trained first layer

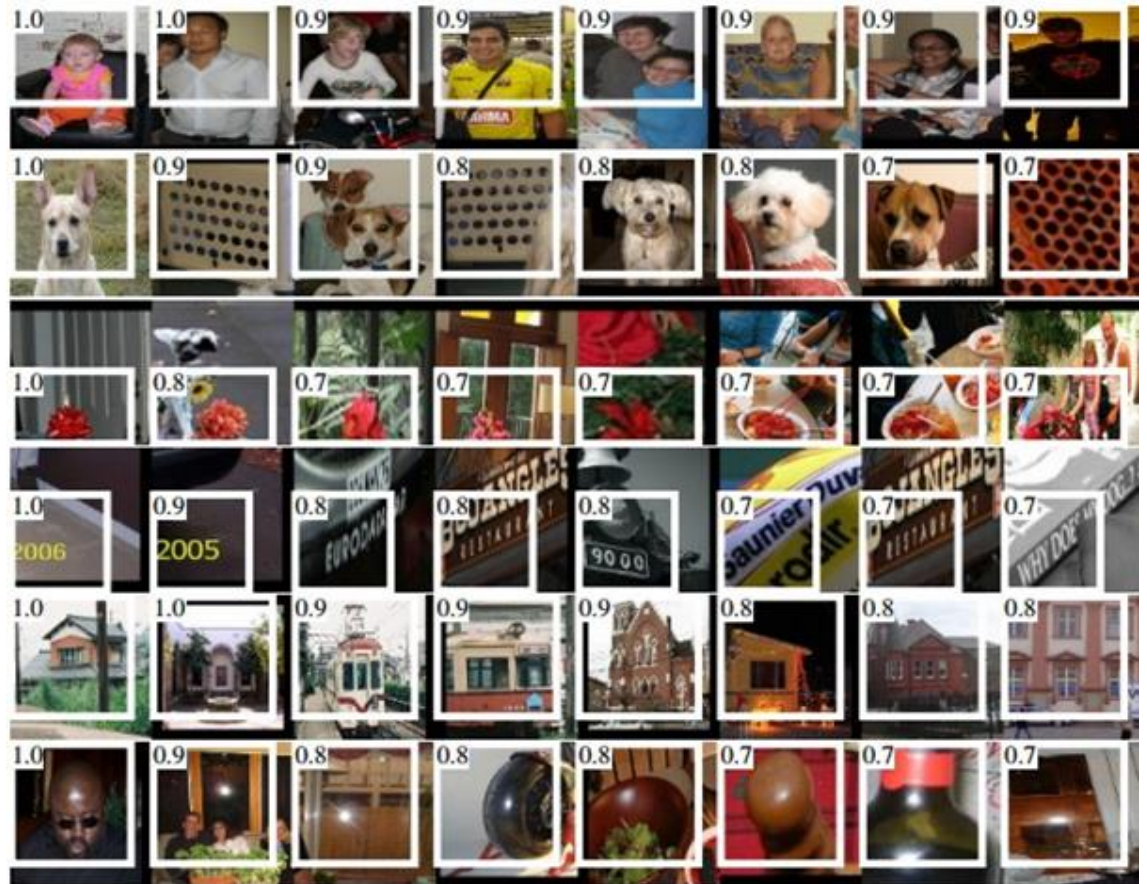
11 x 11
(AlexNet)



<http://cs231n.github.io/understanding-cnn/>

How about higher layers?

- Which images make a specific neuron activate



Ross Girshick, Jeff
Donahue, Trevor
Darrell, Jitendra Malik, "Rich
feature hierarchies for accurate
object detection and semantic
segmentation", CVPR, 2014

