

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/343846181>

A comprehensive survey on convolutional neural network in medical image analysis

Article in *Multimedia Tools and Applications* · August 2020

DOI: 10.1007/s11042-020-09634-7

CITATIONS
20

READS
1,190

4 authors:



Xujing Yao

University of Leicester

9 PUBLICATIONS 214 CITATIONS

[SEE PROFILE](#)



Xinyue Wang

2 PUBLICATIONS 27 CITATIONS

[SEE PROFILE](#)



Shuihua Wang

University of Leicester

791 PUBLICATIONS 23,820 CITATIONS

[SEE PROFILE](#)



Yu-Dong Zhang

University of Leicester

654 PUBLICATIONS 20,953 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Face detection and recognition. Human activity recognition. [View project](#)



ICCSA 2022- Workshop "Advanced Mathematics and Computing Methods in Complex Computational Systems (AMCM2022) [View project](#)

1 **A comprehensive survey on Convolutional Neural Network in**
2 **medical image analysis**

3 Xujing Yao¹, Xinyue Wang¹, Shui-Hua Wang^{1,2,3,*}, Yu-Dong Zhang^{1,2,*},

4 1 School of Informatics, University of Leicester, Leicester, LE1 7RH, UK

5 2 Department of Information Systems, Faculty of Computing and Information Technology, King Abdulaziz University, Jeddah
6 21589, Saudi Arabia

7 3 School of Architecture Building and Civil engineering, Loughborough University, Loughborough, LE11 3TU, UK

8 Email: XY (xy147@le.ac.uk), XW (fireswxy@163.com), SHW (shuihuawng@ieee.org), YDZ (yudongzhang@ieee.org)

9 Correspondence should be addressed to SHW & YDZ

10
11 **Abstract:** CNN is inspired from Primary Visual (V1) neurons. It is a typical deep learning technique and
12 can help teach machine how to see and identify objects. In the most recent decade, deep learning develops
13 rapidly and has been well used in various fields of expertise such as computer vision and natural language
14 processing. As the representative algorithm of deep learning, Convolution Neural Network (CNN) has
15 been regarded as a breakthrough of historic significance in image processing and visual recognition tasks
16 since the astonishing results achieved on ImageNet Large Scale Visual Recognition Competition
17 (ILSVRC). Unlike methods based on handcrafted features, CNN models can build high-level features
18 from low-level ones in a data-driven fashion and have displayed great potential in medical image analysis
19 among the aspects of segmentation of histological images identification, lesion detection, tissue
20 classification, etc. This paper provides a review on CNN from the perspectives of its basic mechanism
21 introduction, structure, typical architecture and main application in medical image analysis through
22 analyzing over 100 references from Google Scholar, PubMed, Web of Science and various sources
23 published from 1958 to 2020.

24
25 **Keywords:** Deep learning; Feedforward Neural Network; Convolutional neural network; Breast
26 Cancer; Lung Nodule; Brain Tumor; Medical image analysis

27
28 **1. BACKGROUND**

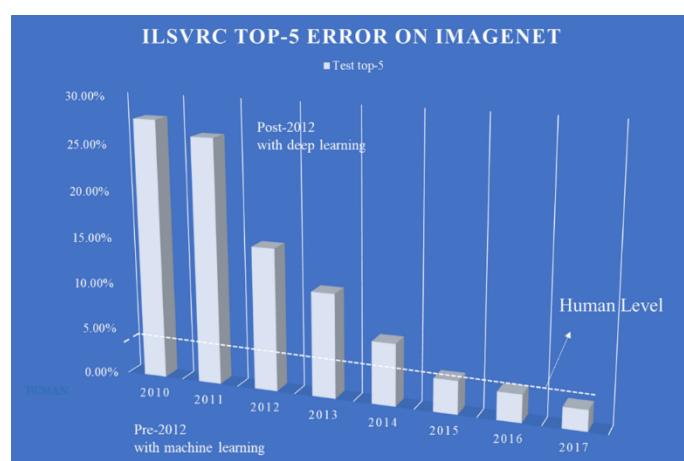
30 Deep Learning (DL) is a promising research direction in the Machine Learning (ML) area. By
31 means of utilizing DL technique, computational models are capable to learn the features and
32 representation of data with a process of extraction. This is an innovation that brings dramatical
33 improvements in image recognition, object detection, voice recognition, genomic science and many other
34 domains in real life [1, 2]

35 The study of machine learning is inspired in part by the exploration of human brain. Around 300B.C.,
36 the concept of understanding the brain has attracted a lot of attention of researchers, among which the
37 celebrated philosopher Aristotle is considered as the forerunner, due to the reason that he opens the
38 history of human trying to learn the brain [3]. In 1943, the team of McCulloch and Pitts [4] brought the
39 earlier neural network model: MCP Neural Model into the attention of public. About one decade later,
40 Rosenblatt [5] introduced the perceptron based on this MCP model to complete the classification task.
41 However, just one decade after the introduction, a significant flaw was identified by Minsky and Papert
42 [6] that the perceptron cannot handle in the nonlinear classification problems, which hindered the

43 subsequent progress of neural network for approximately twenty years [7]. In 1986, a book published by
 44 [Rumelhart, Hinton](#) [8] bring the turnaround for the development of neural network. The book detailedly
 45 introduced the back propagation algorithm as well as the whole derived process of utilizing back
 46 propagation algorithm and the activation function Sigmoid to deal with the nonlinear problems. This not
 47 only led to a great deal of uproar in academia but also attracted a sharp increase number of people into
 48 the research of deep learning. In 1997, the CNN model is presented [9-11]. In the same year, [Hochreiter](#)
 49 and [Schmidhuber](#) [12] discovered Long Short-Term Memory (LSTM) Networks model. After the
 50 introduction of ReLU, neural network has so far been known as the most outstanding tool to deal with
 51 learning problems [13, 14].

52 Because of the efforts of so many great scholars, as shown in Figure 1, deep learning is becoming
 53 more and more accurate and has evolved to perform many complex tasks and achieve loads of
 54 unimaginable breakthroughs [15-18].
 55

56
 57



(a) ILSVRC Top-5 Error on ImageNet



(b) ImageNet Champion results over the years

60 **Figure 1 Breakthroughs in ImageNet competitions over recent years**

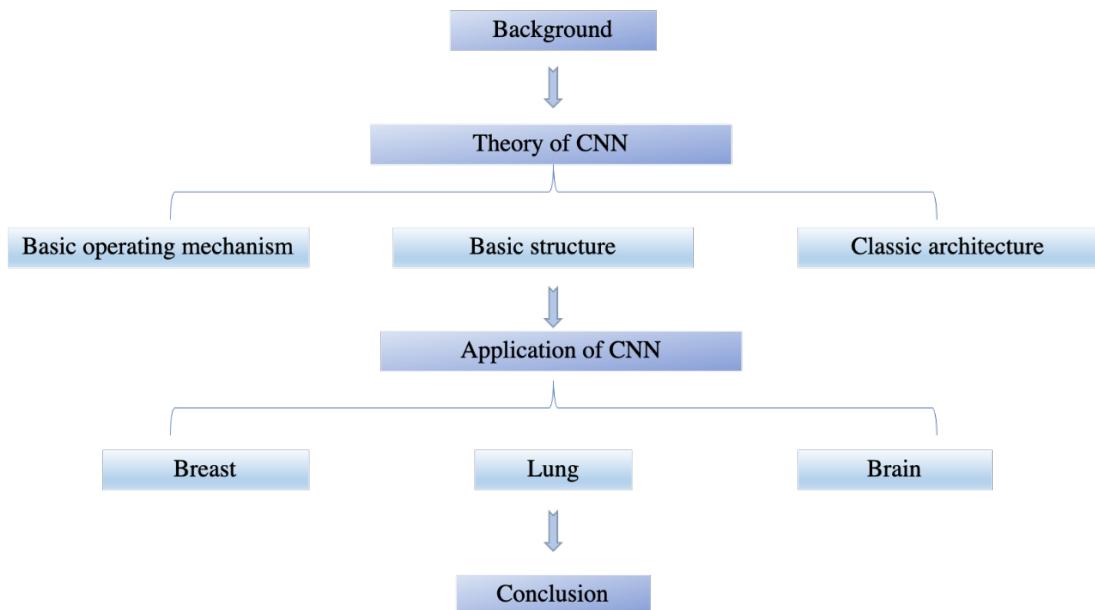
61

62 This survey will introduce the application of CNN technology in analysing medical images. The
 63 reason why we chose to introduce CNN is that it can process image information better than deep neural
 64 network (DNN). DNN is one of the simplest neural networks. In DNN, each neuron belongs to a different
 65 layer and is connected with all neurons in the previous layer. Signal is transmitted one-way from the
 66 input layer to the output layer. While CNN is a feedforward neural network through convolution

67 computation, which is proposed by the biological receptive field mechanism and has translational
68 invariance. It utilizes the convolution kernel to maximize the application of local information and retain
69 the planar structure information. DNN is input in vector form, without considering the structure
70 information of the plane. In the field of image processing, the plane information is very significant, so
71 CNN performs better than DNN in medical image analysis [19].

72 A simple flowchart of this survey is organized as Figure 2. The purpose and contribution of this
73 review is to: (i) introduce the theoretical basis and working principle of CNN in detail; (ii) elaborate
74 description of the popular forms of various CNN basic structures; (iii) detailedly explain classic models
75 and trends from the rise of CNN to the present. (iv) demonstrate the outstanding contributions of CNN
76 models in the medical field from recent years. It can be found from our review that CNN model is
77 increasingly being used in medical fields where it is difficult to find the basis for diagnosis, or hard to
78 achieve high accuracy or fast diagnostic analysis by just utilizing human power. Perhaps in the next few
79 years, humans will be able to use CNN to solve medical problems that are currently difficult to understand
80 and analyze from a medical perspective.

81



82

83

84

Figure 2 Simple flowchart of this survey

85 Section 2, 3 and 4 will explain CNN from three aspects: the basic operating mechanism, the basic
86 structure and typical architectures of CNN. The operating mechanism section introduces the operation
87 unit, operation thought, gradient vanishing, feedforward propagation and feedback propagation. The
88 basic structure section describes principles and details of convolutional layer, pooling layer and fully
89 connected layer. The classic architecture section reviews the theory and advantages of several typical
90 architectures such as AlexNet, VGG-Net, Inception, ResNet, R-CNN, SPP-Net and so on.

91

92

2. BASIC OPERATING MECHANISM

93

2.1 Operation Unit: Neuron

94

95

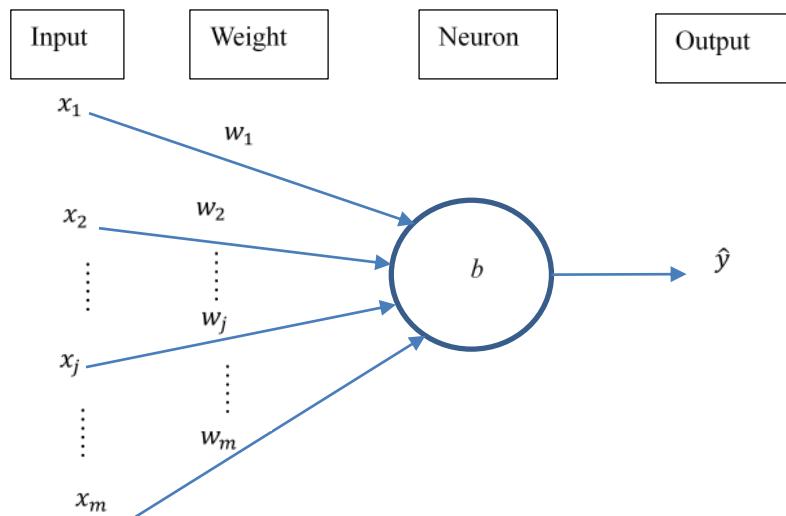
96 In 1943, the MCP Model introduced by [McCulloch and Pitts \[4\]](#) is very similar with biological

97 neural network. From a biological kind of view, every neuron is linked with other neurons. When the
 98 neuron is activated, which means the potential of neuron is over a threshold, it will send chemical
 99 substances to the connected neurons to transform the potential of other neurons. MCP model has the
 100 same mechanism. The circle in Figure 3 represents neuron, which receives the input signals ($x_1 \dots x_m$)
 101 from the other neurons. These input signals will get different weights ($w_1 \dots w_m$) and bias (b) and then
 102 go to the neuron. the output \hat{y} will be got through the activation function [14, 20].

$$\hat{y} = \sigma(z) = \sigma(\sum_{j=1}^m w_j x_j + b) \quad (1)$$

104 Where we assume the input of the activation function as z

105



106

107

108

Figure 3 The McCulloch and Pitts model of a neuron

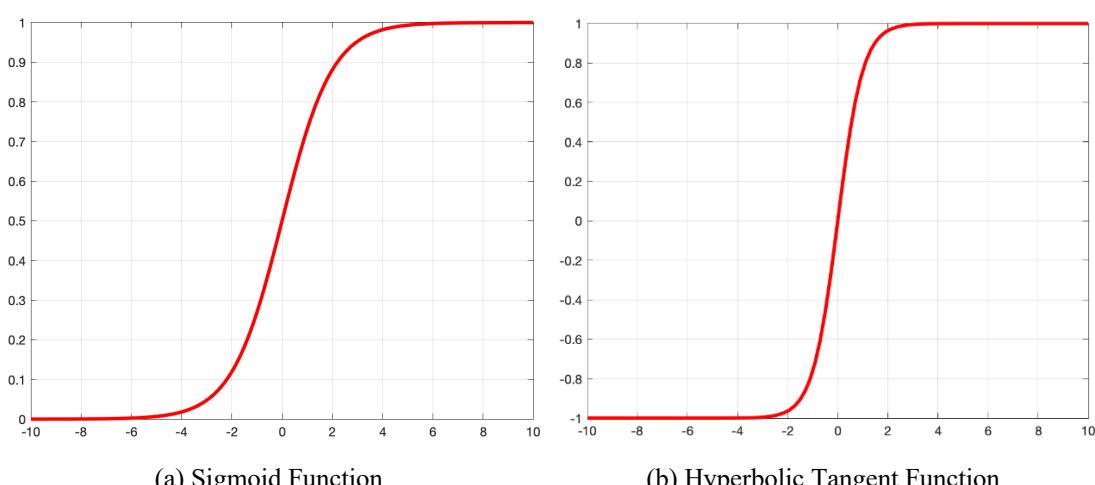
109

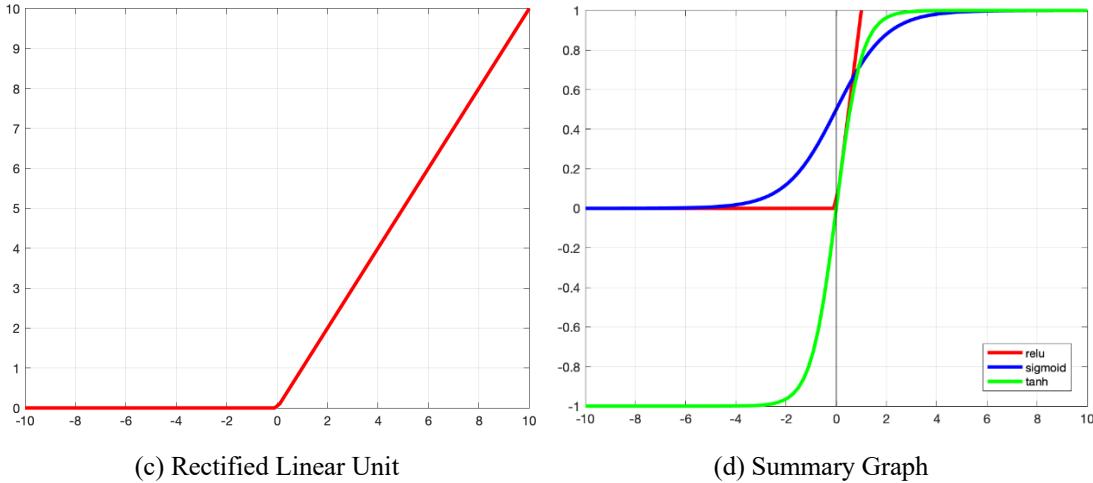
110 In a multi-layer neural network, each neuron node will accept the output of previous layer as its
 111 input value and then pass the value to the following layer. This function relationship between neuron
 112 nodes is defined as the activation function. If you just utilize the linear activation function or do not use
 113 the activation function, all things are about calculating the linear activation function, which is very
 limited, so in order to get more possibilities, the nonlinear activation function is born [21-23].

114

115 2.2 Activation Function and Loss Function

116





(c) Rectified Linear Unit

(d) Summary Graph

Figure 4 Activation functions

117

118

119 **2.2.1. Sigmoid Function**

120

121 Sigmoid Function aims to alter the continuous real value of the input z to the output $\sigma(z)$ of a range
122 between 0 and 1.

123

$$\sigma_{\text{Sigmoid}}(z) = \frac{1}{1+e^{-z}} \quad (2)$$

124

125 It's clear from Equation (2) that if the input is an extremely large negative number, the output is
126 going to be 0. If the input value is an extremely large positive number, the output will be 1. But this
127 function also has obvious disadvantages that it consumes loads of calculation time, and will be easy to
128 produce gradient disappearance and gradient explosion [24]. In addition, as displayed in Figure 4(a), its
129 data is not centralized and symmetric with (0,0.5) as the center, which will bring a negative effect on the
130 gradient calculation of the hidden layer.

131

132 **2.2.2. Hyperbolic Tangent Function (Tanh)**

133

134 Tanh function performs superior to the sigmoid function in hidden layer. It can be viewed as the
135 shifted version of sigmoid as shown in Figure 4(b). The equation of Tanh is shown as:

136

$$\sigma_{\text{Tanh}}(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad (3)$$

137

138 Although it achieves a similar effect of data centralization, in practice, sigmoid and tanh functions
139 tend to reach saturation at the end value, especially in the deep network training, which leads to a slow
140 training speed.

141

142 **2.2.3. Rectified Linear Unit Function (ReLU)**

143

144 It can be seen from Figure 4(c), ReLU solves the gradient disappearance problem when there is a
145 positive input [25]. The equation is shown below

146

$$\sigma_{\text{ReLU}}(z) = \max(0, z) \quad (4)$$

147 Moreover, ReLU costs much less computation time when compared with sigmoid and tanh function.
148 But the gradient vanishing problem remains and the data distribution of ReLU is not centralized when

147 there is a negative input [26].

148

149 2.2.4. Loss Function

150

151 After the activation function, the calculated value \hat{y}_j is got. However, the calculated value cannot
152 equal to the actual value y_j due to some unexplained factors. The loss function aims to minimize the
153 distance between each actual value and the predicted value in the training set, which refers to, the sum
154 of the error $J(w)$ between \hat{y}_j and y_j through utilizing the mean square error principle shown below:

155
$$J(w) = \frac{1}{2} \sum_{j=1}^m (\hat{y}_j - y_j)^2 \quad (5)$$

156 Where w represents weight, which is the coefficient of linear relation

157 However, the mean square error has many disadvantages when it is used in the back propagation
158 algorithm which will be introduced in Section 2.41.5.1.2.4.2.

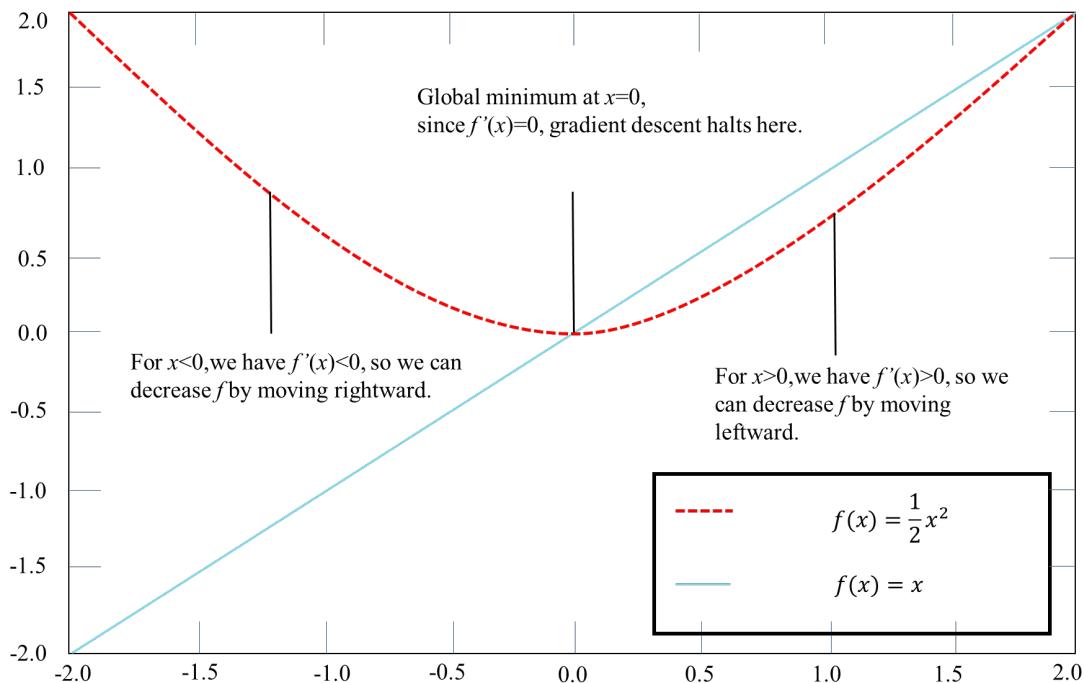
159

160 2.3 Operation Thought: gradient descent

161

162 Most deep learning algorithms involve many kinds of optimizations. Gradient descent method is
163 one of the optimizations which used to minimize the loss function. In this method, the minimum loss
164 function and the value of model parameters are obtained by step-by-step iteration [21].

165



166

167 **Figure 5 The relationship between a function and the derivatives of a function**

168

169 In 1847, Cauchy, Augustin discovered gradient descent. Assume there is a function:

170

$$y = f(x) \quad (6)$$

171 Where both x and y are real numbers

172

173 The derivative of the function is $[dy/dx]$ or $f'(x)$, which presents the slope of $f(x)$ at point x .
Besides, it also shows if the input has a little change like ϵ , the output can also be represented by the

174 slope and the original output, the formula is shown below [27]:

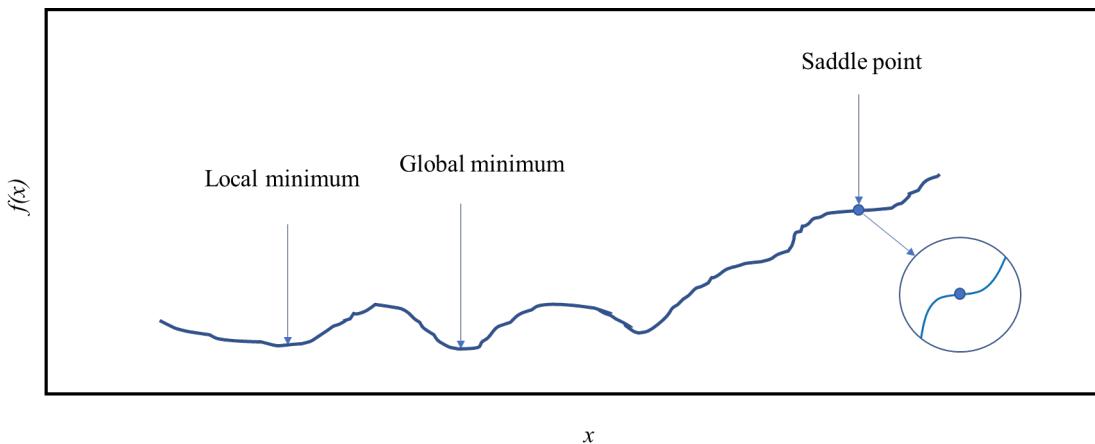
$$f(x) \leftarrow f(x + \epsilon) = f(x) + \epsilon f'(x) \quad (7)$$

175 Where \leftarrow represents the assignment operation.

177 With this equation, it is easy to control the output decrease. When the slope is negative the input
178 should become bigger, and when the slope is positive the input should become smaller in order to get a
179 smaller output value as the Figure 5 shows.

180 Since the loss function should be smallest, the equation will run until it meets the three possibilities
181 as shown in Figure 6. First, the output is the minimum value within limits, which is called local minimum.
182 Second, the output is the lowest point in the whole function, which is called global minimum. And the
183 third possibility is that the output is a saddle point. A saddle point refers to a critical point which is neither
184 maximum nor minimum.

185



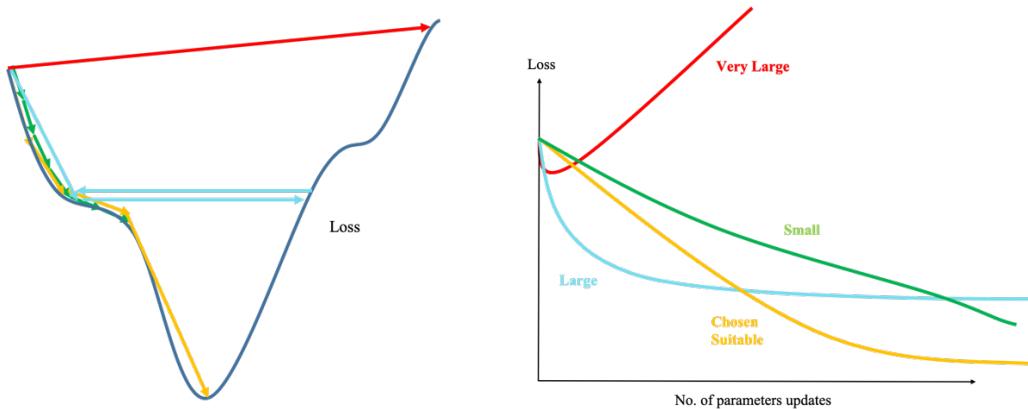
186

187 **Figure 6 Three possibilities the formula will stop when it meets**

188

189 In order to avoid the output which comes from the saddle point, it is important to choose the suitable
190 learning rate ϵ . As displayed in Figure 7, if the learning rate is chosen suitable, the minimum value can
191 be found easily as the orange line. But if ϵ is too small as the green line, the time to run the equation
192 will increase and the output may be generated from a saddle point or local minimum. If ϵ is large as the
193 blue line, it is impossible to find the minimum. If ϵ is too large as the red line, the equation will loss
194 control and the output is far from the right point. Therefore, it is important to control learning rate ϵ . It
195 is convenient and visualized to draw a graph about loss and parameters change. And from Figure 7, it is
196 obvious that the orange line gives us the best learning rate.

197



(a) (b)

Figure 7 The relationship between loss and parameters change

200 Except choosing the suitable learning rate, different kinds of gradient descent can improve the loss
201 function in different levels. We assume the original equation is

$$202 \quad \hat{y} = b + w_1 x_1 + w_2 x_2 + \cdots w_m x_m \quad (8)$$

203 The following introductions and equations come from the loss function based on mean squared
204 error:

$$205 \quad J(w) = \frac{1}{2} \sum_{l=1}^n \sum_{j=1}^m (\hat{y}(x_j^l) - y^l)^2 \quad (9)$$

206 Where:

207 x_j^l the j -th neuron in the l -th layer

208 n the total number of layers

209 Three basic gradient descent variants are introduced next.

211 2.3.1. Batch Gradient Descent (BGD)

212 As shown in Equation (10), BGD requires the calculation of partial derivatives for the whole data
213 set according to Equation (7) and (9). To get the optimum w for achieving the maximum value of $J(w)$,
214 we will repeat the following equation until convergence:

$$216 \quad w_i \leftarrow w_i + \epsilon \sum_{l=1}^n \sum_{j=1}^m (y^l - \hat{y}(x_j^l)) x_{i,j}^l \quad (10)$$

217 Where \leftarrow represents the assignment operation

218 ϵ represents the learning rate

219 i represents the number of iterations

220 However, it also has obvious disadvantages. In one update, it calculates the gradient for the whole
221 data set, so the calculation time is very long. In this way, it will be very difficult to deal with a large scale
222 data set, and update the model in real time with new data [28].

224 2.3.2. Stochastic Gradient Descent (SGD)

226 For a large data set, there may be similar samples, so the BGD method of calculating the gradient
227 with all the data at one time will be redundant. While SGD will update each sample with a gradient every
228 time, which is not redundant and relatively fast, and new samples can be added in. Its equation is shown
229 in Equation (11). However, for accuracy, SGD only uses one sample to determine the direction of the
230 gradient, resulting in a solution that is probably not optimal. However, due to the reason that SGD is
231 updated frequently, it may lead to serious oscillation of cost function [29, 30].

$$232 \quad w_i \leftarrow w_i + \epsilon (y^l - \hat{y}(x_j^l)) x_{i,j}^l \quad \forall j = 1, \dots, m \wedge \forall l = 1, \dots, n \quad (11)$$

234 2.3.3. Mini-Batch Gradient Descent (MBGD)

236 The difference between MBGD and SGD is that MBGD used a small batch of sample data consists
237 of n samples to calculate the loss function during each update. The value range of the hyperparameter n

238 is generally between 50 and 256. In this method, as shown in Equation (12), when updating the
 239 parameters, the variance will be reduced, the convergence will be more stable, and the calculation speed
 240 will be faster. However, if ϵ is extremely small, the convergence rate will be slow. Besides, if the loss
 241 is extremely large, the loss function will oscillate near the minimum value. In addition, the choice of
 242 batch size is also difficult [31, 32].

243 $w_i \leftarrow w_i - \epsilon \sum_{l=t}^{t+n-1} \sum_{j=1}^m (\hat{y}(x_j^l) - y^l) x_{i,j}^l$ (12)

244

245 2.3.4. Gradient vanishing and Gradient explosion

246

247 Gradient disappearance and gradient explosion are essentially the same, both of which are due to
 248 the continuous multiplication effect in gradient back propagation brought by the too deep network layers.
 249 In other words, if the numbers in each layer are smaller than 1, then the positive gradient that you get
 250 when you multiply them is very small, and that's the gradient disappearing. On the contrary, if the
 251 numbers in each layer are bigger than 1, then the forward gradient will be extremely large, that is, the
 252 gradient explosion. There are many ways to solve these two problems, such as switching to activation
 253 functions such as ReLU, LeakyReLU, ELU, utilizing batch normalization, ResNet residual structure,
 254 LSTM structure, pretraining plus finetuning, etc. [33, 34].

255

256 2.4 Propagation Thought

257

258 2.4.1. Feedforward propagation

259

260 Like the introduced in Section 2.1, The core idea of MCP is a linear relationship plus an activation
 261 function in a small local model, which is:

262 $y = \sigma(z) = \sigma_{\text{ReLU}}(b + \sum w * x).$ (13)

263 Where $*$ represents the convolution, b represents bias, σ represents activation function, usually refers to
 264 ReLU.

265 And since MCP is very simple which can only solve binary classification, the Multi-Layer
 266 perceptron (MLP) is introduced which is also called deep neural network (DNN). The major difference
 267 between DNN and MCP is the number of layers. With the core idea of MCP as the feedforward
 268 propagation of DNN, the general formula is:

269 Assume l -th layer has m neurons:

270 $y_j^l = \sigma_{\text{ReLU}}(z_j^l) = \sigma_{\text{ReLU}}(\sum_{k=1}^m w_{jk}^l y_k^{l-1} + b_j^l)$ (14)

271 where w_{jk}^l refers to weight between the j th neurons in the l th layer and the k th neurons in the layer before

272 It can be simplified that the output of l -th layer is:

273 $y^l = \sigma_{\text{ReLU}}(z^l) = \sigma_{\text{ReLU}}(W^l y^{l-1} + b^l)$ (15)

274 Where W^l refers to a $n \times m$ matrix formed by linear coefficient w of layer l

275 y^{l-1} refers to a $m \times 1$ vector formed by output y of layer $l-1$

276 b^l refers to a $n \times 1$ vector formed by bias b of layer l

277 z^l refers to a $n \times 1$ vector formed by linear output before activating of layer l

278 However, DNN has a critical defect that if the input is large, a huge amount of weights and other
 279 parameters require being trained and is likely to lead to overfitting and local optimization. Therefore, the
 280 convolutional neural network (CNN) is discovered by [LeCun, Boser \[35\]](#) in 1989. It has two other kinds
 281 of layer, one is convolutional layer, the other is pooling layer. The detail information will introduce in

282 Section 3.

283

284 2.4.2. Feedback propagation

285

According to the Section 2.1, loss function aims to minimize the difference between the calculated output (\hat{y}_j) and the real output (y_j). And the difference between them is caused by the random weights and bias. In order to get suitable weights and bias to make the calculated output correctly, the feedback propagation is used. The main concept of feedback propagation is to optimize the loss function by gradient descent. We assume W and X represent the weights matrix and input vector in one layer separately.

292 Since the loss function is

$$J(W) = \frac{1}{2} \sum_{l=1}^n (\hat{y}^l - y)^2 \quad (16)$$

294 Where:

295 \hat{y}^l refers to calculated output in the l -th layer

296 y refers to real output in the l -th layer

297 *n* refers to total number of layers

$$\hat{y}^l = \sigma_{\text{ReLU}}(z^l) = \sigma_{\text{ReLU}}(b^l + \sum W^l * X^{l-1}) \quad (17)$$

299 The loss function's partial derivative respect to W^l is

$$300 \quad \frac{\partial J(\mathcal{W})}{\partial W^l} = [(\hat{y}^l - y) \odot \sigma'_{ReLU}(b^l + \sum W^l * X^{l-1})](\hat{y}^{l-1})^T \quad (18)$$

and its partial derivative respect to b^l is

$$\frac{\partial J(W)}{\partial b^l} = \delta^l [(\hat{y}^l - y) \odot \sigma'_{ReLU}(b^l + \sum W^l * X^{l-1})] \quad (19)$$

303 Where \odot represents Hadamard Product, which is defined as the product of the corresponding elements
304 of two matrices

With the two Partial derivative formula, we utilize gradient descent to get the new weight W^l and bias of l -th layer b^l [36]. The formulas are as following:

$$307 \quad W^l \leftarrow W^l - \varepsilon \frac{\partial J(W)}{\partial w^l} \quad (20)$$

$$308 \quad b^l \leftarrow b^l - \varepsilon \frac{\partial J(W)}{\partial b^l} \quad (21)$$

309 Where ε is learning rate

310 Assume the layer number of the last layer which is output layer is L . Then, the equation of gradient
311 expression which is recorded as δ of the L layer is shown below:

$$\delta^L = \frac{\partial J(W, b)}{\partial z^L} = (\hat{y}^L - y) \odot \sigma'_{ReLU}(b^L + \sum W^L * X^{L-1}) \quad (22)$$

With the above formulas, the gradient expression of layer can be calculated by the fully connected layer (FC layer) after it. Assuming the parameters and gradient expression of l -th fully connected layer are known, the gradient expression of $(l - 1)$ -th layer can be got by this formula:

$$316 \quad \delta^{l-1} = \frac{\partial J(W, b)}{\partial z^{l-1}} = \left(\frac{\partial z^l}{\partial z^{l-1}} \right)^T \frac{\partial J(W, b)}{\partial z^l} = \left(\frac{\partial z^l}{\partial z^{l-1}} \right)^T \delta^l \quad (23)$$

317 The relationship between z^l and z^{l-1} is

$$z^l = \hat{y}^{l-1} * W^l + b^l = \sigma_{ReLU}(z^{l-1}) * W^l + b^l \quad (24)$$

The above formula about the FC layer can be shown as following:

$$\delta^l = \left(\frac{\partial z^{l+1}}{\partial z^l} \right)^T \frac{\partial J(W, b)}{\partial z^{l+1}} = (W^{l+1})^T \delta^{l+1} \odot \sigma'_{ReLU}(z^l) \quad (25)$$

When the l -th layer is convolutional layer, the gradient expression of $l - 1$ layer which is hidden layer can be calculated by the formula shown below:

$$\delta^{l-1} = \left(\frac{\partial z^l}{\partial z^{l-1}}\right)^T \delta^l = \delta^l * \text{rot180}(W^l) \odot \sigma'_{ReLU}(z^{l-1}) \quad (26)$$

The convolution kernel has been rotated 180 degrees. So rot180() means flip up and down once, then flip left and right once. Since CNN also includes the pooling layer which compresses data, the data should be restored.

When the l -th layer is pooling layer, the calculation of gradient expression of $(l - 1)$ -th layer is according to equation as:

$$\delta^{l-1} = \text{upsample}(\delta^l) \odot \sigma'_{ReLU}(z^{l-1}) \quad (27)$$

This up-sample function helps accomplishing the logic of pooling error matrix amplification and error redistribution [37].

In conclusion, with the basic operating mechanism about CNN, it is relatively easy to understand the structure of the CNN introduced in the following section.

3. BASIC STRUCTURE

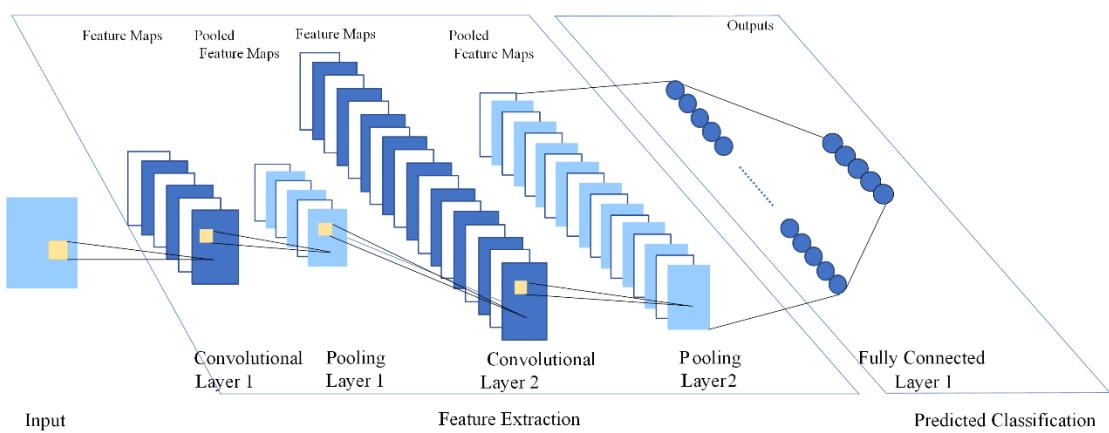


Figure 8 A sample structure of CNN

CNN can optimize the performance and generalization capability of the network by using multilayer structure [38, 39] illustrated in Figure 8. The most commonly used layers are introduced in the next.

3.1 Convolutional layer

The main ability of the convolutional layer is to extract critical features for getting a good prediction from the input images, which means reducing images into forms that are easier to process [40].

3.1.1. Standard Convolution

The element utilized in conducting the convolution operation in a convolutional layer is named the **kernel**.

351 kernels whose parameters needs to be studied. The size of the kernel should be smaller than those of the
 352 input image. In the convolution process, each kernel will be convolved with the input image to compute
 353 an activation map. That is to say, the kernel will slide across the width and height of the input, calculating
 354 the dot product between the input and the kernel at each spatial position. Next, by superimposing the
 355 activation maps of each kernel along the depth dimension, the output image of the convolution layer can
 356 be achieved. It is because each kernel is smaller in width and height when compared to the input, each
 357 neuron in the activation map is associated only to a small local area of the input image, which means size
 358 of each neuron's receptive field is small, equals to size of kernel [41].
 359

360 Figure 9 is considering a 5x5 input image in which matrix is composed of 0.5 and 0 and a 3x3 kernel
 361 composed of 1 and 0. The size of kernel determines the receptive field of the convolution. The most
 362 preference of selecting kernel size for a 2D input image is 3, which means 3x3 pixels.
 363

0.5	0	0.5	0.5	0.5
0	0.5	0	0.5	0
0	0.5	0	0	0.5
0.5	0.5	0.5	0.5	0
0.5	0	0.5	0.5	0

(a) Input Image

1	0	1
0	1	0
1	0	1

(b) 3x3 Kernel

Figure 9 Basic component of Standard Convolutional Layer.