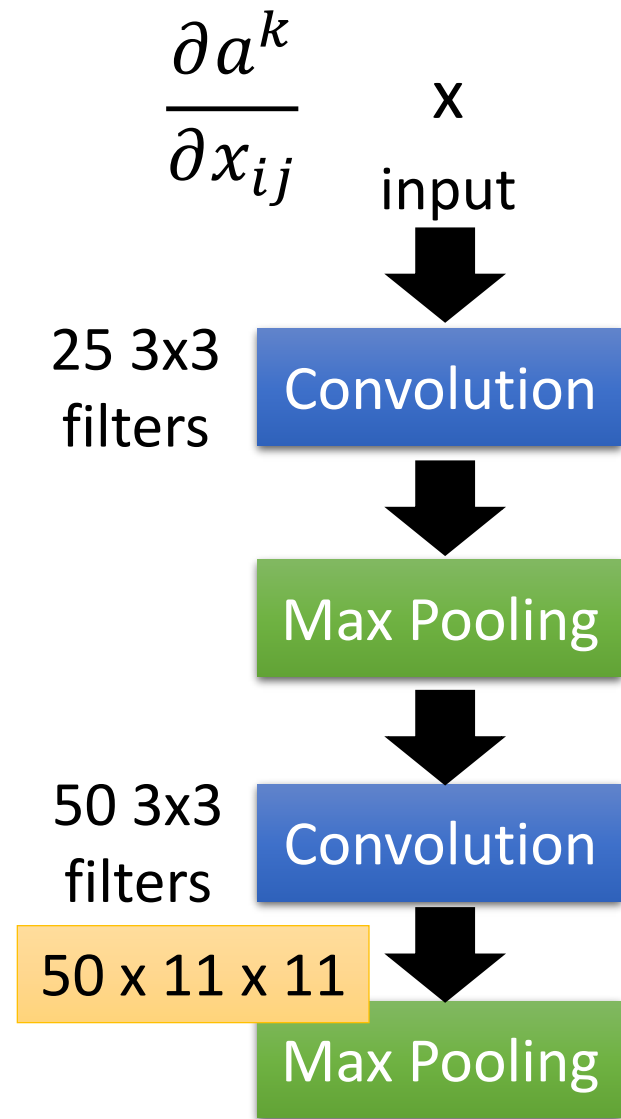
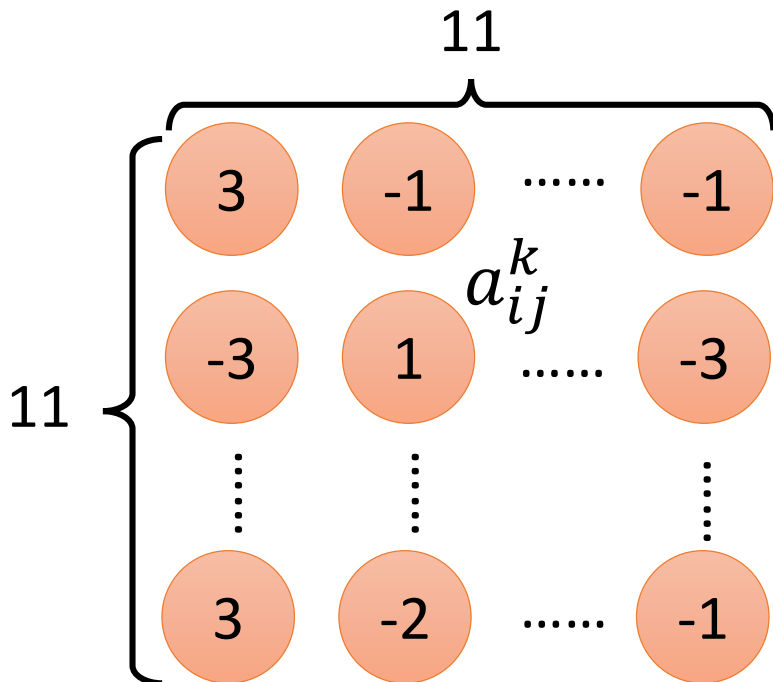
What does CNN learn?

The output of the k-th filter is a 11 x 11 matrix.

Degree of the activation of the k-th filter:

$$a^k = \sum_{i=1}^{11} \sum_{j=1}^{11} a_{ij}^k$$

$x^* = \arg \max_x a^k$ (gradient ascent)



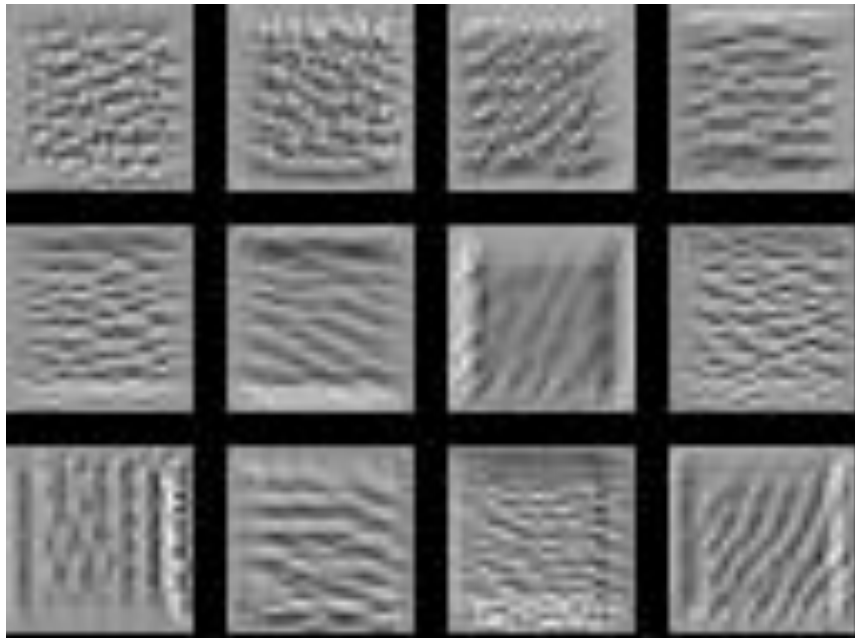
What does CNN learn?

The output of the k-th filter is a 11 x 11 matrix.

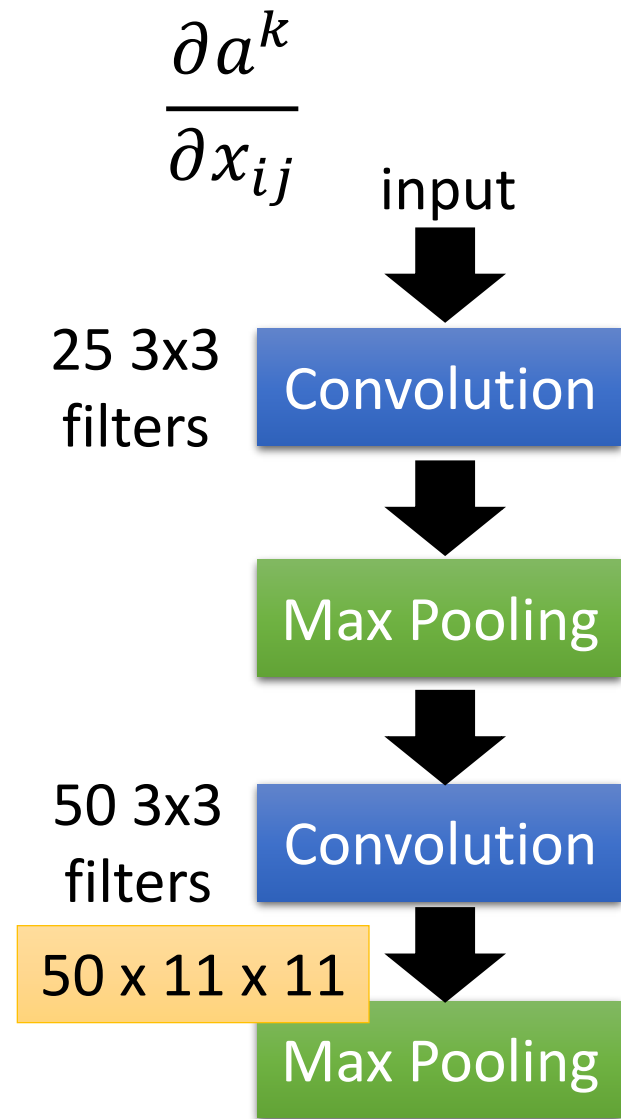
Degree of the activation of the k-th filter:

$$a^k = \sum_{i=1}^{11} \sum_{j=1}^{11} a_{ij}^k$$

$x^* = \arg \max_x a^k$ (gradient ascent)



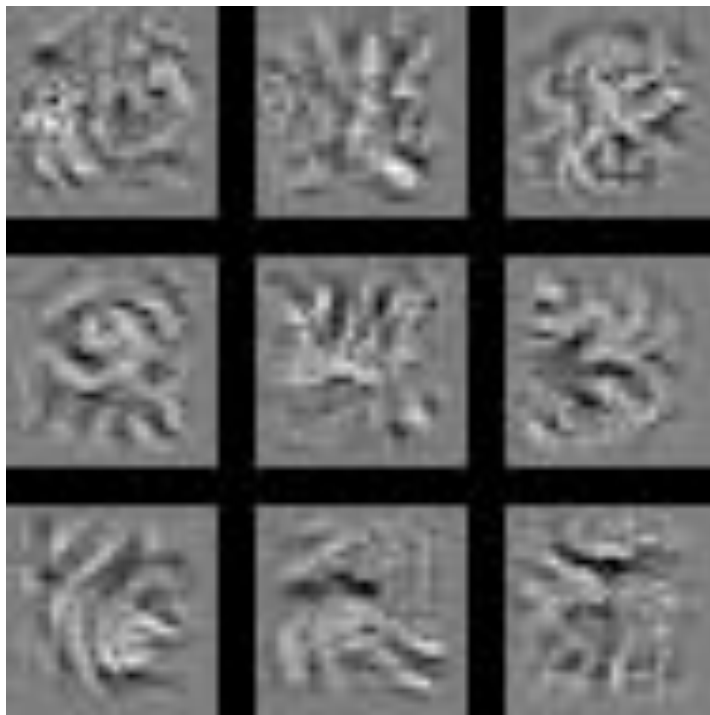
For each filter



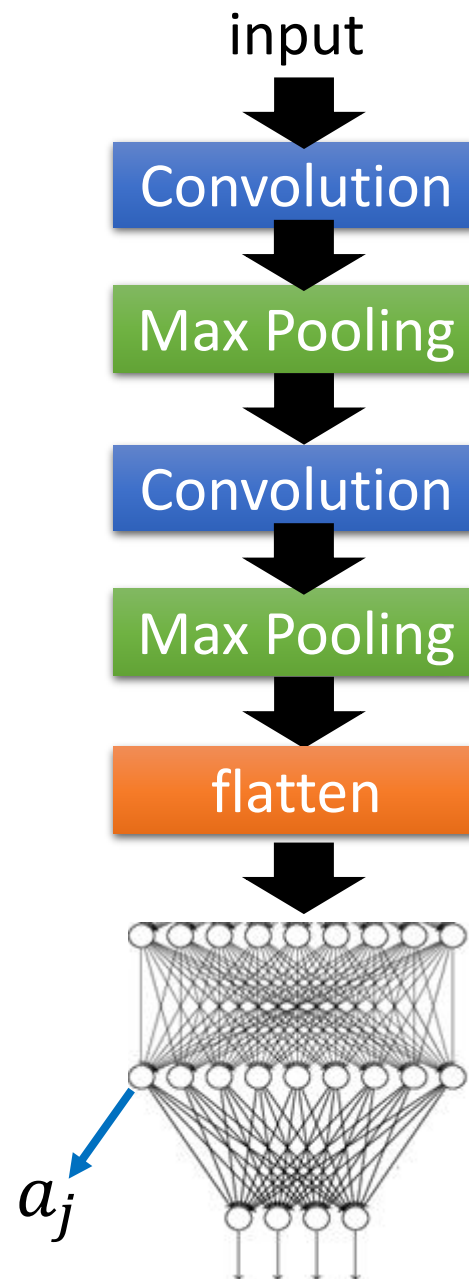
What does CNN learn?

Find an image maximizing the output of neuron:

$$x^* = \arg \max_x a^j$$

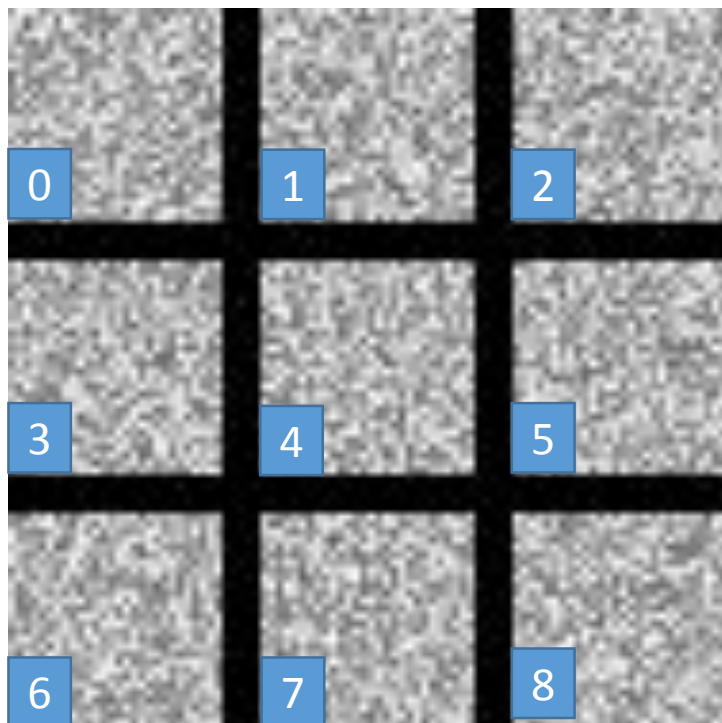


Each figure corresponds to a neuron



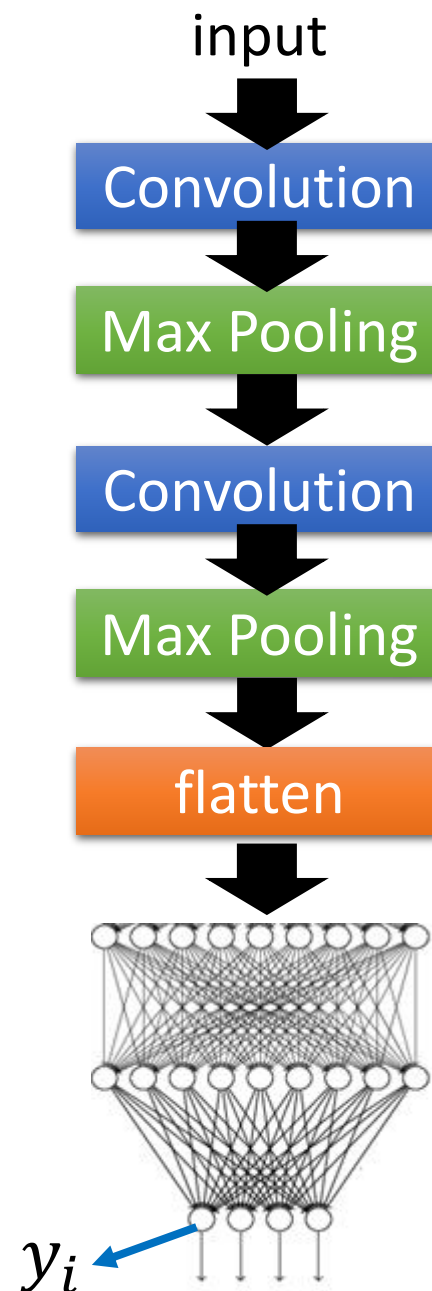
What does CNN learn?

$$x^* = \arg \max_x y^i \quad \text{Can we see digits?}$$



Deep Neural Networks are Easily Fooled

<https://www.youtube.com/watch?v=M2lebCN9Ht4>

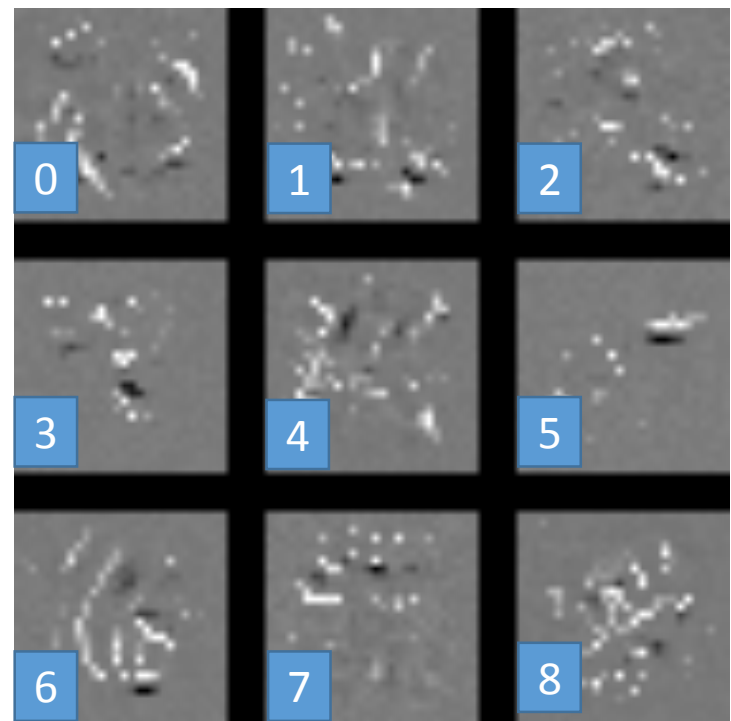
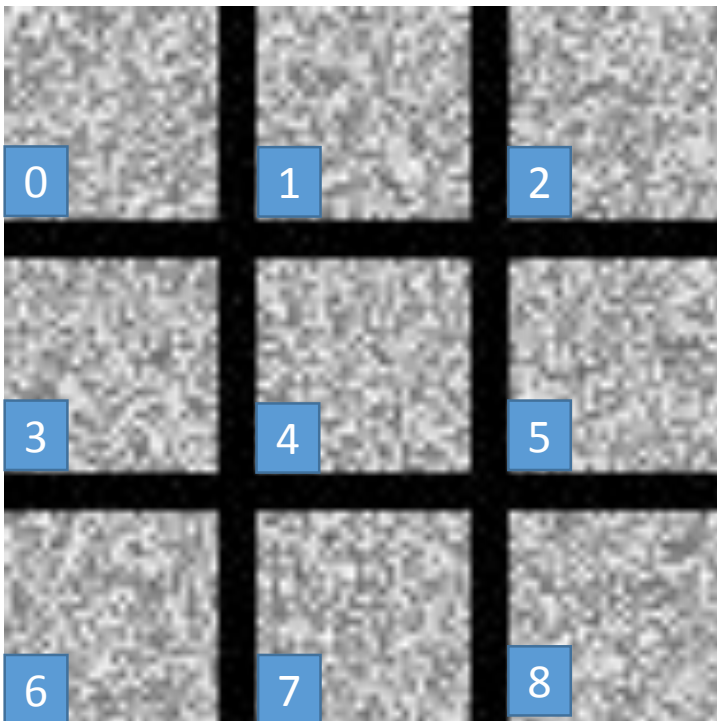


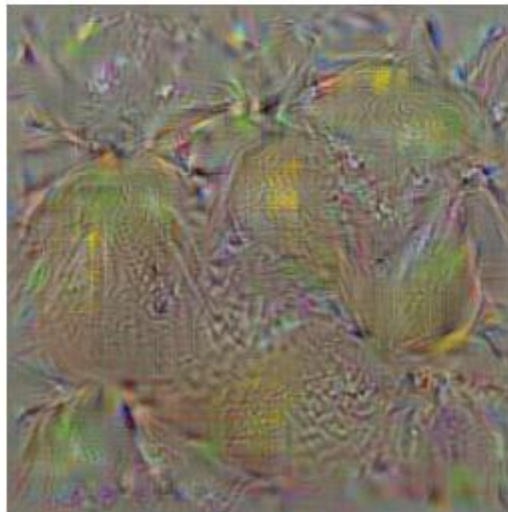
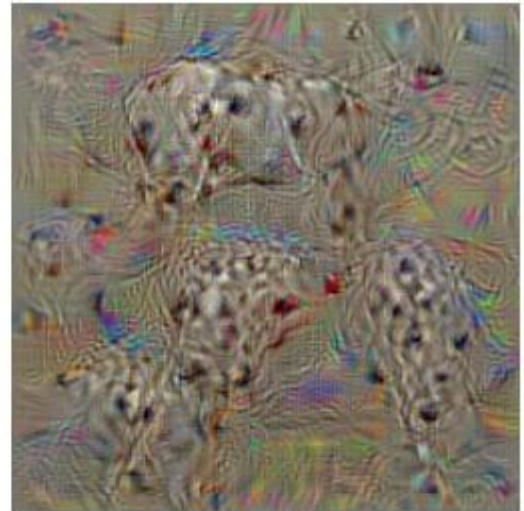
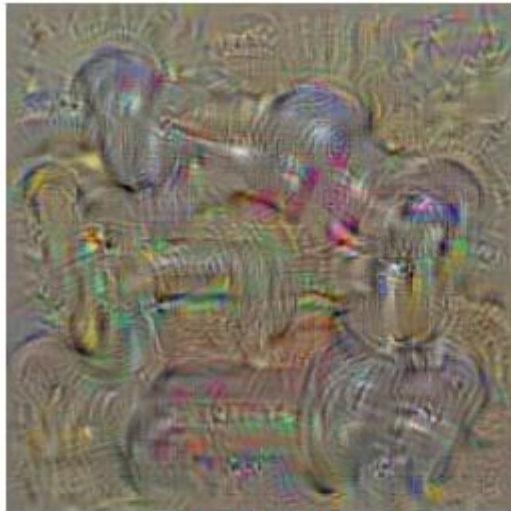
What does CNN learn?

$$x^* = \arg \max_x y^i$$

Over all
pixel values

$$x^* = \arg \max_x \left(y^i - \sum_{i,j} |x_{ij}| \right)$$





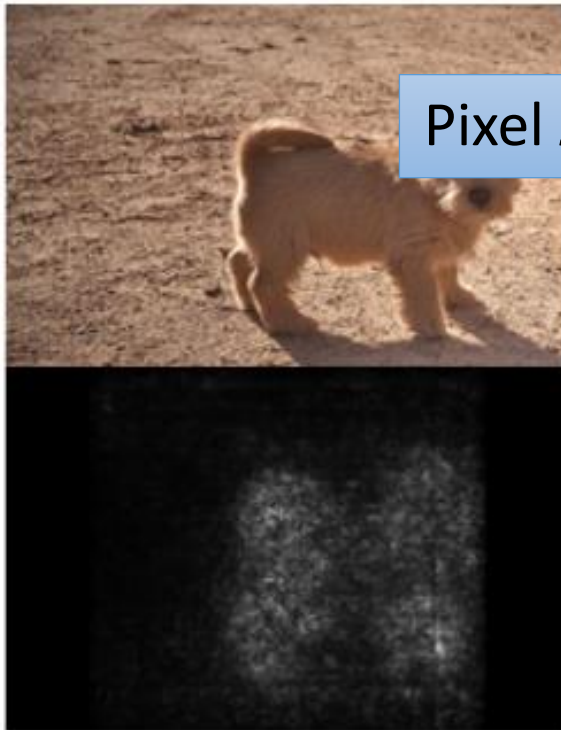
Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR, 2014

$$\left| \frac{\partial y_k}{\partial x_{ij}} \right|$$

y_k : the predicted
class of the model



Pixel x_{ij}

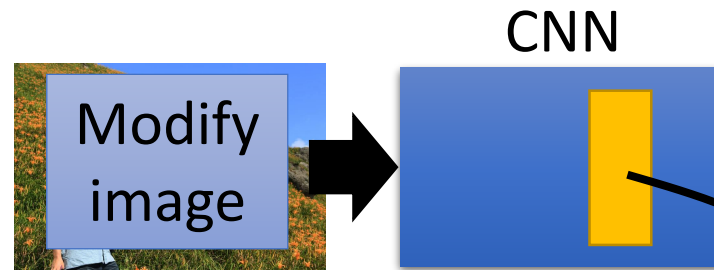


Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR, 2014



Reference: Zeiler, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional networks. In *Computer Vision—ECCV 2014* (pp. 818-833)