363 Here the kernel will convolve over each cell position of the input image and eventually we will get
 364 a convoluted image. For every single point in it, you can take the point and the 3x3 points around it out
 365 and do the convolution. For instance, for the lower right (row 4, column 4) point, the local convolution
 366 of the 5x5 image and the 3x3 matrix could be computed as shown in Figure 10.
 367
 368
 369

Input Image

0.5	0	0.5	0.5	0.5
0	0.5	0	0.5	0
0	0.5	0 × 1	0 × 0	0.5 × 1
0.5	0.5	0.5 × 0	0.5 × 1	0 × 0
0.5	0	0.5 × 1	0.5 × 0	0 × 1

Convolved Feature

0	0	0.5
0	0.5	0
0.5	0	0

Figure 10 3×3 Convolution process

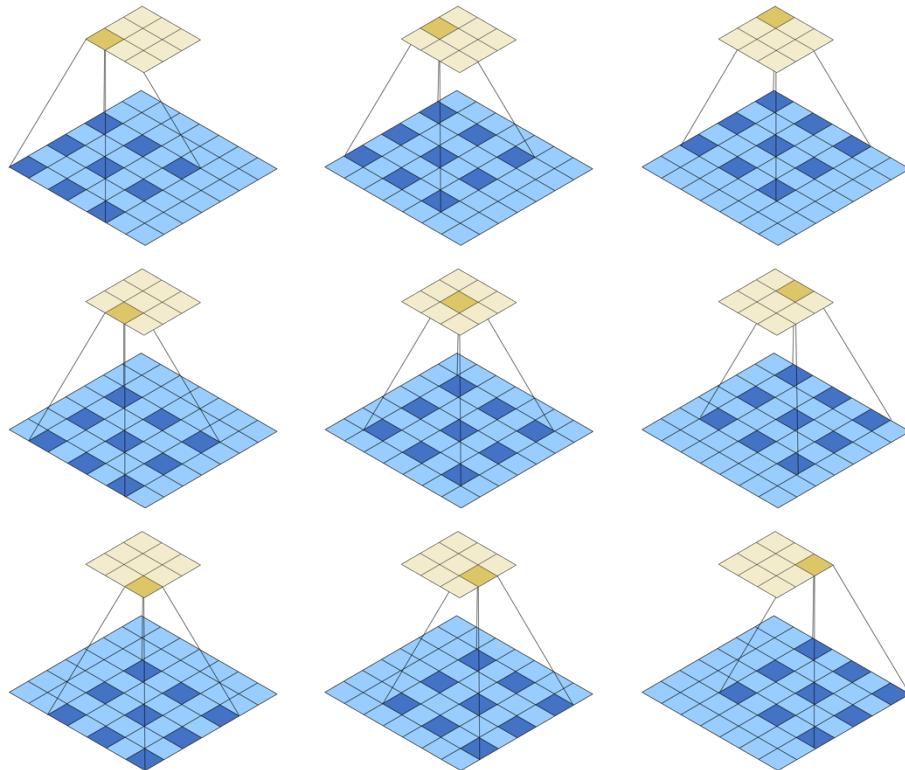
370
 371
 372

373 **3.1.2. Dilated Convolution**

374

375 One disadvantage of standard convolution is its design for up-sampling and pooling. During these
376 processes, there will be internal data structure loss and spatial hierarchical information loss. To deal with
377 this loss problem, one significant hyper-parameter called dilation rate is introduced to convolutional
378 layers. The dilation rate of kernel refers to the spacing between values inside (for a standard convolution,
379 the dilation rate equals to 1) [42]. Through the utilizing of dilated convolution, the receptive field of a
380 3×3 kernel with a dilation rate of 2 and no padding would as large as a 5×5 kernel as illustrated in Figure
381 11.

382



383

384 **Figure 11 The process of convolution utilizing a 3×3 kernel with a dilation rate of 2**

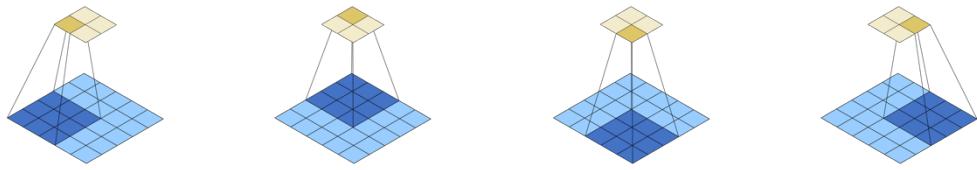
385

386 In this way, a wider reception field would be provided at the same cost of computation which makes
387 dilated convolution a promising solution among the area of real-time segmentation. It can be of great use
388 especially when there is a need of wider field of view but larger kernel size or multiple convolutions
389 cannot be afforded [43].

390

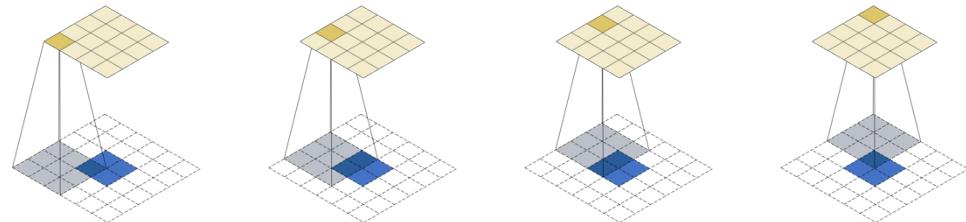
391 **3.1.3. Transposed Convolution**

392



393

394 **Figure 12 Part of the process of convolution utilizing a 3x3 kernel over a 5x5 input using unit-stride**
 395



396

397 **Figure 13 Part of the transpose process of convolution utilizing a 3x3 kernel over a 4x4 input**
 398 **utilizing unit-stride (same to convolution utilizing a 3x3 kernel over a 2x2 input with a zero padding**
 399 **of 2 with unit-stride)**

400

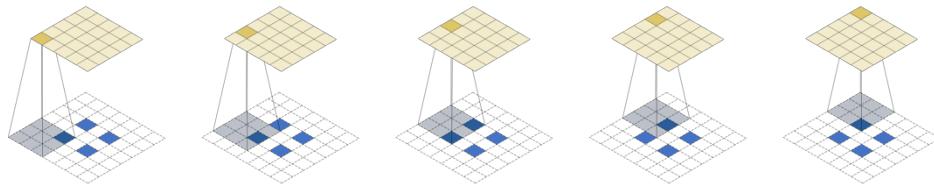
401 For a conventional convolution layer, forward propagation is ‘standard convolution’, which maps
 402 the input feature map into the output feature map as shown in Figure 12, while reverse propagation is
 403 ‘transposed convolution’ [44, 45], which calculates the gradient of the input feature map according to the
 404 gradient of the output feature map and keeps the size of the gradient map same as that of the feature map
 405 illustrated in Figure 13.

406 The process of transposed convolution is to some degree resemble to the deconvolution in that it
 407 produces the same spatial resolution as a deconvolutional layer. But the practical mathematical process
 408 carried out on the values is not the same. The transposed convolution guarantees that the output will be
 409 as large as before, while it is still carrying out a standard convolution process. To accomplish the goal,
 410 some fancy padding on the input is performed.

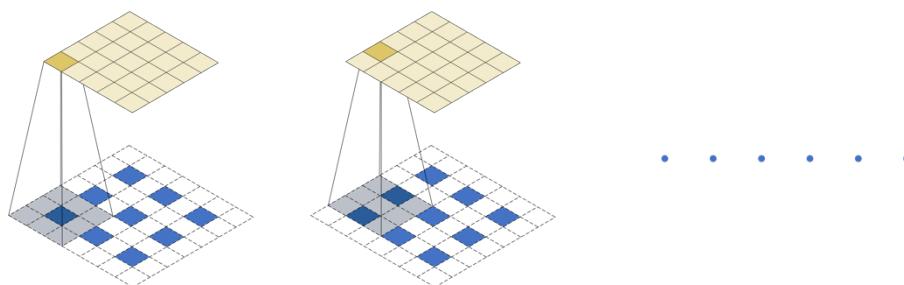
411 The transposed convolution step will not reverse the operation just mentioned, especially not
 412 concerning the numeric values. It only reconstructs the spatial resolution and carries out a convolution.
 413 This is very helpful for Encoder-Decoder architectures. Through utilizing this, upscaling of an image
 414 with a convolution can be achieved at one time, rather than two separate operations.

415 To conclude, if convolution is regarded as the mapping of the map of input size to the map of output
 416 size, transposed convolution is the mapping of the map of output size to the map of input size [46], which
 417 aims to restore the size. Some typical cases of transposed convolution are displayed in Figure 14.

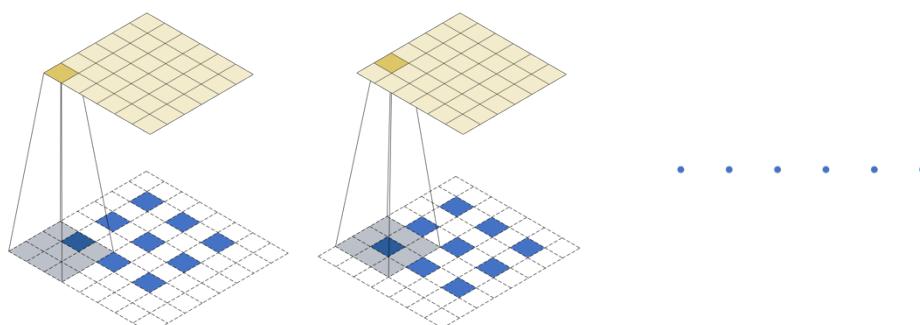
418



419
420 (a) Part of the transpose of convolution utilizing a 3×3 kernel over a 5×5 input utilizing 2×2 strides (same
421 to convolution utilizing a 3×3 kernel over a 2×2 input (gapped) with a zero padding of 2 utilizing unit-
422 stride)



423
424 (b) Part of the transpose of convolution utilizing a 3×3 kernel over a 5×5 input with a zero padding of 1
425 utilizing 2×2 strides (same to convolution utilizing a 3×3 kernel over a 3×3 input (gapped) with a zero
426 padding of 1 utilizing unit-stride)



427
428 (c) Part of the transpose of convolution utilizing a 3×3 kernel over a 6×6 input with a zero padding of 1
429 utilizing 2×2 strides (same to convolution utilizing a 3×3 kernel over a 3×3 input (gapped) with a zero
430 padding of 1 utilizing unit-stride)

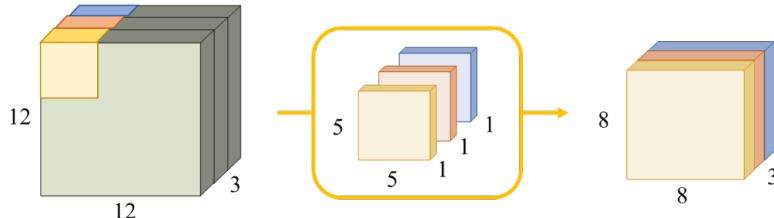
Figure 14 Typical cases of transposed convolution

431 3.1.4. Separable Convolution

432

433 In neural networks, separable convolution is commonly used in order to reduce the computation.
434 The convolution process can be divided into two steps: one depth wise convolution and one point wise
435 convolution [47]. In the depth wise convolution process, we use three convolution kernels of $5 \times 5 \times 1$ to
436 produce an $8 \times 8 \times 3$ image as illustrated in Figure 15.
437

438



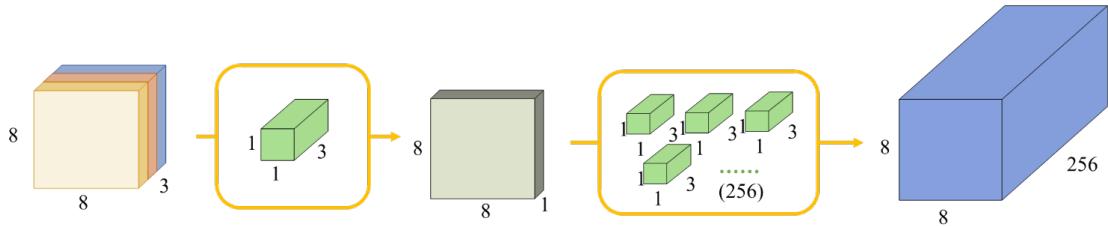
440

441 **Figure 15 Depth wise convolution: transform a 12x12x3 image to an 8x8x3 image by using 3 kernels**

442

443 In the pointwise convolution process, first, the convolution of 1x1x3 is used to check the calculation
 444 of each pixel, and the image of 8x8x1 is obtained. Then, 256 convolution kernels of 1x1x3 are used to
 445 produce the image of 8x8x256 as displayed in Figure 16.

446



447

448 **Figure 16 Pointwise convolution: transform an image from 3 channels to 1 channel, then output an
449 image with 256 channels by using 256 kernels**

450

451 Considering the change in computation, originally, 256 convolution kernels of 5x5x3 were moved
 452 8x8 times, that is,

$$453 \quad 256 \times 5 \times 5 \times 3 \times 8 \times 8 = 1,228,800 \quad (28)$$

454 times of multiplication.

455 During the utilization process of depth wise convolution, 3 convolution kernels of 5x5x1 move 8x8
 456 times, requiring

$$457 \quad 3 \times 5 \times 5 \times 1 \times 8 \times 8 = 4,800 \quad (29)$$

458 times of multiplication.

459 During the utilization process of pointwise convolution, 256 convolution kernels of 1x1x3 move
 460 8x8 times, requiring

$$461 \quad 256 \times 1 \times 1 \times 3 \times 8 \times 8 = 49,152 \quad (30)$$

462 times of multiplication calculation, which adds up to 53,952 times of calculation. Thus, the complete
 463 separable convolution process only needs 53,952 times of calculation which is much smaller than the
 464 computation of original process [48]. Separable convolution helps a lot in the application of mobile
 465 devices due to it uses the parameters efficiently [49].

466

467

3.2 Pooling layer

468

469 The function of pooling layer is reducing the resolution of feature maps. In other words, it aims to
 470 reduce the size of matrix generated by a convolutional layer. Typical categories of pooling will be
 471 introduced in the next [50].

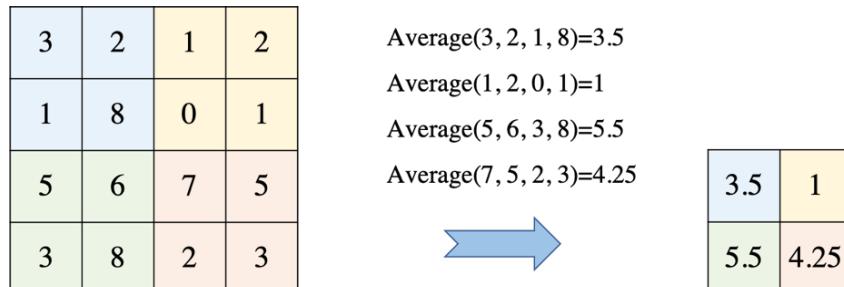
472

473 **3.2.1. Average Pooling**

474

475 In a simple term, average pooling means taking the average of a small range. Figure 17 indicates an
 476 instance of average pooling utilizing a kernel of 2x2 with a stride of 2. After pooling, the average value
 477 in the 2x2 kernel will be calculated and reserved in a rectified feature map. In this way, the 4x4 input is
 478 compressed to 2x2.

479



480 **Figure 17 Average Pooling using a 2x2 kernel with a stride of 2**

481

482 Error of feature extraction mainly caused by two sources. On the one hand, the increasing variance
 483 of the estimated value led from the limitation of neighborhood size. On the other hand, the parameter
 484 error from the convolutional layer will result in the deviation of the estimated mean value. Generally,
 485 average pooling to some degree solve the first problem by retaining more background information of the
 486 input image. Average pooling emphasizes the sub-layer sampling of the overall feature information,
 487 reducing the parameter dimension while ensuring the integrity of the information transferred. In a large
 488 model such as Dense Net, average pooling is frequently utilized to shrink the dimension and at the same
 489 time facilitate the transfer of information to the following module for extracting features [51].

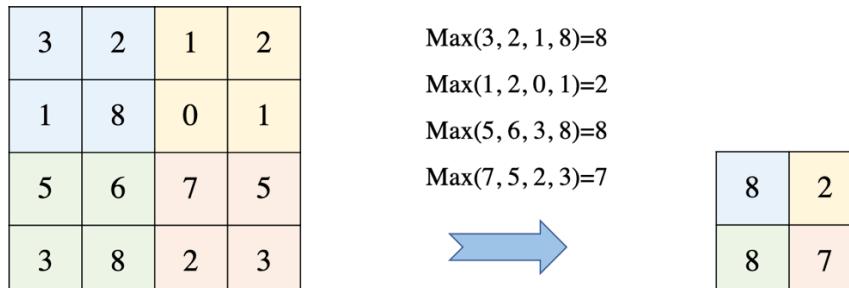
490

491 **3.2.2. Max Pooling**

492

493 In a simple term, max pooling means taking the maximum of a small range. Figure 18 indicates an
 494 instance of max pooling utilizing a kernel of 2x2 with a stride of 2. After pooling, the maximum value in
 495 the 2x2 kernel will be left in a rectified feature map. Eventually, the 4x4 input image is shrunk to 2x2.

496



497 **Figure 18 Max Pooling using a 2x2 kernel with a stride of 2**

498

499 Generally speaking, max pooling is more effective and commonly used in image processing. As it
 500 feels more like doing a feature selection. It selects features with better classification and identification as
 501 well as providing nonlinearity. Besides, max pooling is able to deal with the second error mentioned in

504 the average pooling paragraph, which means it can preserve more texture information.

505 For instance, if the feature we need to find is a dog, then as long as there is a dog in one area of the
506 image, it means that there is a dog in the whole image. Therefore, we should take the maximum matching
507 degree between all areas in the image and the characteristics of the dog, that is, maximum pooling should
508 be utilized.

509

510 3.2.3. Stochastic Pooling

511

512 In stochastic pooling, the pool mapping response is obtained by sampling the multinomial
513 distribution formed by the activation of each pool area. To be more exact, we compute the probabilities
514 p for each region j via normalizing the activations within the region using:

515
$$p_i = \frac{a_i}{\sum_{k \in R_j} a_k} \quad (31)$$

516 Where a refers to value of element

517 k refers to the index of element in region of feature map

518 We then sample from the multinomial distribution according to p to choose a location within the
519 region. The procedure is in Figure 19.

520

SUM(3, 0, 1, 1, 0, 0, 0, 5, 0)=10 Divide each by 10 (SUM) Sample a location from P()

3	0	1
1	0	0
0	5	0

0.3	0	0.1
0.1	0	0
0	0.5	0

0.3

521

522 **Figure 19 The transformation process of stochastic pooling**

523

524 The samples for each pooling region in each layer for each training example are drawn
525 independently to one another. When back propagation is carried out over the network, the gradient is
526 guided back to the pool area using the same selected location, similar to using the maximum pool for
527 back propagation. When passing information through the network, max pooling only focuses on the
528 strongest activation of each region. While, sometimes non-maximal activations should also be taken into
529 consideration. Stochastic pooling does well in this as those additional activations could be taken into
530 account.

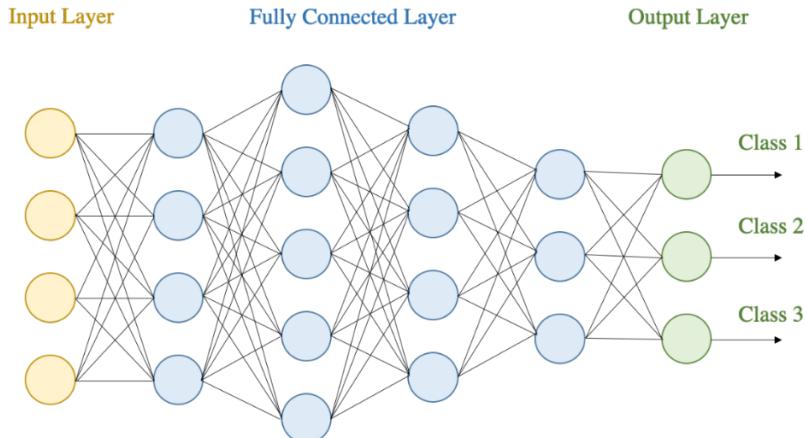
531 In a word, stochastic pooling is between the average pooling and max pooling. Probability is
532 assigned to pixel points according to the value size, and then sub-sampling is conducted based on the
533 probability. In the average meaning, it obeys average pooling criterion, and in the local meaning, it is
534 similar to max pooling [52].

535 Although pooling operation has little effect on improving the overall accuracy, it plays an obvious
536 role in decreasing parameters, reducing overfitting, optimizing model performance and saving
537 computational force, so it is an essential process in neural network design [53].

538

539 **3.3 Fully connected layer**

540



541

542 **Figure 20 A sample structure of FC layer**

543

544 Fully connected (FC) layers often appear near the end of CNN system and are utilized to sum the
545 features of the previous design. If we consider convolution and pooling in the front as process of feature
546 engineering, local amplification and extraction of local features. The FC layer in the back can be regarded
547 as feature weighting. A sample structure of FC layer as displayed in Figure 20 is usually a fast method
548 for learning nonlinear combinations of the high-level properties generated from convolutional layer. In
549 that space, the FC layer will learn about a possible nonlinear function.