Deep Dream



- Given a photo, machine adds what it sees



$\begin{bmatrix} 3.9 \\ -1.5 \\ 2.3 \\ \vdots \end{bmatrix}$

Green arrow pointing up next to 3.9
Orange arrow pointing down next to -1.5
Green arrow pointing up next to 2.3

Deep Dream

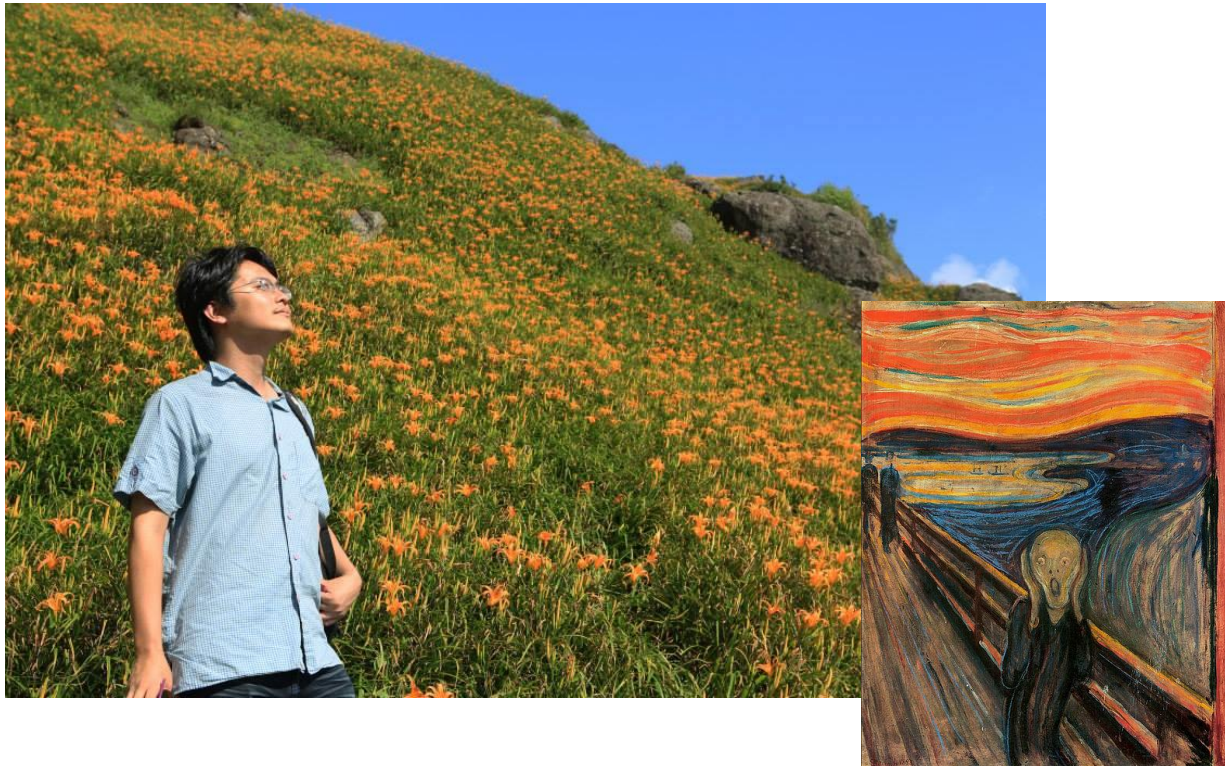
- Given a photo, machine adds what it sees



<http://deepdreamgenerator.com/>

Deep Style

- Given a photo, make its style like famous paintings



<https://dreamscopeapp.com/>

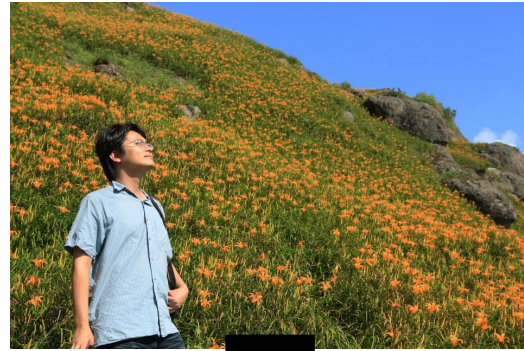
Deep Style

- Given a photo, make its style like famous paintings



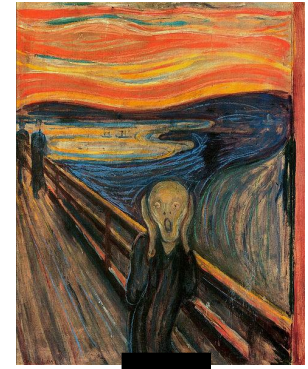
<https://dreamscopeapp.com/>

Deep Style



CNN

content



CNN

style

A Neural
Algorithm of
Artistic Style

<https://arxiv.org/abs/1508.06576>



CNN

?

More Application: Playing Go



Black: 1
white: -1
none: 0



Network



Next move
(19 x 19
positions)

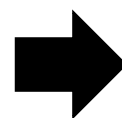
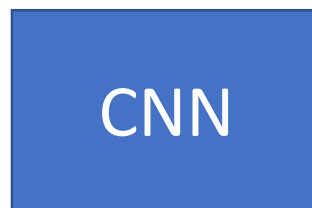
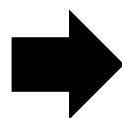
19 x 19 vector

Fully-connected feedforward
network can be used

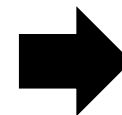
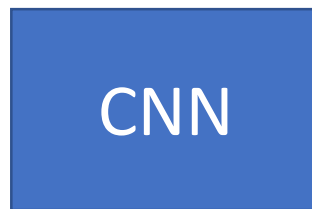
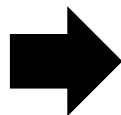
But CNN performs much better.

More Application: Playing Go

Training: record of previous plays 黒: 5之五 → 白: 天元 → 黒: 五之5 ...



Target:
“天元” = 1
else = 0

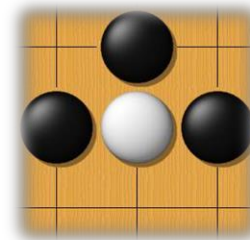


Target:
“五之5” = 1
else = 0

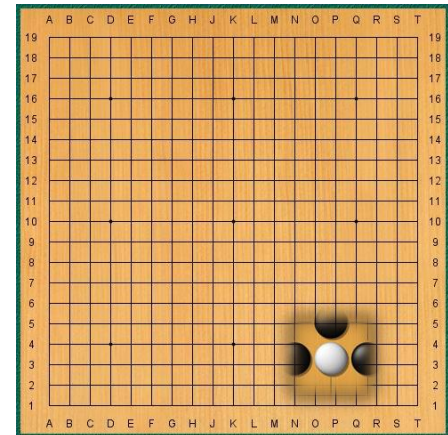
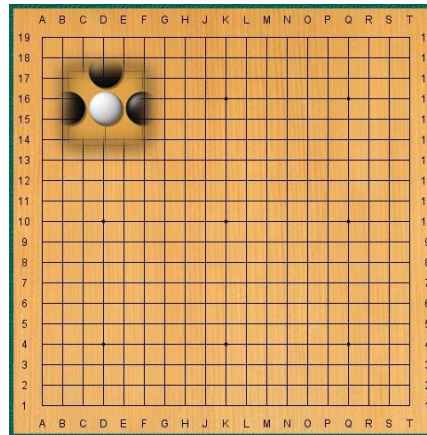
Why CNN for playing Go?

- Some patterns are much smaller than the whole image

Alpha Go uses 5 x 5 for first layer



- The same patterns appear in different regions.



Why CNN for playing Go?