550 The basic procedure of learning is as follow. Firstly, flatten the image that has already been
551 converted to a form suitable for a multilevel perceptron into a column vector and feeding it back to the
552 feedforward neural network. Then, apply the propagation after planarization to each iteration of training.
553 In this way, the model will have the ability to distinguish between the main features and some low-level
554 features in the image and classify them through classification techniques such as SoftMax [54].

555 In some circumstances, FC layer can be realized by conducting convolution operation, that is, we
556 can utilize a convolutional layer with 1x1 kernel size to replace the FC layer that is fully connected to
557 the front layer. Within normal FC layer, the important feature weight is increased to obtain the category
558 of the identified object. While in convolution operation, like an intentional weakening of fully connection,
559 the weak influence outside the local area is directly erased to zero on the basis of the inspiration of local
560 vision. There is also a bit of enforcement that the parameters utilized in different parts are the same. The
561 weakening reduces the parameters, saves the calculation, and specializes in local areas rather than seek
562 for more. While the enforcement can help further reduce the parameters.

563 After introducing the basic structure of CNN, the next section will focus on some classic CNN
564 architectures.

565

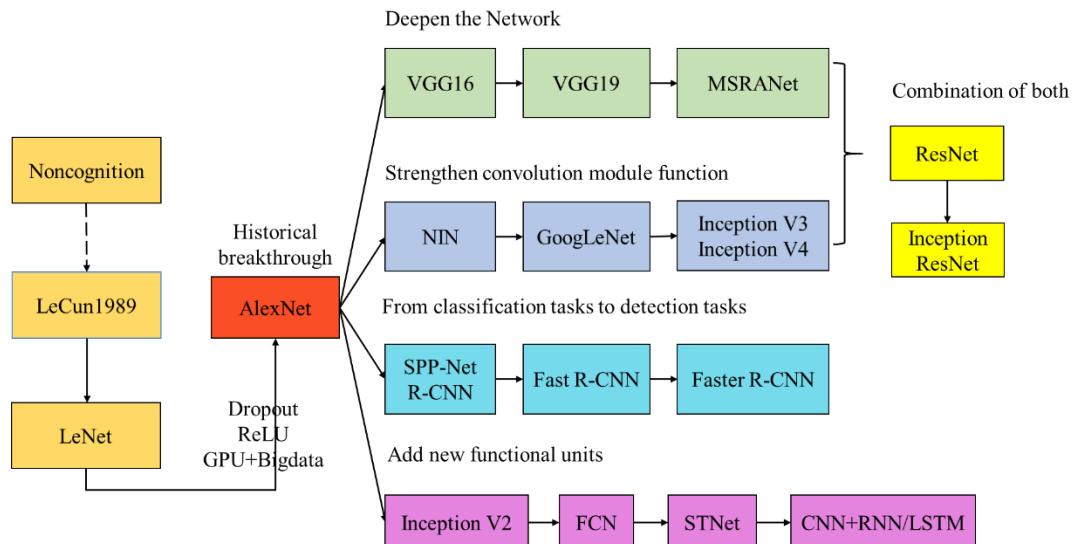
566 **4. CLASSIC ARCHITECTURE**

567

568 In the early stage, because of the lack of training data and computational ability, it is impossible to
569 train high performance CNN under no over-fitting circumstance. However, in the past few years, with
570 the introduction of ReLU and dropout, along with the historical opportunity achieved by GPU and big
571 data, the research on CNN has experienced several historic breakthroughs, began to obtain first-class
572 results and attracted attention of the world. Figure 21 displays several classic architectures of CNN in its

573 development history [55].

574



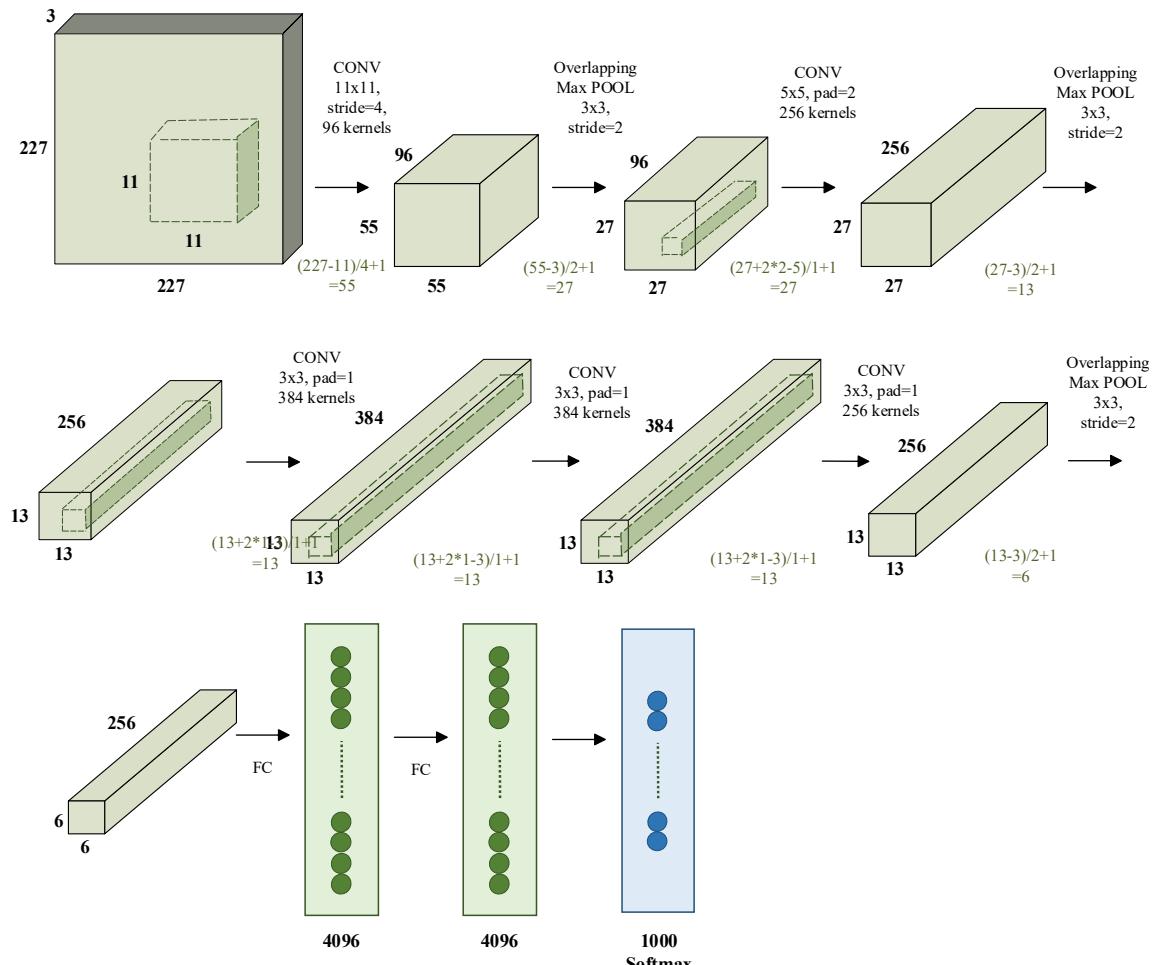
575

576 **Figure 21 Classic CNN architecture in development history**

577

578 4.1 AlexNet

579



580

581

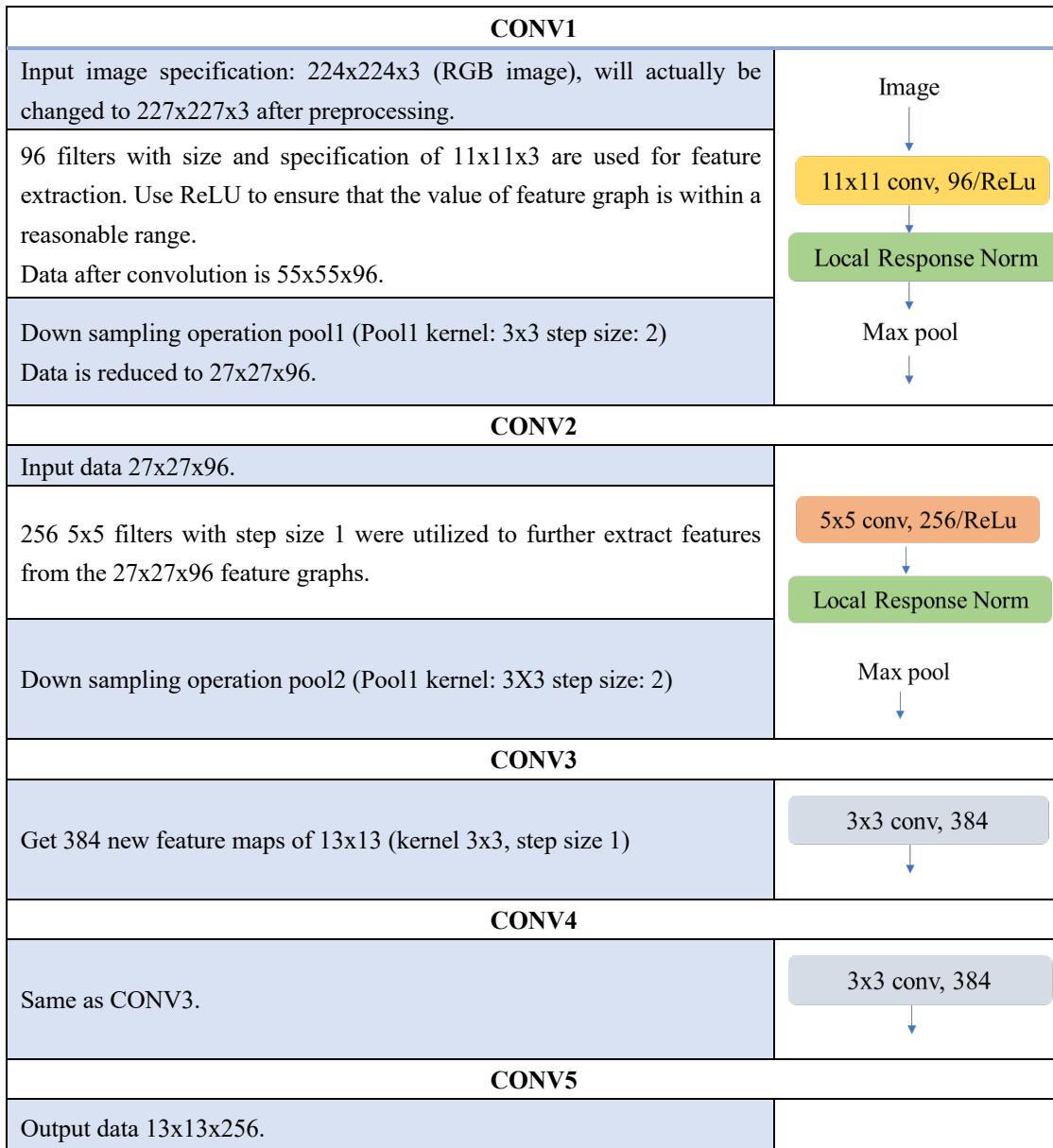
Figure 22 A sample structure for AlexNet

582

583 AlexNet came from a team led by [Krizhevsky, Sutskever \[16\]](#) in 2012 and its brief structure is
 584 illustrated in Figure 22. Due to the limited GPU capacity at that time, dual-GPU training was used, and
 585 the network was cut into two parts accordingly. However, with the development of GPU, single GPU can
 586 easily train AlexNet, so the two parts can be combined as shown in the figure. (Note that this leads to
 587 more connections and more computation) From this figure, it is obvious that the entire network structure
 588 of AlexNet consists of 5 convolutional layers and 3 FC layers, with a total depth of 8 layers [\[56, 57\]](#).

589 There is a clear process of direct conversion between layers and a simple interpretation of each layer
 590 is displayed in Figure 23. And then during a dropout process, where you randomly drop some information,
 591 you will get 4096 new neurons. The use of dropout can reduce overfitting and enhance the generalization
 592 of the model because it makes the model depend less on some local features. The dropout rate usually
 593 equals to 0.5.

594



Pool3 for down sampling operation (Pool3 kernel: 3x3 step size: 2) Data reduced to 6x6x256 [(13-3)/2+1=6]		
FC6	A full connection was made with 4096 neurons to 256 feature maps of size 6 by 6.	
FC7	Similar to FC6.	
FC8	A thousand neurons were used to fully connect the 4,096 neurons in FC7, which then went through a gaussian filter to get 1,000 float values, which is what we consider as predictive probability.	

Figure 23 Interpretation of each layer in AlexNet

595

596

597 4.2 VGG-Net

598

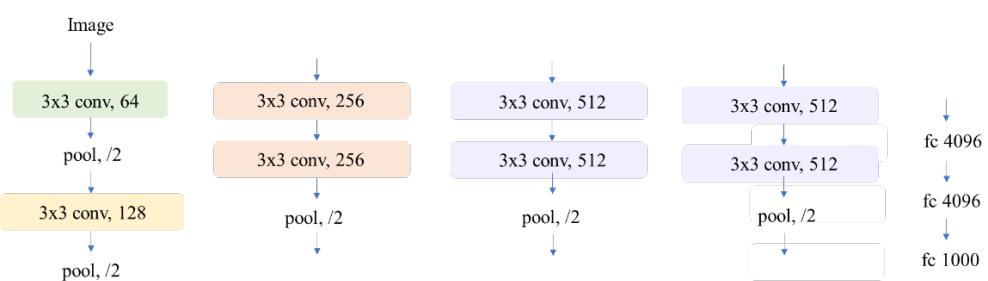
599 VGG-Net is a well-known deep CNN developed by [Simonyan and Zisserman \[58\]](#). It writes about
600 the relationship between the depth and performance of CNN [59]. Through repeatedly stacking the
601 convolutional layer and the maximum pooling layer, VGG-Net creates a 16~19 layers CNN which
602 achieves the runner-up in ILSVRC 2014 and the champion in positioning, with an error rate of 7.5% on
603 the top5. So far, VGG-Net has still been a widely used network to extract image features.

604

605 The innovation of VGG-Net is it uses 3x3 convolution kernel and 2x2 pooling kernel to improve
606 performance through deepening network structure. The growing number of network layers does not result
607 in an explosion in the number of parameters due to the reason that the number of parameters is
608 concentrated in the final three FC layers. Meanwhile, the series of two 3x3 convolutional layers can be
609 regarded as one 5x5 convolutional layer, and the series of three 3x3 convolutional layers can be regarded
610 as one 7x7 convolutional layer, that is, the size of the field of three 3x3 convolutional layers is equivalent
611 to one 7x7 convolutional layer. But the number of convolutional parameters of 3x3 is only about half of
612 7x7, and the former can have 3 nonlinear operations while the latter only has 1 nonlinear operation, which
613 makes the former have relatively stronger learning ability for features than the latter.

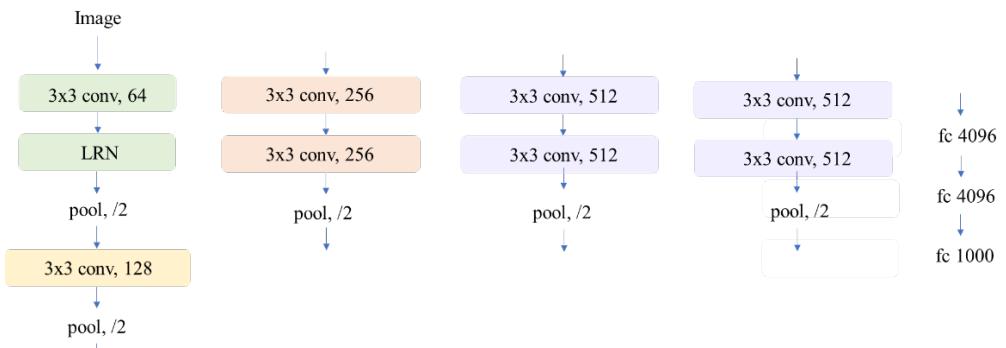
614

615 During training, VGG-Net first trains network of level A, and then initializes the later complex
616 models by repeatedly using the weight of network A, in this way, the convergence speed improves. The
configuration of VGG-Net is shown in Figure 24.