- Subsampling the pixels will not change the object



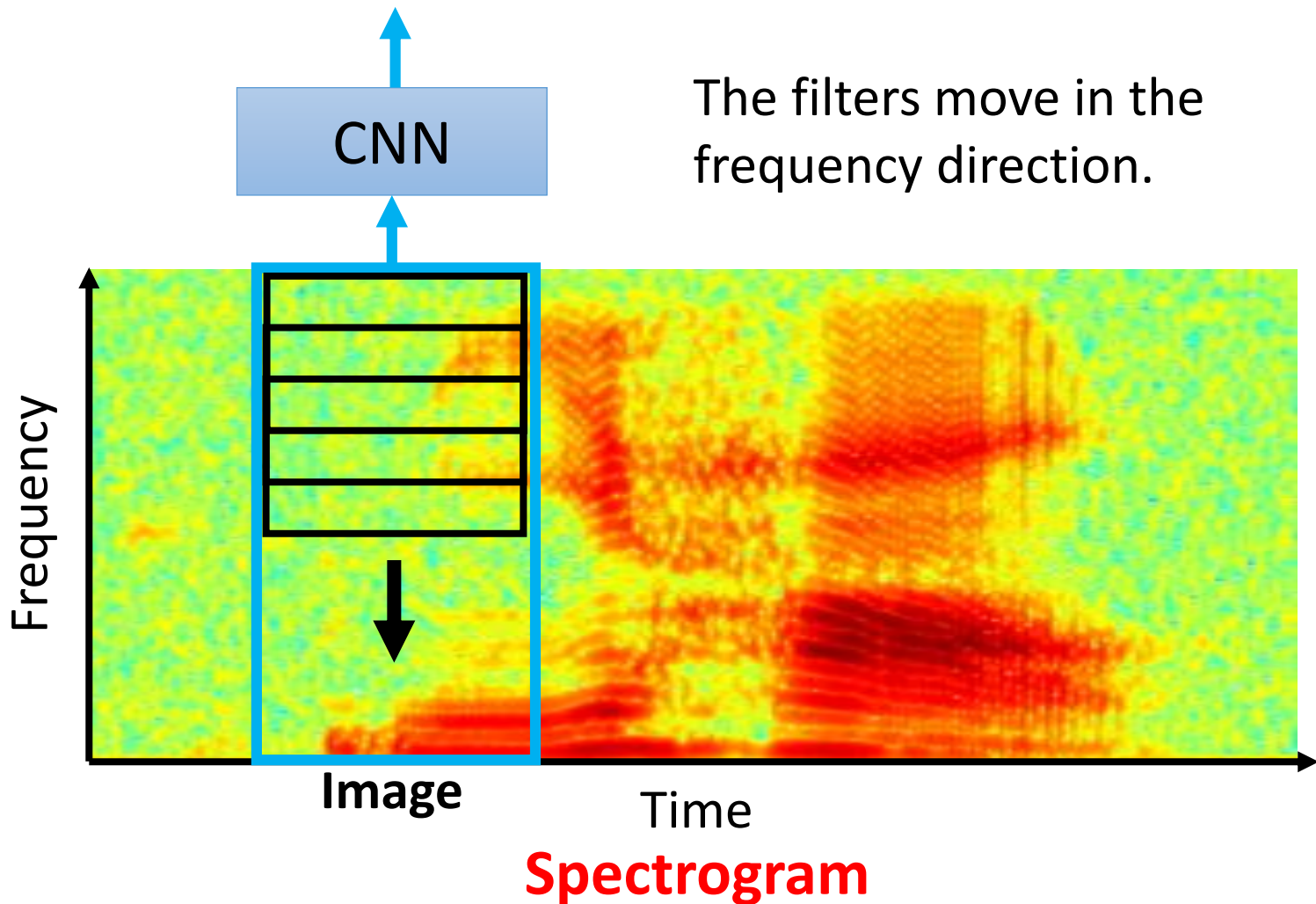
Max Pooling

How to explain this???

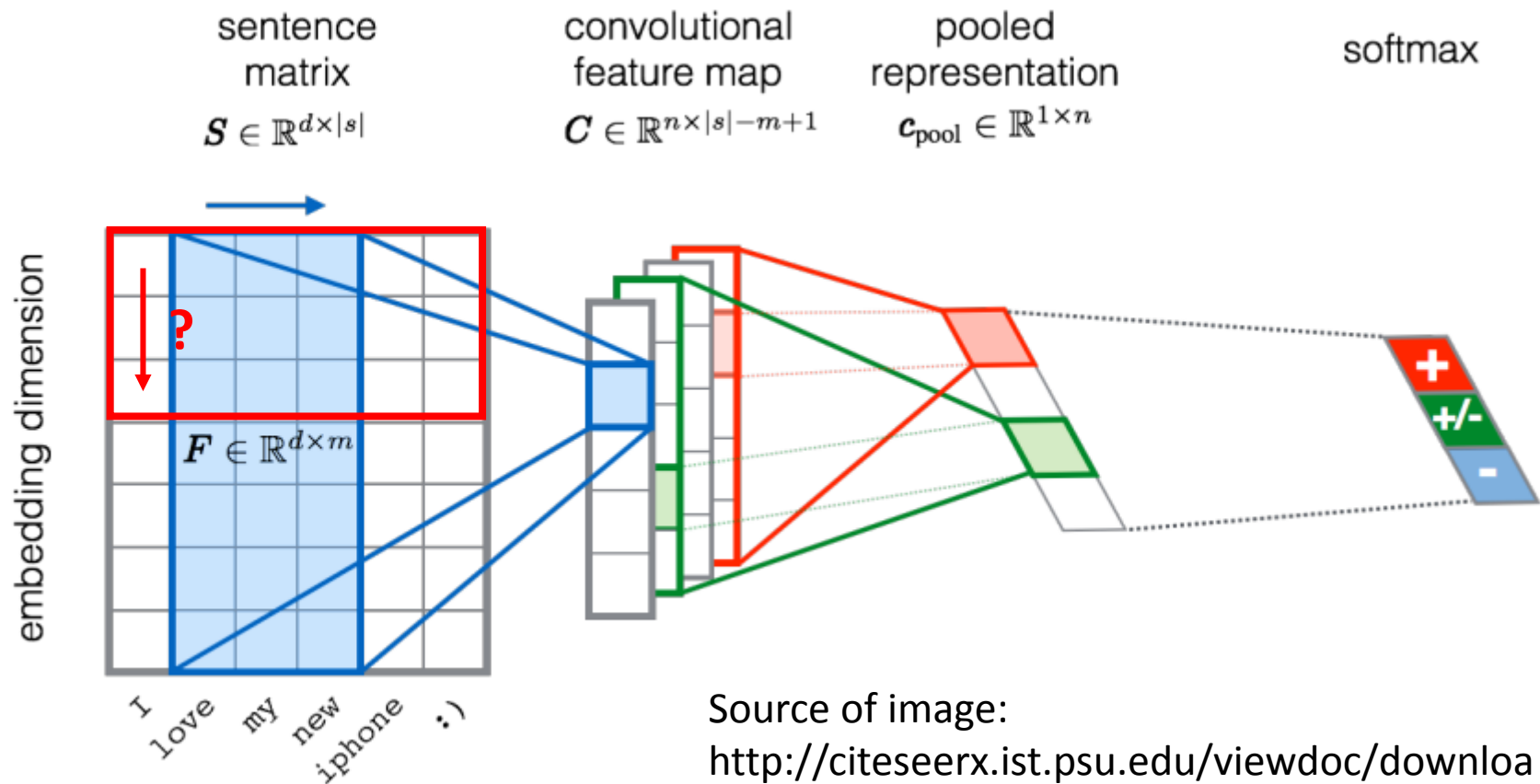
Neural network architecture. The input to the policy network is a $19 \times 19 \times 48$ image stack consisting of 48 feature planes. The first hidden layer zero pads the input into a 23×23 image, then convolves k filters of kernel size 5×5 with stride 1 with the input image and applies a rectifier nonlinearity. Each of the subsequent hidden layers 2 to 12 zero pads the respective previous hidden layer into a 21×21 image, then convolves k filters of kernel size 3×3 with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size 1×1 with stride 1, with a different bias for each position, and applies a softmax function. The

Alpha Go does not use Max Pooling Extended Data Table 3 additionally show the results of training with $k = 128, 256$ and 384 filters.

More Application: Speech



More Application: Text



Source of image:
<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.703.6858&rep=rep1&type=pdf>

Acknowledgment

- 感謝 Guobiao Mo 發現投影片上的打字錯誤