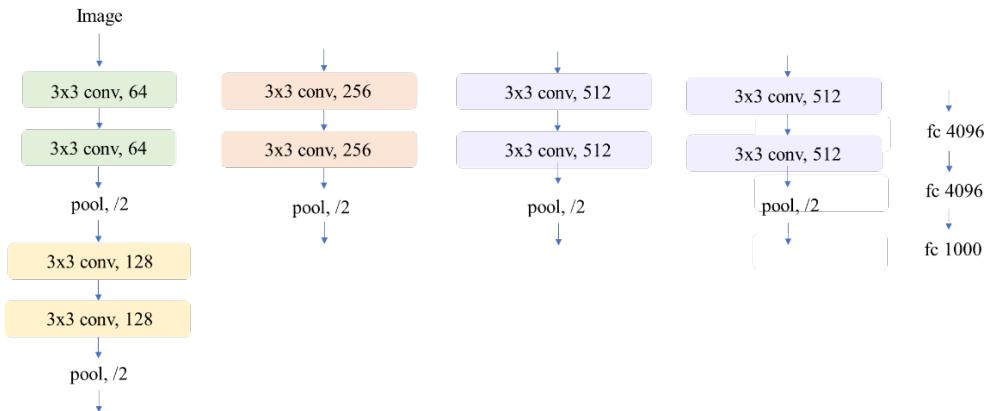


618

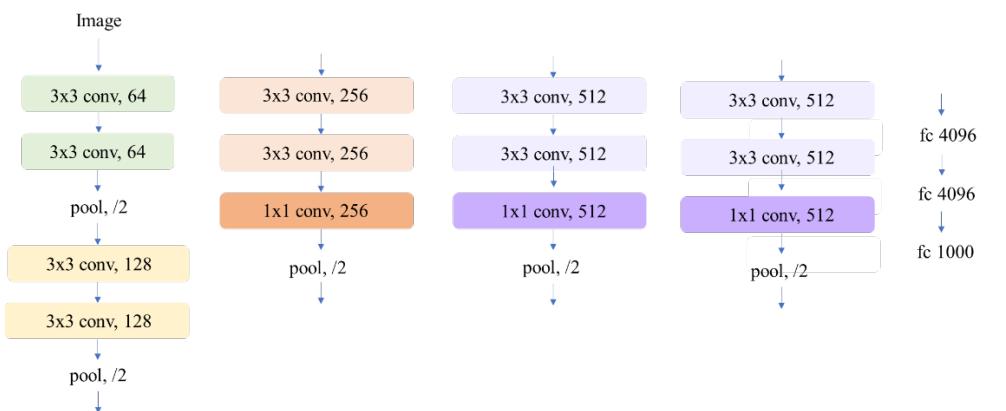
(a) A: 11 weight layers

619
620

(b) A-LRN: 11 weight layers

621
622

(c) B: 13 weight layers

623
624

(d) C: 16 weight layers

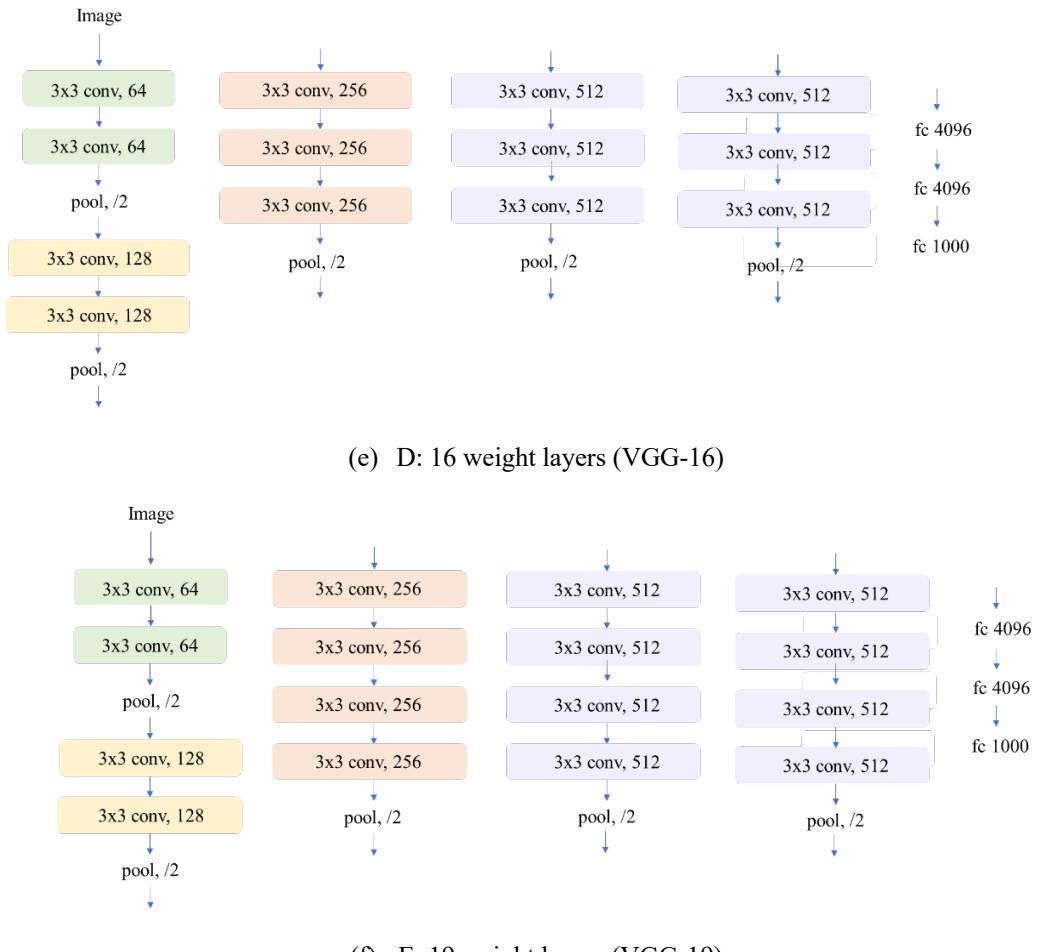


Figure 24 Configuration of VGG-Net

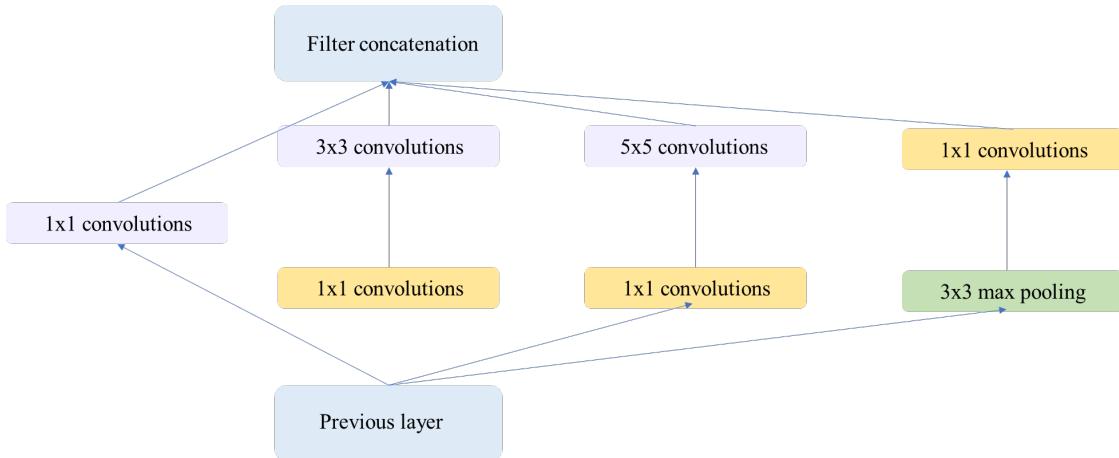
The authors of VGG-Net conclude that LRN layer has little effect. The experiment confirms that, when the network is going deeper, performance of network would be improved as well. The convolution of 1 by 1 is also efficient, allowing linear transformations of pixels without changing the number of output channels and the parameters of the neural network. But it doesn't work as well as the convolution of 3 by 3, because the convolution of 3 by 3 learns more about spatial features.

4.3 Inception (GoogLeNet)

Inception (also known as GoogLeNet) is a classic deep learning structure proposed by [Szegedy](#), [Ioffe](#) [17]. Before this, many structures mentioned before all utilize adding the depth of the network to achieve better training performance, however, utilizing more layers will at the same time bring loads of negative effects, for instance, overfitting, gradient disappearance, gradient explosion, etc. Inception improves the training results by proposing a new way which is making more efficient use of computing resources, in other words, extracting more features for the same amount of computation [60, 61].

The basic structure is illustrated in Figure 25. The entire inception structure is made up of several such inception modules in series. The inception structure has two main contributions: one is using 1×1 convolution to raise and lower the dimension, the other is carrying out convolution and reintegration on multiple scales at the same time.

649

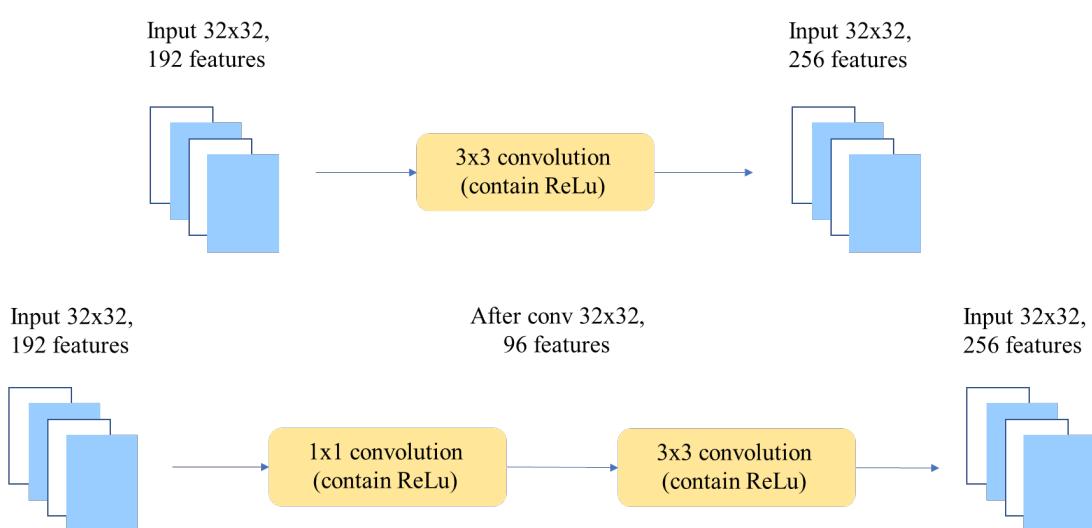


650

651 **Figure 25 Basic structure of the inception module**

652

653 There are multiple yellow 1x1 convolution blocks in the figure, which have two effects.

654 The first effect is they can help extracting richer features by adding more convolution to the same
655 size of the feeling field. This idea comes from the Network in the Network (NIN) which indicates that if
656 you put two convolutions in series, you can combine more nonlinear features. The structure of NIN is
657 likely to the multi-layer structure in the traditional neural network. The latter multi-layer spans the
658 receptive field of different sizes by adding pooling layer between layers, so as to extract features at a
659 higher scale. NIN structure is multi-layer in the same scale (without pooling layer in the middle), so it
660 can extract stronger nonlinearity in the same receptive field [62].661 The second effect is 1x1 convolution can be used to reduce the dimension, which helps decreasing
662 the computational complexity. More specifically, when a convolution layer has loads of input features,
663 the convolution operation on the input will produce much computation. However, if shrinking the
664 dimensionality of input first, and doing the convolution after the number of features are reduced, the
665 computation will significantly be reduced as well. Figure 26 compares the multiplication times of the
666 two schemes before and after optimization.
667

668

669 **Figure 26 The comparison between the multiplication times of the two schemes before and after**

670 optimization

671

672 With the same input of a set of 192 features, 32x32 size, and the output of 256 features, the first
673 scheme is directly realized by 3x3 convolution, which requires

674 $192 \times 256 \times 3 \times 3 \times 32 \times 32 = 452984832$ multiplications (32)

675 In the second scheme, the convolution of 1x1 was first utilized to decrease the number of features
676 to 96, and then 3x3 was used to restore 256 features, which required

677 $192 \times 96 \times 1 \times 1 \times 32 \times 32 + 96 \times 256 \times 3 \times 32 \times 32 = 245366784$ multiplications (33)

678 The dimensionality reduction method of 1x1 was used to save half of the computation. Besides, if
679 the number of features in the final output remains the same (256 sets), the dimensionality reduction in
680 the middle is similar to the effect of compression and will not influence the result of the eventual training.

681 Apart from the 1x1 convolution blocks, the input layer has four branches, which were convolved or
682 pooled with filters of different sizes, and finally integrated together on the feature dimension. Szegedy
683 explained the benefits of this new structure from a number of perspectives.

684 The first explanation is the features of different scales can be extracted by convolution on multiple
685 scales at the same time. The richer the features, the more accurate the final classification.

686 The second explanation is utilizing the principle of decomposition of sparse matrix into dense matrix
687 to accelerate the convergence. For instance, the first matrix in Equation (34) is a sparse matrix (many
688 elements are 0 and are unevenly distributed in the matrix), which is convolved with a 2x2 matrix and
689 needs to be calculated for each element in the sparse matrix. However, if the sparse matrix is decomposed
690 into 2 sub dense matrices and then convolved with the 2x2 matrix, the regions with more 0 in the sparse
691 matrix can be eliminated from calculation, greatly reducing the amount of calculation. Applying this
692 principle to Inception needs decomposing on the feature dimension. The input data of the traditional
693 convolution layer convolved with the convolution kernel of one scale (such as 3x3) and output fixed
694 dimensions data (such as 256 features). All 256 output features were uniformly distributed on the scale
695 of 3x3, which can be understood as the output of a sparsely distributed feature set. While Inception
696 module extract features on multiple scales (such as 1x1, 3x3, 5x5), thus the output features of 256 is not
697 evenly distributed, but the features that have a strong correlation gather together (such as 96 1x1 features
698 together, 96 3x3 features together and 64 5x5 features together). This can be understood as multiple
699 densely distributed sub-feature sets. In this feature set, as the more relevant features cluster together, the
700 non-relevant non-key features weaken, so that the Inception method outputs the same 256 features with
701 less ‘redundant’ information. With such ‘pure’ feature set transferred layer by layer and eventually used
702 as input of reverse calculation, the natural convergence rate becomes faster.

703
$$\begin{bmatrix} 3 & 4 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 3 & 0 & 0 \\ 0 & 0 & 0 & 1 & 4 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \otimes \begin{bmatrix} 3 & 1 \\ 2 & 4 \end{bmatrix} \Leftrightarrow \begin{bmatrix} [3 & 4] \otimes [3 & 1] \\ [2 & 3] \otimes [2 & 4] \\ [0 & 1] \otimes [3 & 1] \\ [1 & 0] \otimes [2 & 4] \end{bmatrix}$$
 (34)

704 Where \otimes represents Kronecker Product which is used for operations between two matrices of any size.
705 \Leftrightarrow represents Substantively equivalent.

706 The third explanation is according to Hebbin's principle [63]. The Hebbin's principle is a theory in
707 neuroscience that explains what happens to neurons in the brain during the learning process: fire together,
708 wire together. Hebbin believes that ‘Two neurons, or systems of neurons, that are always firing at the
709 same time form a ‘combination’ in which the firing of one neuron promotes the firing of the other’. For

example, dogs drool when they see meat. After repeated stimulation, the neurons in the brain that recognize meat and those responsible for salivating will promote each other and ‘entwine’ with each other, making them drool faster when they see meat in the future. The usage of this principle in Inception structures is to bring together highly relevant features. It's a little bit like the second explanation above, where you separate the features of 1x1, 3x3, and 5x5. Since the ultimate goal of training convergence is to extract independent features, the convergence speed can be accelerated as long as the highly relevant features be gathered in advance.

717

718 4.4 ResNet

719

720 ResNet was proposed by [He, Zhang \[64\]](#) from Microsoft research institute in Dec, 2015. Through
721 utilizing residual Unit, ResNet successfully trained the 152-layer neural network, and achieved the best
722 result in ILSVRC 2015. Meanwhile, the number of parameters was lower than VGG-Net, with
723 outstanding effect. The structure of ResNet is able to greatly speed up the training of ultra-DNN, as well
724 as perfect the model accuracy.

725 The initial inspiration of ResNet comes from this problem: if depth of neural network is
726 continuously increasing, a degradation problem will occur. This means the accuracy will first rise and
727 then reach saturation, afterwards the continuous increase of depth will cause the decline of accuracy. The
728 difference between this problem and overfitting is, in this problem, the error increases in both testing set
729 and training set.

730 If a relatively shallow network reaches saturated accuracy, then adding several $y = x$ congruent
731 mapping layers to the end will at least not increase the error, which means, the deeper network will not
732 cause the increasing of error on training set. The idea of utilizing congruent maps to transfer the previous
733 layer of output directly to next layer was the inspiration for ResNet. If we directly pass the input to the
734 output as the beginning value, the target we require to study at this point will be:

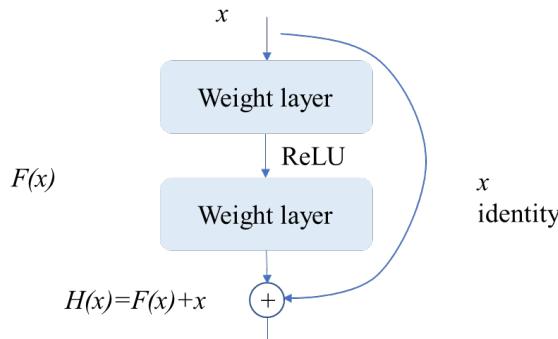
$$735 \quad F(x) = H(x) - x \quad (35)$$

736 Where x refers to input of neural network

737 $H(x)$ refers to expected output

738 The Residual Unit of ResNet is illustrated in Figure 27. ResNet altered the learning objective from
739 learning a complete output to learning only the difference between output and input, thus simplifying the
740 learning objectives and difficulty.

741



742

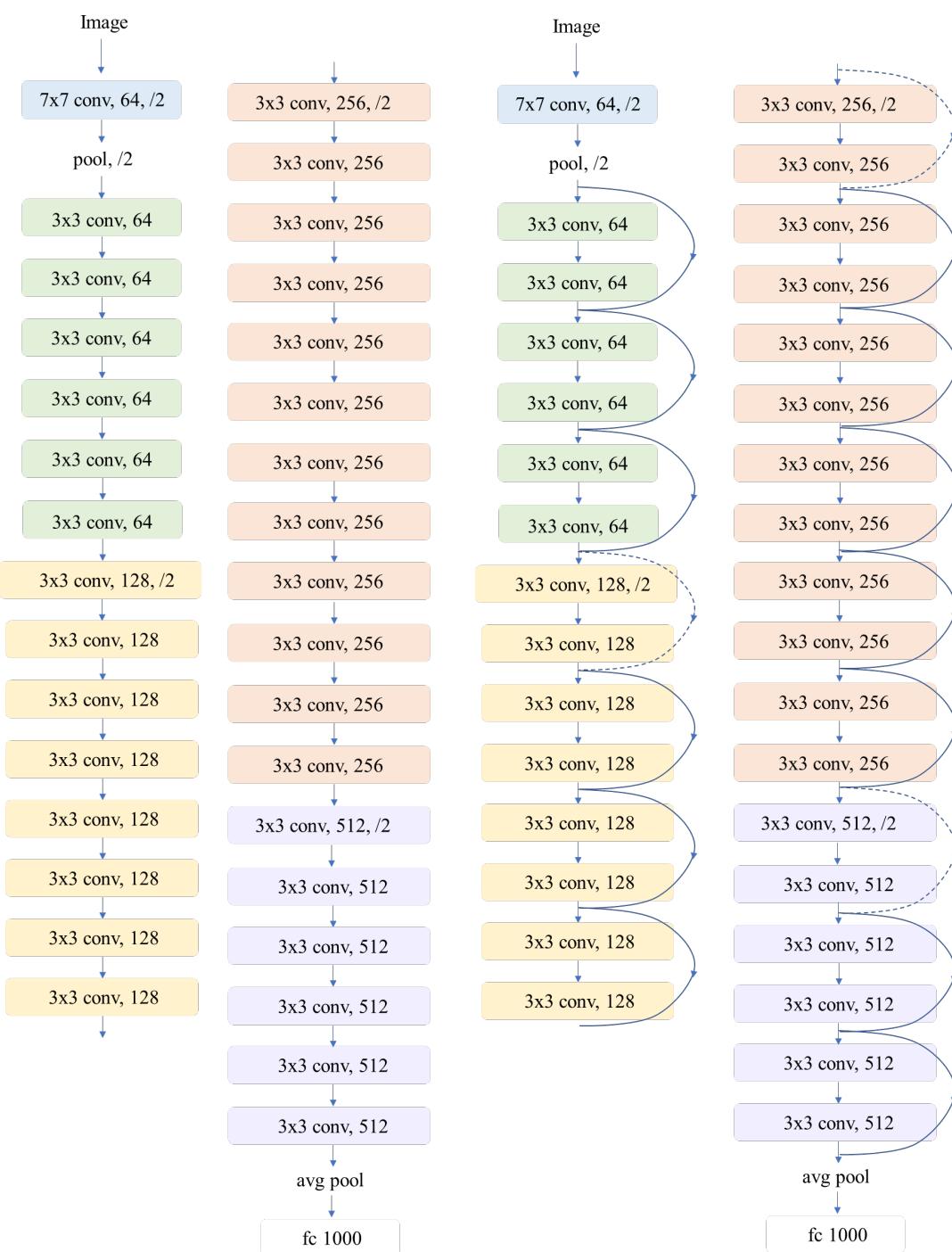
743 **Figure 27 A sample structure of residual block**

744

745 The Residual block is realized by shortcut connection as displayed in Figure 27. Through utilizing

746 shortcut, this block performs an element-wise superposition of the input and output. This simple addition
 747 will not increase additional parameters and computation to network. However, it could efficiently
 748 improve the training speed and training effect of the model. In addition, if the layer number of model is
 749 increasing, this structure is also able to efficiently handle in the problem of degradation.

Figure 28 shows a comparison between a 34-layer ordinary CNN and a 34-layer ResNet network. It can be seen clearly that the biggest difference between a straight-through CNN and ResNet is that ResNet has a large amount of bypass branches which transmit the input straightway to the subsequent layer, allowing the subsequent layer to directly study the residual. This category of structure is as well-known as shortcut or skip connection.



(a)

(b)

Figure 28 Comparison between a 34-layer ordinary CNN and a 34-layer ResNet network.

4.5 U-Net

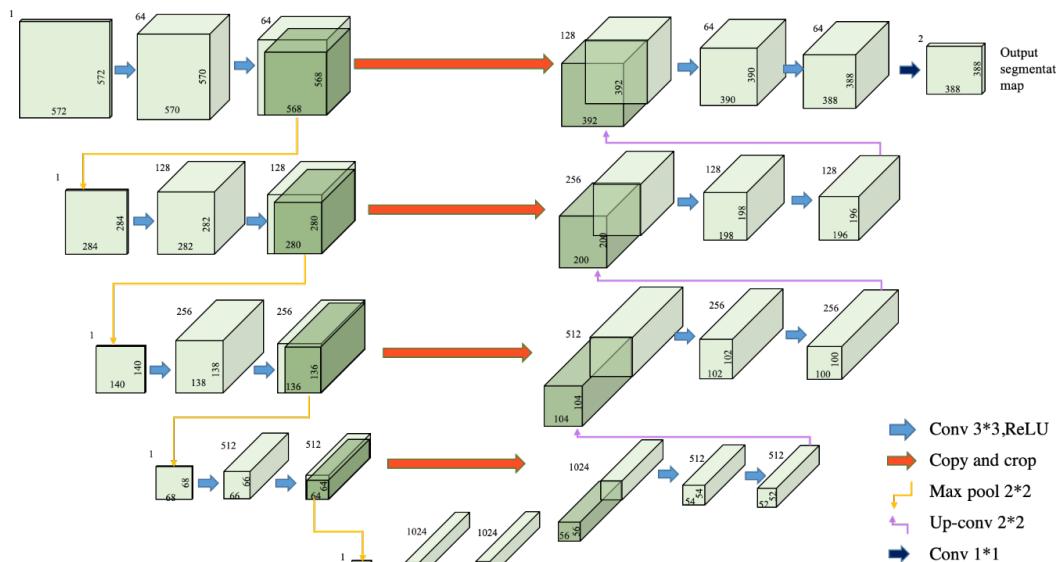


Figure 29 The basic structure of U-Net

U-Net is a classical full convolutional network. This U-shaped structure can capture semantic information in the process of down-sampling and accurately locate the corresponding up-sampling process. This network can train end-to-end with a small number of images, but the effect exceeds that of the previous best model, CNN based on sliding Windows [65]. The U-shaped structure is shown in Figure 29.

Assuming that the input image tile is of 572x572. The left side of the network composed of four blocks is called contracting path, which is a series of reduction sampling operations formed by convolution and max pooling. Finally, a feature map with a size of 32x32 is obtained.

Network on the right-hand side is named expansive path. It is as well composed of four blocks. Before the start of each block, it will use deconvolution to enlarge two times the size of the feature map, and the number in half. Then, it is merged with the feature map of the symmetric contracting path on the left due to the reason that the size of the feature map of the contracting path on the left is different from that of the expansive path on the right. The normalization of U-Net is conducted by clipping feature map of the contracting path to the same size as feature map of the expansive path (i.e., the dark green block in Figure 29). Effective convolution operation is still used in convolution operation of extension path, and the final size of feature map is 388x388.

U-Net mainly uses SoftMax cross entropy loss to deal with the problem that similar target boundaries of medical images are not easy to be recognized. Besides, weights are added to each pixel in the objective function to make the network more able to distinguish the boundaries. The equations are as follows:

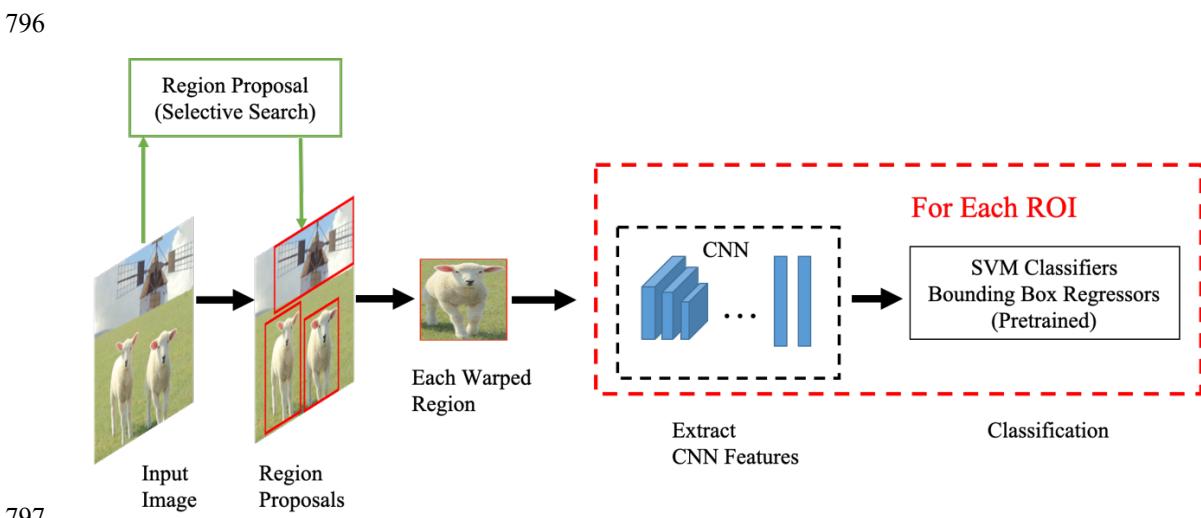
$$E = \sum_{x \in \Omega} w(x) \log(p_{l(x)}(x)) \quad (36)$$

$$w(x) = w_c(x) + w_0 \cdot \exp\left(-\frac{(d_1(x)+d_2(x))^2}{2\sigma^2}\right) \quad (37)$$

785 Where $p_{l(x)}(x)$ is loss function of SoftMax
 786 l refers to the label value of a pixel
 787 w refers to the weight value of a pixel
 788 w_c is the weight that balances the proportion of the category
 789 d_1 is the distance from the pixel to the nearest cell
 790 d_2 is the distance from the pixel to the second closest cell.
 791 w_0 and σ are constants.
 792

793 4.6 R-CNN Series

794 4.6.1. R-CNN



797
798 **Figure 30 The basic structure of R-CNN**
799

800 The general procedure for target detection of R-CNN is illustrated in Figure 30: (i) Extract regional
 801 proposals from input images using selective search algorithm and warp them to a fixed size. Selective
 802 search synthesizes the methods of both continues search and segmentation. It aims at identifying possible
 803 target locations for object recognition. (ii) Input the normalized regional proposals into the CNN network
 804 and extract features. (iii) SVM classification is used for identification of these features, and linear
 805 regression is applied for fine-tuning the border position and size, with each category separately training
 806 a bounding box regression. By using less and higher quality region of interest (ROI), R-CNN performs
 807 better than the system using traditional sliding window approach in the evaluation of accuracy and speed
 808 [66].

809 However R-CNN framework still exist the following problems: (i) it is very time-consuming and
 810 has lots of repeated computation as each candidate box needs to go through CNN alone; (ii) its training
 811 is divided into several stages and the steps are tedious; (iii) it takes up disk space: the convolved
 812 characteristics of the data needs to be kept separately [67].

813 4.6.2. SPP-Net

815

816 In order to solve the problems of R-CNN, SPP-Net made the following improvements.

817 One of the main ideas of SPP-Net is to remove the operations such as crop/ warp on the original
 818 image and replace them with the spatial pyramid pooled layer (SPP) on the convolution feature. In R-
 819 CNN, the proposal was warped before it was sent to the net. Part of the reason for the introduction of
 820 SPP layer is that different levels of crop/ warp on images in R-CNN will not only take time and effort
 821 but also cause some problems as shown in Figure 31. For example, the crop in the figure will cause the
 822 object to be incomplete, and warp will cause serious deformation after the object is stretched [68].
 823

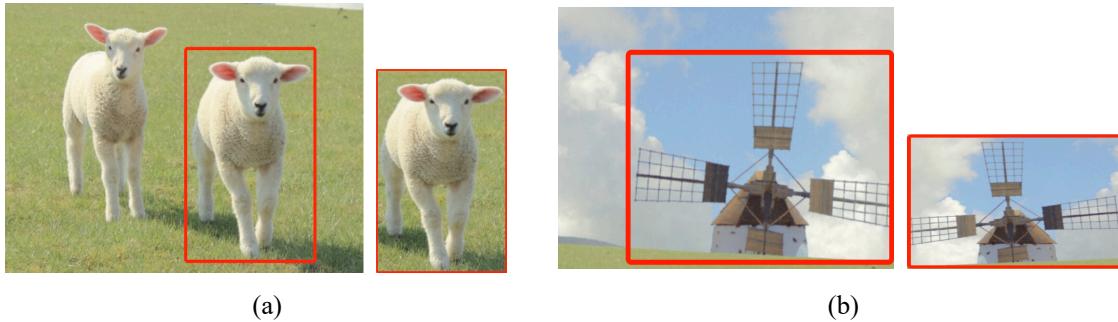


Figure 31 (a) Crop; (b) Warp

824
 825 However, the more important reason is that in R-CNN, each regional proposal is subject to a CNN
 826 feature extraction operation, which brings a lot of double calculations. In SPP-Net, we only performed
 827 CNN feature extraction on the whole image once, and then map the proposal to the feature map generated
 828 by the convolutional network, which is equivalent to feature map obtained by image in a proposal going
 829 through the layer by layer convolution process.
 830

831

832 **4.6.3. Fast R-CNN**

833

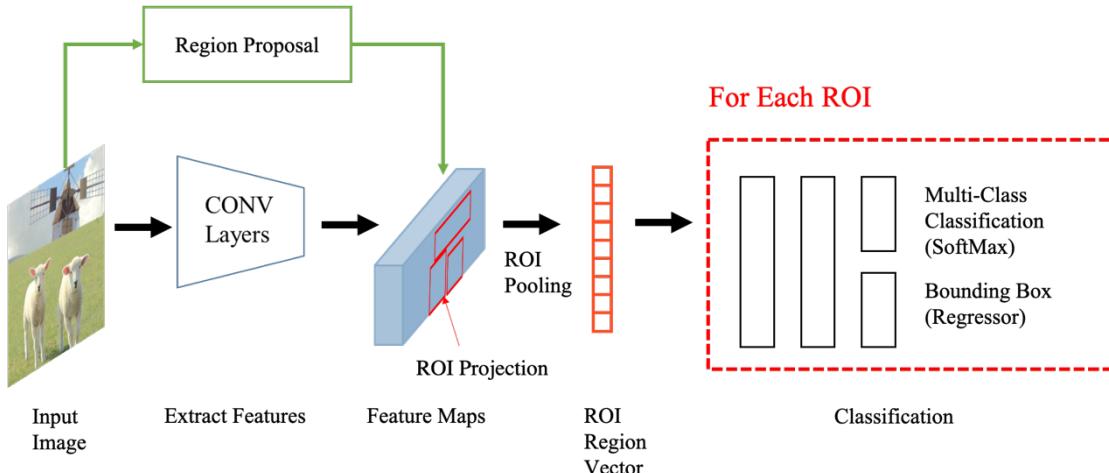


Figure 32 The basic structure of fast R-CNN

834

835 It could be obviously seen in Figure 32 that, fast R-CNN is also a method that makes some
 836 perfections based on R-CNN. Similar to SPP-Net, it only conducts CNN feature extraction for the whole
 837 image once, and then utilizes an ROI pooling layer. While, due to the reason that it is not a fixed size
 838 input, the size of the pooling grid shall be adjusted dynamically each time to achieve regional
 839 normalization. Fast R-CNN uses SoftMax to replace SVM for classification, as well as multi-task loss
 840
 841

842 function to carry out border regression and classification together [69].

843 The main steps are as follows: (i) Take the whole picture as the input and use CNN to get the feature
844 layer of the image. (ii) Extract regional candidate boxes from the original image by selective search and
845 projected onto the final feature layer one by one. (iii) ROI pooling operation is conducted for each area
846 candidate box on the feature layer to achieve feature representation of fixed size. ROI pooling is able to
847 accelerate the speed of training and testing. (iv) After two FC layers, utilize SoftMax to conduct multi-
848 target classification respectively, and apply regression model in fine-tuning the border position and size.

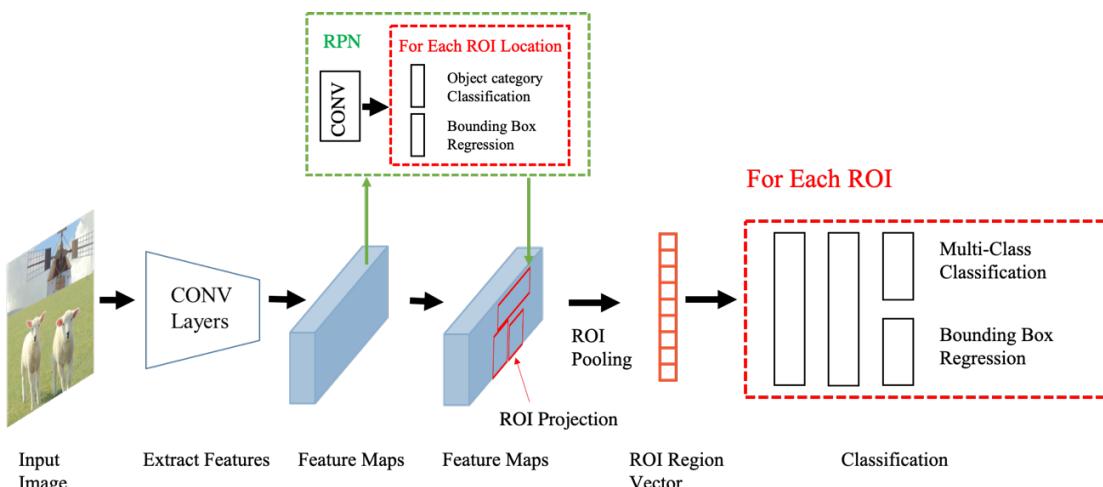
849 The disadvantage of fast R-CNN is that region proposal is extracted by selective search, and most
850 of the target detection time is spent on it (normally, if a region proposal is proposed 2~3s once, the feature
851 classification will be proposed only 0.32s once), which could not reach the requirements of real-time
852 application.

853 In general, fast R-CNN uses feature CNN as extractors to extract features of the whole image first.
854 And then it applies the method of creating candidate regions directly to the extracted feature map. For
855 example, fast R-CNN chose feature map output by convolutive layer Conv5 in VGG-16 to generate ROI.
856 These ROI would then be combined with corresponding feature graphs to cut out feature blocks and used
857 in target detection tasks. We use ROI pooling to convert feature blocks to fixed sizes and feed them to
858 the FC layer for classification and positioning. Because fast R-CNN has no repeated feature extraction
859 operation, it saves lots of processing time. The training speed of fast R-CNN is 10 times faster than R-
860 CNN, and the predicting speed is 150 times that of R-CNN.

861

862 4.6.4. Faster R-CNN

863



864

865

866

867 As displayed in Figure 33, faster R-CNN utilizes the proposal achieved by Region Proposal
868 Network (RPN) to replace the proposal achieved by selective search.

869 The function of RPN network is outputting a batch of rectangular candidate regions with the input
870 of an image, which is somewhat alike to the selective search in the previous target detection. The network
871 structure is according to a CNN, but includes the SoftMax and Bounding box regression multi-task
872 models [70].

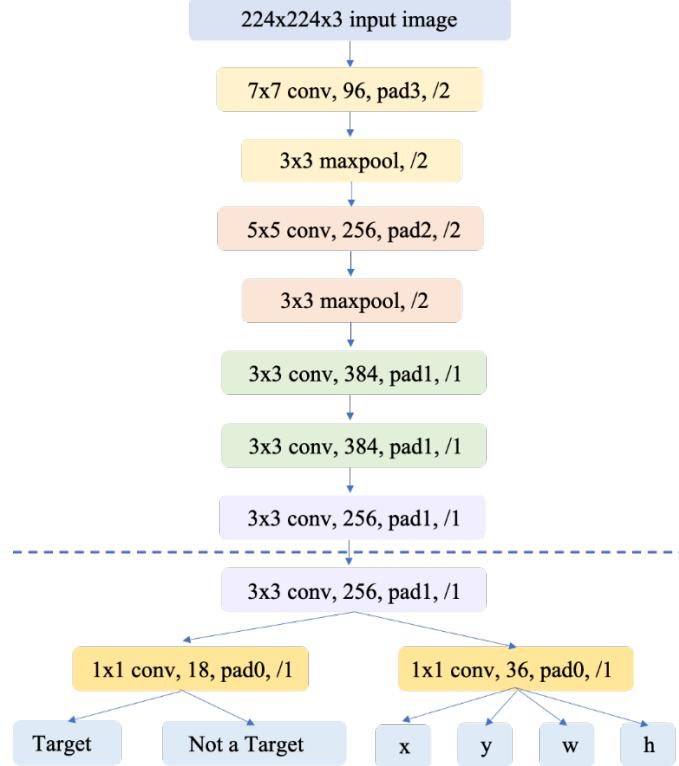


Figure 34 The basic structure of RPN Network

873

874

875

876 In Figure 34, above the dashed line is the structure in front of the last convolutional layer of ZF
 877 network. Under the dashed line is the particular structure of RPN network. For the particular structure,
 878 the convolution is 3×3 , and then the output is divided into two ways through 1×1 convolution, one way
 879 of which outputs the probability of target and non-target, and the other way outputs the four parameters
 880 related to box.

881 The procedure of faster R-CNN is: (i) Same as fast R-CNN. (ii) K different Anchor boxes are used
 882 in the final convolution feature layer to conduct a region proposal, that is, object/non-object binary
 883 classification is carried out for the corresponding area of each anchor box, and k regression models are
 884 used for fine-tuning the location and size of the candidate box. Regional proposal (via RPN) is generated
 885 on the feature map, with about 300 suggestion windows per image. iii) Same as fast R-CNN. iv) Classify
 886 the objects and do border regression.

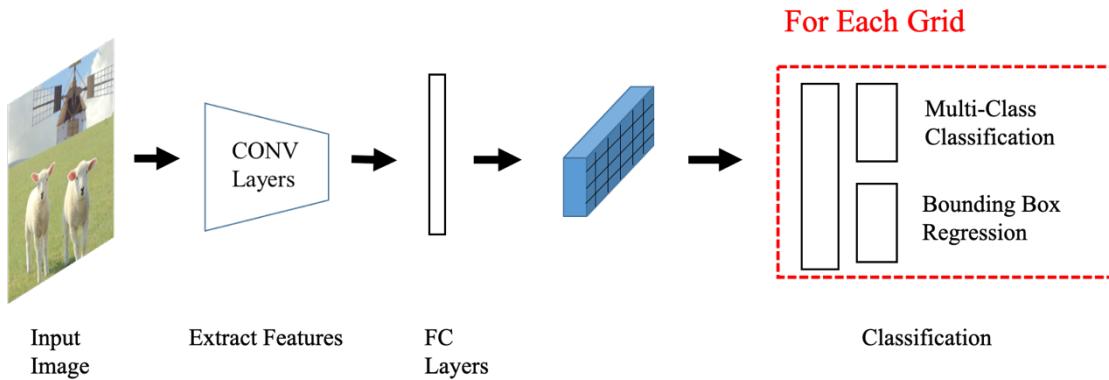
887

888 4.7 YOLO

889

890 YOLO series algorithm is a typical one-stage target detection algorithm. Its basic structure is
 891 illustrated in Figure 35.

892



893

894

895

896 2.3.6.1 YOLOv1

897

898 YOLOv1 is an end-to-end target detection approach, and the determination of position and category
 899 is accomplished by direct regression without need for intermediate regional proposal to find the target.
 900 If two small targets fall into a cell at the same time, the model can only predict one. Large object
 901 intersection over union (IOU) error and small object IOU error are close to the loss contribution value in
 902 network training. Thus, for small object, small IOU error will also bring a great influence on the network
 903 optimization, reducing positioning accuracy of object detection. YOLO uses multiple lower sampling
 904 layers, and the object features obtained from network are not precise, which will also affect the detection
 905 effect [71].

906 Compared with faster R-CNN, YOLOv1 has more obvious advantages. Firstly, faster R-CNN needs
 907 an RPN network to find the candidate region, while for YOLOv1, after the input image has been divided,
 908 it directly takes the region as the candidate region. YOLOv1 utilizes the end-to-end regression method
 909 rather than use those regional proposal steps, so direct regression accomplishes the determination of
 910 position and category. Therefore, compared with fast R-CNN, the main advantage of YOLOv1 lies in (i)
 911 its speed; (ii) global processing makes background errors less than regional based methods; (iii) good
 912 generalization performance. However, YOLOv1 performs not very well in target positioning process,
 913 which causes low detection accuracy of YOLO.

914 To sum up, YOLOv1 proposes a new idea of a single stage. Compared with the two-stage method,
 915 YOLOv1 has obvious advantages in speed and impressive real-time characteristics. However, YOLOv1
 916 also has some problems, such as rough grid division and the number of boxes that could be generated by
 917 each grid, which limits the detection of small-scale objects and objects that are close to each other.

918

919 2.3.6.2 YOLOv2

920

921 The differences between YOLOv2 and YOLOv1 are as follows:

- 922 1) The dropout layer was removed, and Batch Normalization was added in each convolutive layer,
 923 with a 2% improvement in Mean Average Precision (mAP).
- 924 2) High-resolution classifier. Input data of high-resolution changed from 224*224 to 448*448, and
 925 mAP improved by 4%.
- 926 3) YOLOv2 deletes a pooling layer to get high resolution data. Use anchor to replace FC layer.
- 927 4) Anchor box is determined by k-mean clustering and distance calculation of clustering is used

928 instead of Euclidean distance.

929 5) YOLOv2 uses anchor box to predict the object category rather than use grid to predict the
930 probability of each grid in YOLOv1.

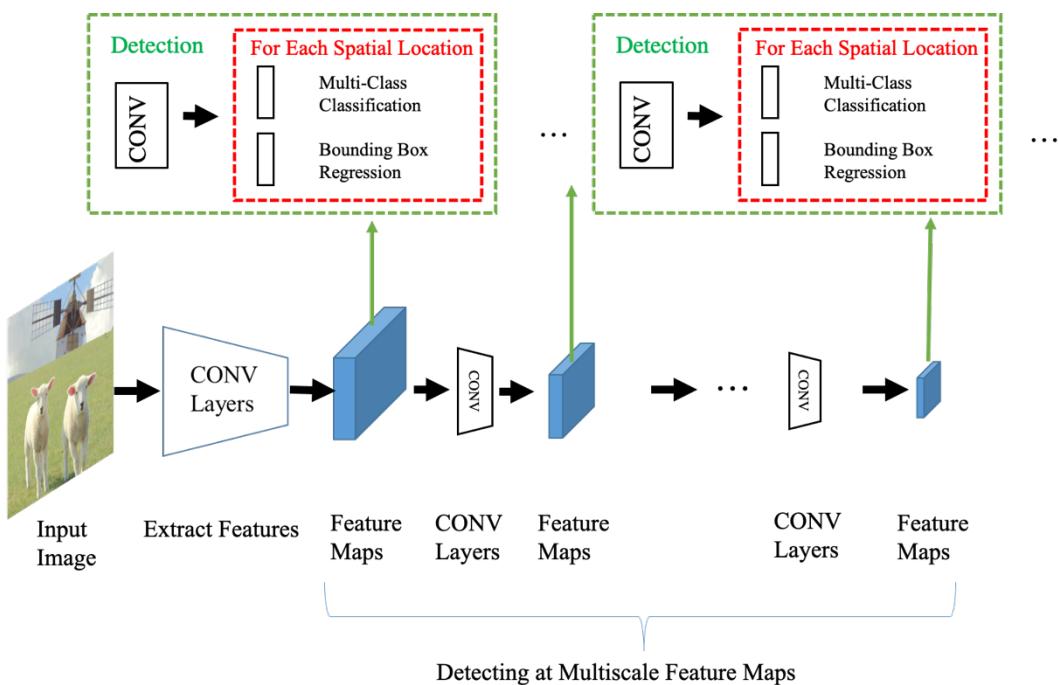
931 6) The bounding Box coordinate continues YOLOv1's prediction of Bounding boxes' coordinates
932 with respect to the upper left corner of the grid. To ensure that the center of bounding Box falls within
933 the grid (i.e. its coordinates fall between 0 and 1), Bounding Box center is sent into the activation function.

934 7) YOLOv2 adds a passthrough layer to improve the detection of small objects. This layer conducts
935 information fusion of high-resolution feature and low-resolution feature. Through the stacked in adjacent
936 features rather than the spatial location of different channel, it integrates characteristic of high and low
937 resolution, making the mAP value increased by 1%.

938 8) YOLOv2 only uses convolutional layer and pooling layer, thus the size of input image could be
939 dynamically adjusted. Besides, the accuracy and speed can be adjusted according to own requirements.
940

941 4.8 SSD

942



943

944 **Figure 36 The basic structure of SSD**

945

946 YOLO predicts a location utilizing full-image features, SSD predicts a location utilizing features
947 around that location [72].

948 According to Figure 36, it seems that SSD model is somewhat alike to RPN, for instance, the Default
949 Bounding Box in SSD is very likely to the anchor in RPN. However, SSD takes into account different
950 scales in different feature layers, and RPN thinks about different scales in one feature layer.

951 SSD outputs a series of discretization boxes, produced on feature maps on different layers, with
952 different aspect ratio.

953 This section introduced the classic architectures in CNN from an early stage in nearly 10 years ago
954 till now. The CNN model has been continuously evolved and eliminated in the careful study of countless
955 scientists, thus forming various model systems with their own merits. The following section will focus

956 on how different model made remarkable achievements and outstanding contributions in different
957 medical fields.

958

959 **5. APPLICATION**

960

961 Since CNN have some many models and it can be used in lots of fields of image classification like
962 identifying human face and patient condition. The following content will introduce the application of
963 CNN for classifying the patient conditions in different body parts.

964

965 **5.1. Breast**

966

967 According to the statistical analysis defined from Institute for Health Metrics and Evaluation
968 (IHME), there was a noticeable increase in breast cancer mortality in most regions of the world between
969 1990 and 2015, with a total mortality rate of 0.7 per 100,000. Data also indicates that the five-year
970 survival rate of early stage breast cancer patients exceeds 95%, while the five-year survival rate of late
971 period breast cancer patients is less than 20%. The World Health Organization (WHO) has confirmed
972 that early diagnosis and treatment of breast cancer patients could not only decrease mortality, but also
973 help some patients avoid psychological suffering and financial loss from mastectomy and axillary lymph
974 node dissection [73-75]. Thus, early detection is a critical step in prevention and treatment of breast
975 cancer.

976 Mammography is a specific imaging method to screen and diagnose breast cancer. [Lévy and Jain](#)
977 [76] and his group used CNN to straightway classify pre-detected breast masses in mammograms as
978 benign or malignant. They utilized AlexNet, GoogleNet model and achieve outstanding precision rate
979 result as shown in Table 1. Another project focuses on applying CNN to detect microcalcifications in
980 mammograms, since CNN shows excellent performance in classifying image data. In order to solve the
981 problem of large class imbalance between pixels from microcalcifications and other breast tissue, the
982 research team of [Mordang, Janssen](#) [77] applied a CNN architecture which is inspired by OxfordNet to
983 classify the microcalcification pixels which are difficult to be identified. The researchers also used
984 cascade classifier for comparison and find out that CNN classifier had significantly better results for
985 different data set with a sensitivity increase of about 6%. Except the Mammography, breast cancer
986 cytological specimen can also be analyzed in order to diagnose benign and malignant cases. To improve
987 the diagnosis of breast cancer, team of [Gao, Wu](#) [78] presented a network in combination of a shallow
988 CNN and a deep CNN. They developed the shallow one to achieve ‘virtual’ recombined images from
989 input dataset and apply the deep CNN to extract features. Finally, the recombined and ‘virtual’
990 recombined images would help the model to classify benign or cancer category of brain circumstance. In
991 the research of [Sun, Tseng](#) [79], the group applied CNN into a semi-supervised learning scheme for breast
992 cancer detection, which only needs very few labelled data for training. This scheme allows the usage of
993 unlabelled data if a small amount of labelled data is already learned by the system. The final result
994 achieves 88.18% with unlabelled data and 82.43% with the mixed and unlabelled data, which shows a
995 significant improvement than previous research in accuracy. In the pixel-wise region segmentation
996 scheme proposed by [Su, Liu](#) [80], fast scanning deep CNN was applied to relieve the computational
997 burden by reducing the redundant computations in the convolutional layers and max pooling layers. [Ting,](#)
998 [Tan](#) [81] introduced a model for breast cancer diagnosis and named it as Convolutional Neural Network
999 Improvement for Breast Cancer Classification (CNNI-BCC). This algorithm achieves 89.47% sensitivity

1000 and 90.50% accuracy. [Bardou, Zhang \[82\]](#) tried to figure out the breast cancer classification on the basis
 1001 of histology images utilizing CNN. The result reveals that CNN performs better than the handcrafted
 1002 feature-based classifier with 98.33% for binary classification and 82.33% for multi-class classification.
 1003

1004 **Table 1 Results of application of convolutional neural network in breast area**

Authors	Goals	Model	Results	
			Acc	Sen
D Lévy et al. [76]	Classify the mass, benign or malignant	AlexNet	0.8900	0.8680
		GoogLeNet	0.9290	0.9340
J Mordang et al. [77]	Detect microcalcifications and then classify the mass, benign or malignant	Two CNNs	/	0.7035±0.0068
F Gao et al. [78]	Classify the benign or cancer	Shallow-Deep CNN	0.90±0.06	/
W Sun et al. [79]	Allow the usage of unlabelled data in breast cancer diagnosis	A CNN model including 3 pairs of convolutional layers	0.8818	/
H Su et al. [80]	Relieve computational burden in segmentation by reducing the redundant computations	Fast scanning CNN	0.9100	/
F Ting et al. [81]	Classify incoming medical images as malignant, benign or normal patient	CNNI-BCC	0.9050	0.8947
D Bardou et al. [82]	Binary classification	CNN using an ensemble model	0.9833	/
	Multi-class classification		0.8233	

1005 (Acc= Accuracy; Sen = Sensitivity)

1006

1007 **5.2. Lung**

1008

1009 Lung cancer is a major global health problem today and for decades to come, as it is the major cause
 1010 of cancer deaths worldwide, far outnumbering breast cancer, accounting for 18.4% of all cancer deaths.
 1011 The data show that about 70% of patients are diagnosed with advanced lung cancer, so only 15% of
 1012 patients with lung cancer survive five years after being diagnosed. This means that early diagnosis and
 1013 improved diagnostic accuracy are crucial to the treatment of lung cancer and this area of technology has
 1014 long been the subject of research [\[83, 84\]](#).

1015 [Alakwaa, Nassef \[85\]](#) with his team used CT scans as input and 3D-CNN as model to detect and
 1016 classify the lung unmarked nodules. They first segmented out lung tissue by thresholding, and then, a
 1017 modified U-Net that can determine nodule candidates in CT scans is used to locate the regions that has
 1018 relatively high possibility of containing nodule candidates. Finally, these regions with nodule candidates
 1019 were sent to 3D-CNN to be identified if there exists lung cancer, since the 3D-CNN could make fully
 1020 usage of the spatial 3D context information of lung nodules. The accuracy of model achieved 86.6%. In
 1021 another similar research displayed in Table 3 in detail, the group of [Kang, Liu \[86\]](#) added the concept of
 1022 multi-view strategy into the usage of 3D-CNN to improve the model performance. They conducted a
 1023 binary classification and a three-type classification on CT images and compared results among different
 1024 methods. Compared with 2D CNN, 3D CNN is no doubt better in error rate whatever in binary
 1025 classification or ternary classification. Moreover, compared with one-view-one-network strategy, the
 1026 multi-view-one-network strategy also achieves a lower error rate. The final best results achieved a

1027 superior error rate of 4.59% for the binary classification and 7.70% for the ternary classification. One
 1028 breakthrough brought by [Shen, Zhou \[87\]](#) and his team was that, in terms of regional analysis, instead of
 1029 depending on nodule segmentation, they directly modeled raw nodule patches without any prior
 1030 definition of nodule morphology. A Multi-scale Convolutional Neural Networks (MCNN) model was set
 1031 up to capture nodule heterogeneity and the input also comes from CT which shows both lung nodule
 1032 screening and nodule annotations. The best accuracy rate achieved 86.84 % with the RF classifier as
 1033 displayed in Table 2. The first model for dealing lung pattern classification problem was presented by
 1034 [Anthimopoulos, Christodoulidis \[88\]](#) and proven to be effective when compared to previous methods on
 1035 a public dataset. In the research of [Hamidian, Sahiner \[89\]](#), they first presented a fully convolution
 1036 network (FCN) transferred from 3D CNN and used it to produce scoring results for all the volumes rather
 1037 than just volumes of interest. Then, utilizing the difficult negative examples produced by FCN to train a
 1038 new discriminant CNN model. The whole system was a combination of FCN and the new CNN. An
 1039 excellent result in lung nodule characterization was achieved by Tumor Net in the work of [Hussein,](#)
 1040 [Gillies \[90\]](#). The network they proposed had nearly 6 percent improvement in accuracy when compared
 1041 to methods without utilizing CNN on the same dataset.
 1042

1043 **Table 2 Results of application of convolutional neural network in lung area**

Authors	Goals	Model	Results	
			Acc	Sen
W Alakwaa et al. [85]	Detect the lung unmarked nodules	U-Net and 3D CNN	0.867	0.853
G Kang et al. [86]	Binary classification	MVO-3D CNN-CNN1	0.9541	/
	Ternary classification		0.923	/
W Shen et al. [87]	Model raw nodules to detect nodule heterogeneity	Multi-scale CNN	0.8684	/
A Chris et al. [88]	Lung Pattern Classification	Self-created	0.8561	/
H Sahiner et al. [89]	Automatically detect lung nodules	3D FCN+CNN	/	0.95
H Gillies et al. [90]	Lung nodule characterization	TUMORNET	0.9231	/

1044 (MVO = multi-view-one; Acc= Accuracy; Sen = Sensitivity)

1045 *(CT1= CT image channel 1; CT3= CT image channel 3)

1046

1047 **Table 3 Details in research of G Kang et al. [86]**

Input data	Goals	Model			Acc	Sen	Spc
CT 1	binary classification	3D CNN	CNN1	0.9518	0.9556	0.9301	
			CNN2	0.9503	0.9551	0.9299	
			CNN3	0.9481	0.9845	0.8976	
		2D CNN	CNN1	0.9454	0.8888	0.9967	
			CNN2	0.944	0.8893	0.9972	
			CNN3	0.9413	0.894	0.9874	
		3D CNN	CNN1	0.9525	0.956	0.9394	
			CNN2	0.9524	0.9561	0.9389	
CT 3							

				CNN3	0.9362	0.9417	0.8973	
CT	ternary classification		2D CNN	CNN1	0.9459	0.8892	1	
				CNN2	0.9446	0.89	0.9976	
				CNN3	0.944	0.8897	0.9968	
				MVO	CNN	0.9541	/	
CT 1			3D CNN	CNN1	0.9472	/	/	
				CNN2	0.9439	/	/	
				3D CNN	CNN1	0.9188	/	
CT 3			MVO	CNN2	0.9158	/	/	
				CNN3	0.8906	/	/	
				2D CNN	CNN1	0.8234	/	
CT			3D CNN	CNN2	0.8373	/	/	
				CNN3	0.8479	/	/	
				2D CNN	CNN1	0.9222	/	
			MVO	CNN2	0.9175	/	/	
				CNN3	0.9069	/	/	
				2D CNN	CNN1	0.8184	/	
			OVO	CNN2	0.8195	/	/	
				CNN3	0.8681	/	/	
				3D CNN	CNN	0.923	/	
			MVO	CNN1	0.9187	/	/	
				CNN2	0.9093	/	/	

1048 (MVO = multi-view-one; OVO = one-view-one; Acc= Accuracy; Sen = Sensitivity; Spc = Specificity)

1049 *(CT1= CT image channel 1; CT3= CT image channel 3)

1050

1051 5.3. Brain

1052

1053 In the application field of brain, CNN has also made outstanding contributions as illustrated in Table
1054 4. It is well known that there are huge differences in appearance between brain tumors, which
1055 undoubtedly poses difficult challenges for pathologists or doctors to accurately segment and identify
1056 brain tumors. Thus, making use of deep learning to segment and identify brain tumors becomes a
1057 promising way of thinking in human medicine [91-100].

1058 [Khan, Shah \[101\]](#) and his team proposed a fully automated CNN-based method for segmenting brain
1059 tumors that required a combination of handcrafted features, and achieved excellent results in terms of
1060 complete, core, and enhancing tumor. More recently, a fine-tuned VGG-19 (FT-VGG-19) approach to
1061 classifying multistage brain tumors that can help radiologists better analyze MRI has yielded compelling
1062 results, according to [Sajjad, Khan \[102\]](#) and his team. This classification system not only made use of
1063 deep learning to segment the tumor region in the image, but also made usage of data expansion to solve

1064 the problem of insufficient data, which showed improved performance in comparison to other existing
 1065 methods. In addition, to help doctors detect gliomas at an early stage, [Anaraki, Ayati \[103\]](#) and his team
 1066 came up with a glioma classifier that combines genetic algorithms with a CNN model (CNN-GA) and
 1067 gets it right more than 90 percent of the time. This classifier is very flexible because it can be effectively
 1068 classified in multiple case studies on the classification of brain tumors. [Moeskops, Viergever \[104\]](#)
 1069 introduced an automatic segmentation method based on CNN, which aims to segment MR brain images
 1070 into lots of tissue classes. This method tried to learn to recognize the significant information for
 1071 classification on the basis of training data rather than just dependent on explicit features. The team applied
 1072 the methods to 5 datasets and achieved good Dice similarity score (Dss) over segmented tissue classes
 1073 with the highest result of 0.91, which proves that this automatic segmentation method using CNN has
 1074 robustness to datasets of different age and acquisition protocol. A brain extraction tool based on the auto-
 1075 context CNN invented by [Salehi, Erdogmus \[105\]](#) outperforms several results recently reported in brain
 1076 analysis research. It was superior than other brain extraction methods not only in high Dss results, but
 1077 also due to its characteristic of registration-free, which means this CNN technique relied less on accuracy
 1078 of alignment between atlases and anatomy in brain and did not require making assumptions in geometry
 1079 of brain image. The team of [Kamnitsas, Chen \[106\]](#) introduced a multi-scale CNN solution to segment
 1080 brain lesion and achieved better performance than previous method on a small-scale dataset. [Hussain,](#)
 1081 [Anwar \[107\]](#) presented several architectures based on CNN to segment Glioma tumors in brain and
 1082 reaches higher results when compared to state-of-art methods. [Mehta and Sivaswamy \[108\]](#) introduced
 1083 M-net to help segmenting human brain structure. M-net attracted attention in medical research field as it
 1084 ran 3 times faster in speed when compared to traditional methods.

1085

1086 **Table 4 Results of application of convolutional neural network in brain area**

Authors	Goals	Model	Results			
			Acc	Dss		
				Com	Co	E
K Shah et al. [101]	Fully automatically segment the brain tumor	TP-CNN**	/	0.81	0.76	0.73
S Khan et al. [102]	Classify multi-grade brain tumor	FT-VGG-19	0.9458	/		
A Ayati et al. [103]	Classify different grades of glioma	CNN-GA	0.909	/		
			0.942	/		
P Moeskops et al. [104]	Segment MR brain images into tissue classes	Self-created	/	0.91		
S Erdogmus et al. [105]	Brain extraction	Auto-context CNN	/	0.9773		
K Chen et al. [106]	Segment lesion	Multi-scale CNN	/	0.66		
H Anwar et al. [107]	Segment Glioma tumors in brain	Self-created	/	0.86	0.87	0.90
M and S[108]	Segment deep brain structure	M-net	/	0.8578		

1087 (Dss= Dice Similarity Score; Acc=Accuracy; Com=Complete; Co=Core; E=Enhancing)

1088 * (TP-CNN is a combination of CSGP+MRILP+MRIGP)

1089

1090 Overall, we can see that CNN does bring excellent performance when compared to other state-of-

1091 art methods in various areas of medical image analysis. It is often possible to quickly achieve accuracy
1092 of up to 90% or more than 95% from data sets that human pathologists may expend considerable effort
1093 and time on. There is no doubt that CNN will be of great help to the diagnosis in the medical image
1094 analysis area.

1095

6. CONCLUSION

1097

1098 This review presents a comprehensive view on Convolutional Neural Network in the medical field.
1099 The content including its developing history, basic operating mechanism, basic structure, classic
1100 architecture and applications in medical image analysis.

1101 We first outlined the basic operating mechanism of convolutional neural network, with the survey
1102 in neuron, activation function, loss function, gradient descent, feedforward and feedback propagation.
1103 And then detailedly introduced 4 categories of convolutional layer, 3 categories of pooling layer, and
1104 fully connected layer. The major benefits and limitations of several classic architecture: AlexNet, VGG-
1105 Net, Inception, Resnet, U-Net, R-CNN Series, YOLO and SSD are concluded. The applications of CNN
1106 among breast, lung and brain area are finally considered from summarizing a number of relevant core
1107 journals. We hope that by reading our review, readers can understand the basic principles and popular
1108 use framework of CNN, be able to utilize popular format of convolutional or pooling layer like transposed
1109 convolutional layer, separable convolutional layer, stochastic pooling in appropriate circumstances, and
1110 feel the rapid development of CNN and outstanding contributions brought by CNN in the field of
1111 medicine.

1112 A growing number of cases prove that CNN can explore and discover the mysteries of some
1113 unknown medical fields. It is expected that the CNN could effectively has the ability to aid diagnostic
1114 applications for evaluating more and more diseases in humans in the future.

1115

Acknowledgement

1117

1118 This work was partially supported by Royal Society International Exchanges Cost Share Award, UK
1119 (RP202G0230); Medical Research Council Confidence in Concept Award, UK (MC_PC_17171); Hope
1120 Foundation for Cancer Research, UK (RM60G0680); British Heart Foundation Accelerator Award, UK;
1121 Guangxi Key Laboratory of Trusted Software (kx201901); Fundamental Research Funds for the Central
1122 Universities (CDLS-2020-03); Key Laboratory of Child Development and Learning Science (Southeast
1123 University), Ministry of Education.

1124

Reference

1126

- 1127 1. LeCun, Y., Y. Bengio, and G. Hinton, *Deep learning*. nature, 2015. **521**(7553): p. 436-444.
- 1128 2. Schmidhuber, J., *Deep learning in neural networks: An overview*. Neural networks, 2015. **61**: p. 85-117.
- 1129 3. Reddy, S., *Use of artificial intelligence in healthcare delivery*, in *eHealth-Making Health Care Smarter*. 2018, IntechOpen.
1130 p. p. 81-97.
- 1131 4. McCulloch, W.S. and W. Pitts, *A logical calculus of the ideas immanent in nervous activity*. The bulletin of mathematical
1132 biophysics, 1943. **5**(4): p. 115-133.
- 1133 5. Rosenblatt, F., *The perceptron: a probabilistic model for information storage and organization in the brain*. Psychological
1134 review, 1958. **65**(6): p. 386.

- 1135 6. Minsky, M. and S.A. Papert, *Perceptrons: An introduction to computational geometry*. 2017: MIT press.
- 1136 7. Zhang, Z. and M. Sarhadi. *A modified neuron activation function which enables single layer perceptrons to solve some linearly*
 1137 *inseparable problems*. in *Proceedings of 1993 International Conference on Neural Networks (IJCNN-93-Nagoya, Japan)*.
 1138 1993. IEEE.
- 1139 8. Rumelhart, D.E., G.E. Hinton, and R.J. Williams, *Learning representations by back-propagating errors*. nature, 1986.
 1140 323(6088): p. 533-536.
- 1141 9. Lawrence, S., et al., *Face recognition: A convolutional neural-network approach*. IEEE transactions on neural networks, 1997.
 1142 8(1): p. 98-113.
- 1143 10. Mikolov, T., et al. *Recurrent neural network based language model*. in *Eleventh annual conference of the international speech*
 1144 *communication association*. 2010.
- 1145 11. Schuster, M. and K.K. Paliwal, *Bidirectional recurrent neural networks*. IEEE transactions on Signal Processing, 1997. 45(11):
 1146 p. 2673-2681.
- 1147 12. Hochreiter, S. and J. Schmidhuber, *Long short-term memory*. Neural computation, 1997. 9(8): p. 1735-1780.
- 1148 13. Baudat, G. and F. Anouar, *Generalized discriminant analysis using a kernel approach*. Neural computation, 2000. 12(10): p.
 1149 2385-2404.
- 1150 14. Agarap, A.F., *Deep learning using rectified linear units (relu)*. arXiv preprint arXiv:1803.08375, 2018.
- 1151 15. Gauen, K., et al. *Low-power image recognition challenge*. in *2017 22nd Asia and South Pacific Design Automation*
 1152 *Conference (ASP-DAC)*. 2017. IEEE.
- 1153 16. Krizhevsky, A., I. Sutskever, and G.E. Hinton. *ImageNet classification with deep convolutional neural networks*. in *Advances*
 1154 *in neural information processing systems*. 2012.
- 1155 17. Szegedy, C., et al. *Inception-v4, inception-resnet and the impact of residual connections on learning*. in *Thirty-first AAAI*
 1156 *conference on artificial intelligence*. 2017.
- 1157 18. Howard, A.G., *Some improvements on deep convolutional neural network based image classification*. arXiv preprint
 1158 arXiv:1312.5402, 2013.
- 1159 19. Li, J., et al. *A comparison of deep learning methods for environmental sound detection*. in *2017 IEEE International Conference*
 1160 *on Acoustics, Speech and Signal Processing (ICASSP)*. 2017. IEEE.
- 1161 20. Marsalli, M. *Mcculloch-pitts neurons*. in *The 2008 Annual Meeting of the consortium on cognitive science instruction (ccsi)*.
 1162 2006.
- 1163 21. Goodfellow, I., Y. Bengio, and A. Courville, *Deep Learning [Sl]*. 2016, MIT press: 1 Rogers Street in Cambridge, MA 02142.
- 1164 22. Nielsen, M.A., *Neural networks and deep learning*. Vol. 2018. 2015: Determination press San Francisco, CA, USA:.
- 1165 23. Scardapane, S., et al., *Kafnets: Kernel-based non-parametric activation functions for neural networks*. Neural Networks, 2019.
 1166 110: p. 19-32.
- 1167 24. Chung, H., S.J. Lee, and J.G. Park. *Deep neural network using trainable activation functions*. in *2016 International Joint*
 1168 *Conference on Neural Networks (IJCNN)*. 2016. IEEE.
- 1169 25. Baldassi, C., E.M. Malatesta, and R. Zecchina, *Properties of the Geometry of Solutions and Capacity of Multilayer Neural*
 1170 *Networks with Rectified Linear Unit Activations*. Physical Review Letters, 2019. 123(17): p. 6.
- 1171 26. Pedamonti, D., *Comparison of non-linear activation functions for deep neural networks on MNIST classification task*. arXiv
 1172 preprint arXiv:1804.02763, 2018.
- 1173 27. Lemaréchal, C., *Cauchy and the gradient method*. Doc Math Extra, 2012. 251: p. 254.
- 1174 28. Ruder, S., *An overview of gradient descent optimization algorithms*. arXiv preprint arXiv:1609.04747, 2016.
- 1175 29. Bottou, L., *Large-scale machine learning with stochastic gradient descent*, in *Proceedings of COMPSTAT'2010*. 2010,
 1176 Springer. p. 177-186.
- 1177 30. Bottou, L., *Stochastic gradient descent tricks*, in *Neural networks: Tricks of the trade*. 2012, Springer. p. 421-436.
- 1178 31. Khirirat, S., H.R. Feyzmahdavian, and M. Johansson. *Mini-batch gradient descent: Faster convergence under data sparsity*.

- 1179 in *2017 IEEE 56th Annual Conference on Decision and Control (CDC)*. 2017. IEEE.
- 1180 32. Hinton, G., N. Srivastava, and K. Swersky, *Overview of mini-batch gradient descent*. Neural Networks for Machine Learning,
- 1181 2012. **575**.
- 1182 33. Bengio, Y., P. Simard, and P. Frasconi, *Learning long-term dependencies with gradient descent is difficult*. IEEE transactions
- 1183 on neural networks, 1994. **5**(2): p. 157-166.
- 1184 34. Pascanu, R., T. Mikolov, and Y. Bengio, *Understanding the exploding gradient problem*. CoRR, abs/1211.5063, 2012. **2**: p.
- 1185 417.
- 1186 35. LeCun, Y., et al., *Backpropagation applied to handwritten zip code recognition*. Neural computation, 1989. **1**(4): p. 541-551.
- 1187 36. Che, Z.-G., T.-A. Chiang, and Z.-H. Che, *Feed-forward neural networks training: a comparison between genetic algorithm*
- 1188 *and back-propagation learning algorithm*. International journal of innovative computing, information and control, 2011. **7**(10):
- 1189 p. 5839-5850.
- 1190 37. Schmidt, W.F., et al. *Initializations, back-propagation and generalization of feed-forward classifiers*. in *IEEE International*
- 1191 *Conference on Neural Networks*. 1993. IEEE.
- 1192 38. Rastegari, M., et al. *Xnor-net: Imagenet classification using binary convolutional neural networks*. in *European conference*
- 1193 *on computer vision*. 2016. Springer.
- 1194 39. Rouhi, R., et al., *Benign and malignant breast tumors classification based on region growing and CNN segmentation*. Expert
- 1195 Systems with Applications, 2015. **42**(3): p. 990-1002.
- 1196 40. Dasgupta, A. and S. Singh. *A fully convolutional neural network based structured prediction approach towards the retinal*
- 1197 *vessel segmentation*. in *2017 IEEE 14th International Symposium on Biomedical Imaging (ISBI 2017)*. 2017. IEEE.
- 1198 41. Juefei-Xu, F., V. Naresh Bodetti, and M. Savvides. *Local binary convolutional neural networks*. in *Proceedings of the IEEE*
- 1199 *conference on computer vision and pattern recognition*. 2017.
- 1200 42. Li, Y., X. Zhang, and D. Chen. *Csrnet: Dilated convolutional neural networks for understanding the highly congested scenes*.
- 1201 in *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018.
- 1202 43. Yu, F. and V. Koltun, *Multi-scale context aggregation by dilated convolutions*. arXiv preprint arXiv:1511.07122, 2015.
- 1203 44. Zeiler, M.D., G.W. Taylor, and R. Fergus. *Adaptive deconvolutional networks for mid and high level feature learning*. in *2011*
- 1204 *International Conference on Computer Vision*. 2011. IEEE.
- 1205 45. Zeiler, M.D. and R. Fergus. *Visualizing and understanding convolutional networks*. in *European conference on computer*
- 1206 *vision*. 2014. Springer.
- 1207 46. Dumoulin, V. and F. Visin, *A guide to convolution arithmetic for deep learning*. arXiv preprint arXiv:1603.07285, 2016.
- 1208 47. Chollet, F. *Xception: Deep learning with depthwise separable convolutions*. in *Proceedings of the IEEE conference on*
- 1209 *computer vision and pattern recognition*. 2017.
- 1210 48. Chen, L.-C., et al. *Encoder-decoder with atrous separable convolution for semantic image segmentation*. in *Proceedings of*
- 1211 *the European conference on computer vision (ECCV)*. 2018.
- 1212 49. Sheng, T., et al. *A quantization-friendly separable convolution for mobilenets*. in *2018 1st Workshop on Energy Efficient*
- 1213 *Machine Learning and Cognitive Computing for Embedded Applications (EMC2)*. 2018. IEEE.
- 1214 50. Sun, M., et al., *Learning pooling for convolutional neural network*. Neurocomputing, 2017. **224**: p. 96-104.
- 1215 51. Zubair, S., F. Yan, and W. Wang, *Dictionary learning based sparse coefficients for audio classification with max and average*
- 1216 *pooling*. Digital Signal Processing, 2013. **23**(3): p. 960-970.
- 1217 52. Zeiler, M.D. and R. Fergus, *Stochastic pooling for regularization of deep convolutional neural networks*. arXiv preprint
- 1218 arXiv:1301.3557, 2013.
- 1219 53. Yu, D., et al. *Mixed pooling for convolutional neural networks*. in *International conference on rough sets and knowledge*
- 1220 *technology*. 2014. Springer.
- 1221 54. Huang, J.-T., J. Li, and Y. Gong. *An analysis of convolutional neural networks for speech recognition*. in *2015 IEEE*
- 1222 *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 2015. IEEE.

- 1223 55. Zar\'andy, \AA., et al. *Overview of CNN research: 25 years history and the current trends*. in *2015 IEEE International Symposium*
1224 *on Circuits and Systems (ISCAS)*. 2015. IEEE.
- 1225 56. Iandola, F.N., et al., *SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and < 0.5 MB model size*. arXiv preprint
1226 arXiv:1602.07360, 2016.
- 1227 57. Yuan, Z.-W. and J. Zhang. *Feature extraction and image retrieval based on AlexNet*. in *Eighth International Conference on*
1228 *Digital Image Processing (ICDIP 2016)*. 2016. International Society for Optics and Photonics.
- 1229 58. Simonyan, K. and A. Zisserman, *Very deep convolutional networks for large-scale image recognition*. arXiv preprint
1230 arXiv:1409.1556, 2014.
- 1231 59. Mateen, M., et al., *Fundus image classification using VGG-19 architecture with PCA and SVD*. *Symmetry*, 2019. **11**(1): p. 1.
- 1232 60. Szegedy, C., et al. *Rethinking the inception architecture for computer vision*. in *Proceedings of the IEEE conference on*
1233 *computer vision and pattern recognition*. 2016.
- 1234 61. Szegedy, C., et al. *Going deeper with convolutions*. in *Proceedings of the IEEE conference on computer vision and pattern*
1235 *recognition*. 2015.
- 1236 62. Lin, M., Q. Chen, and S. Yan, *Network in network*. arXiv preprint arXiv:1312.4400, 2013.
- 1237 63. Salekin, M.S., A.B. Jelodar, and R. Kushol. *Cooking state recognition from images using inception architecture*. in *2019*
1238 *International Conference on Robotics, Electrical and Signal Processing Techniques (ICREST)*. 2019. IEEE.
- 1239 64. He, K., et al. *Identity mappings in deep residual networks*. in *European conference on computer vision*. 2016. Springer.
- 1240 65. Ronneberger, O., P. Fischer, and T. Brox. *U-net: Convolutional networks for biomedical image segmentation*. in *International*
1241 *Conference on Medical image computing and computer-assisted intervention*. 2015. Springer.
- 1242 66. Gkioxari, G., R. Girshick, and J. Malik. *Contextual action recognition with r* cnn*. in *Proceedings of the IEEE international*
1243 *conference on computer vision*. 2015.
- 1244 67. Lenc, K. and A. Vedaldi, *R-cnn minus r*. arXiv preprint arXiv:1506.06981, 2015.
- 1245 68. Purkait, P., C. Zhao, and C. Zach, *SPP-Net: Deep absolute pose regression with synthetic views*. arXiv preprint
1246 arXiv:1712.03452, 2017.
- 1247 69. Girshick, R. *Fast r-cnn*. in *Proceedings of the IEEE international conference on computer vision*. 2015.
- 1248 70. Ren, S., et al. *Faster r-cnn: Towards real-time object detection with region proposal networks*. in *Advances in neural*
1249 *information processing systems*. 2015.
- 1250 71. Du, J. *Understanding of object detection based on CNN family and YOLO*. in *Journal of Physics: Conference Series*. 2018.
1251 IOP Publishing.
- 1252 72. Liu, W., et al. *Ssd: Single shot multibox detector*. in *European conference on computer vision*. 2016. Springer.
- 1253 73. Schairer, C., et al., *Probabilities of death from breast cancer and other causes among female breast cancer patients*. *Journal*
1254 *of the National Cancer Institute*, 2004. **96**(17): p. 1311-1321.
- 1255 74. Azamjah, N., Y. Soltan-Zadeh, and F. Zayeri, *Global Trend of Breast Cancer Mortality Rate: A 25-Year Study*. *Asian Pacific*
1256 *journal of cancer prevention: APJCP*, 2019. **20**(7): p. 2015.
- 1257 75. Waks, A.G. and E.P. Winer, *Breast cancer treatment: a review*. *Jama*, 2019. **321**(3): p. 288-300.
- 1258 76. L\'evy, D. and A. Jain, *Breast mass classification from mammograms using deep convolutional neural networks*. arXiv preprint
1259 arXiv:1612.00542, 2016.
- 1260 77. Mordang, J.-J., et al. *Automatic microcalcification detection in multi-vendor mammography using convolutional neural*
1261 *networks*. in *International Workshop on Breast Imaging*. 2016. Springer.
- 1262 78. Gao, F., et al., *SD-CNN: A shallow-deep CNN for improved breast cancer diagnosis*. *Computerized Medical Imaging and*
1263 *Graphics*, 2018. **70**: p. 53-62.
- 1264 79. Sun, W., et al., *Enhancing deep convolutional neural network scheme for breast cancer diagnosis with unlabeled data*.
1265 *Computerized Medical Imaging and Graphics*, 2017. **57**: p. 4-9.
- 1266 80. Su, H., et al. *Region segmentation in histopathological breast cancer images using deep convolutional neural network*. in

- 1267 2015 IEEE 12th International Symposium on Biomedical Imaging (ISBI). 2015. IEEE.
- 1268 81. Ting, F.F., Y.J. Tan, and K.S. Sim, *Convolutional neural network improvement for breast cancer classification*. Expert Systems with Applications, 2019. **120**: p. 103-115.
- 1269
- 1270 82. Bardou, D., K. Zhang, and S.M. Ahmad, *Classification of breast cancer based on histology images using convolutional neural networks*. IEEE Access, 2018. **6**: p. 24680-24693.
- 1271
- 1272 83. Bray, F., et al., *Erratum: Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries*.
- 1273
- 1274 84. de Koning, H.J., et al., *Reduced lung-cancer mortality with volume CT screening in a randomized trial*. New England Journal of Medicine, 2020. **382**(6): p. 503-513.
- 1275
- 1276 85. Alakwaa, W., M. Nassef, and A. Badr, *Lung cancer detection and classification with 3D convolutional neural network (3D-CNN)*. Lung Cancer, 2017. **8**(8): p. 409.
- 1277
- 1278 86. Kang, G., et al., *3D multi-view convolutional neural networks for lung nodule classification*. PloS one, 2017. **12**(11): p. 1-21.
- 1279
- 1280 87. Shen, W., et al. *Multi-scale convolutional neural networks for lung nodule classification*. in *International Conference on Information Processing in Medical Imaging*. 2015. Springer.
- 1281
- 1282 88. Anthimopoulos, M., et al., *Lung pattern classification for interstitial lung diseases using a deep convolutional neural network*. IEEE transactions on medical imaging, 2016. **35**(5): p. 1207-1216.
- 1283
- 1284 89. Hamidian, S., et al. *3D convolutional neural network for automatic detection of lung nodules in chest CT*. in *Medical Imaging 2017: Computer-Aided Diagnosis*. 2017. International Society for Optics and Photonics.
- 1285
- 1286 90. Hussein, S., et al. *Tumornet: Lung nodule characterization using multi-view convolutional neural network with gaussian process*. in *2017 IEEE 14th International Symposium on Biomedical Imaging (ISBI 2017)*. 2017. IEEE.
- 1287
- 1288 91. Toğuçar, M., B. Ergen, and Z. Cömert, *BrainMRNet: Brain tumor detection using magnetic resonance images with a novel convolutional neural network model*. Medical Hypotheses, 2020. **134**: p. 109531.
- 1289
- 1290 92. Xue, Y., et al., *A multi-path 2.5 dimensional convolutional neural network system for segmenting stroke lesions in brain MRI images*. NeuroImage: Clinical, 2020. **25**: p. 102118.
- 1291
- 1292 93. Jiang, Y., et al., *A novel distributed multitask fuzzy clustering algorithm for automatic MR brain image segmentation*. Journal of medical systems, 2019. **43**(5): p. 118.
- 1293
- 1294 94. Jiang, Y., et al. *A Novel Negative-Transfer-Resistant Fuzzy Clustering Model with a Shared Cross-Domain Transfer Latent Space and its Application to Brain CT Image Segmentation*. IEEE/ACM Transactions on Computational Biology and Bioinformatics, 2020. DOI: 10.1109/TCBB.2019.2963873.
- 1295
- 1296 95. Jiang, Y., et al. *EEG-Based Driver Drowsiness Estimation Using an Online Multi-View and Transfer TSK Fuzzy System*. IEEE Transactions on Intelligent Transportation Systems, 2020. 1-13 DOI: 10.1109/TITS.2020.2973673.
- 1297
- 1298 96. Jiang, Y.Z., et al., *Seizure Classification From EEG Signals Using Transfer Learning, Semi-Supervised Learning and TSK Fuzzy System*. IEEE Transactions on Neural Systems and Rehabilitation Engineering, 2017. **25**(12): p. 2270-2284.
- 1299
- 1300 97. Jiang, Y.Z., et al., *Recognition of Epileptic EEG Signals Using a Novel Multiview TSK Fuzzy System*. IEEE Transactions on Fuzzy Systems, 2017. **25**(1): p. 3-20.
- 1301
- 1302 98. Jiang, Y.Z., et al., *Collaborative Fuzzy Clustering From Multiple Weighted Views*. IEEE Transactions on Cybernetics, 2015. **45**(4): p. 688-701.
- 1303
- 1304 99. Jiang, Y.Z., et al., *Seizure Recognition Using a Novel Multitask Radial Basis Function Neural Network*. Journal of Medical Imaging and Health Informatics, 2019. **9**(9): p. 1865-1870.
- 1305
- 1306 100. Jiang, Y.Z., et al., *A Novel Synthetic CT Generation Method Using Multitask Maximum Entropy Clustering*. Ieee Access, 2019. **7**: p. 119644-119653.
- 1307
- 1308 101. Khan, H., et al., *Cascading handcrafted features and Convolutional Neural Network for IoT-enabled brain tumor segmentation*. Computer Communications, 2020.
- 1309
- 1310

- 1311 102. Sajjad, M., et al., *Multi-grade brain tumor classification using deep CNN with extensive data augmentation*. Journal of
1312 computational science, 2019. **30**: p. 174-182.
- 1313 103. Anaraki, A.K., M. Ayati, and F. Kazemi, *Magnetic resonance imaging-based brain tumor grades classification and grading*
1314 *via convolutional neural networks and genetic algorithms*. Biocybernetics and Biomedical Engineering, 2019. **39**(1): p. 63-
1315 74.
- 1316 104. Moeskops, P., et al., *Automatic segmentation of MR brain images with a convolutional neural network*. IEEE transactions on
1317 medical imaging, 2016. **35**(5): p. 1252-1261.
- 1318 105. Salehi, S.S.M., D. Erdogmus, and A. Gholipour, *Auto-context convolutional neural network (auto-net) for brain extraction in*
1319 *magnetic resonance imaging*. IEEE transactions on medical imaging, 2017. **36**(11): p. 2319-2330.
- 1320 106. Kamnitsas, K., et al., *Multi-scale 3D convolutional neural networks for lesion segmentation in brain MRI*. Ischemic stroke
1321 lesion segmentation, 2015. **13**: p. 46.
- 1322 107. Hussain, S., S.M. Anwar, and M. Majid, *Segmentation of glioma tumors in brain using deep convolutional neural network*.
1323 Neurocomputing, 2018. **282**: p. 248-261.
- 1324 108. Mehta, R. and J. Sivaswamy. *M-net: A convolutional neural network for deep brain structure segmentation*. in *2017 IEEE*
1325 *14th International Symposium on Biomedical Imaging (ISBI 2017)*. 2017. IEEE.
- 1326