# Image Captcha Recognition

**Xiaoguang Wang**
A53246911
xiw503@ ucsd.edu

**Qian Wang**
A53235276
qiw018@ ucsd.edu

**Hao Zhang**
A53234732
haz351@ ucsd.edu

## Abstract

Image captcha has long been used as a common protective tool for some websites to tell computer and human apart. Here, we employed the pre-trained AlexNet and combined with joint training method to build a model that can automatically recognize the numbers and characters on the captcha image and received a recognition accuracy of 96% on the most complex occasion with numbers, upper and lower characters. Also, we tried to borrow some technology from traditional computer vision like median blur to denoise the captcha images and improved the performance of our recognition pipeline.

## 1   Introduction

Captcha here is an acronym for "**C**ompletely **A**utomatic **P**ublic **T**uring test to tell **C**omputers and **H**umans **A**part". It is useful since computers can hardly read the contents like numbers and characters from images like a human. However, thanks to the fast development in the research of computer vision and especially artificial neural network, machines now can better understand the specific semantic information from images. For example, LeNet[1] has proved great success in handwritten digit recognition.



Figure 1: Some examples for captcha images

Inspired from the above ideas and forerunners, we proposed a novel system which can automatically learn useful features from the captcha images and make predictions of the numbers and characters on each digit separately. Figure 2 provides the broad framework for our system.

In this project, we proposed the novel captcha recognition network above and implemented the system on TensorFlow. The main contributions of our work are outlined as follows:

● Modified the existing AlexNet with multiple output layers to make predictions for each digit of captcha code separately and used multi-task learning method to tune the model.

● Combined with some image smoothing methods from traditional computer vision to pre-process the input images to the network and improved performance.
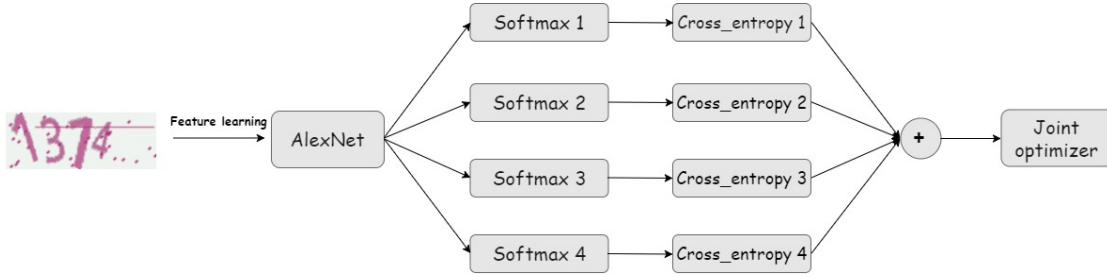
Figure 2: Network structure for the recognition system (4-digits captcha)

## 2    Approach

In this section, we present some specific description pf the designed network and some implementation details. Specifically, although the method proposed in this project can be easily adapted to any length captcha recognition, we implemented a system for 4-digit here for the convenience to verify the feasibility of our ideas.

### 2.1    Model architecture overview

Figure 2 contains the architecture of the model we built in this project. It is a transformation from traditional AlexNet.

**Feature learning**: We firstly used a pre-trained AlexNet for image feature learning since it has been well trained on a large image dataset and proved to represent images features well with lower dimension feature vectors. We leave the net parameters untouched during the training process. Also, we explored some more modern networks like GoogLeNet, but unluckily, the GPU server seems out of resources to run this large network well. So, we chose AlexNet at last and believe this network can extract necessary features well from the simple captcha images.
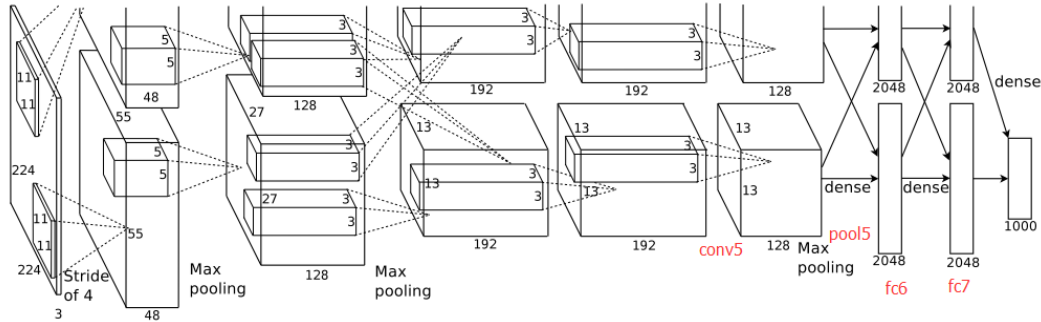


Figure 3: Original AlexNet structure for comparison

**Prediction**: we used one *softmax* layer to predict each digit separately (here we set the captcha length to be 4, so we implemented 4 *softmax* layers). To train this model, we use a method called multi-task learning (we will give specific explanation in the next part).

### 2.2 Multi-task learning

Normally in multiple categories classification problem, we use *softmax* layer to map image features from their feature space to a probability distribution. Then backpropagate the calculated the cross-entropy loss value to tune net parameters. However, here we have an output and get a cross-entropy loss value for each output *softmax* layer.

66  To handle this problem, we used a method from multi-tasking learning called **joint training**.
67  That is in the training process, instead of backpropagating these loss values four times
68  separately, we backpropagate the sum or the mean value of these loss values and assume that
69  we can get the best prediction if we can minimize this 'complex loss value' through some
70  methods like gradient descent.

71

## 72  2.3 Image smoothing

73  Smoothing is a widely used effect in graphics software, typically to reduce image noise and
74  reduce detail. The visual effect of this blurring technique is a smooth blur resembling that of
75  viewing the image through a translucent screen, distinctly different from the bokeh effect
76  produced by an out-of-focus lens or the shadow of an object under usual illumination. Here,
77  we mainly tried Gaussian blur and median blur.

78  **Gaussian blur**: Gaussian blur (also known as Gaussian smoothing) is the result of blurring
79  an image by a Gaussian Function. Mathematically, applying a Gaussian blur to an image is
80  the same as convolving the image with a Gaussian Function. Since the Fourier transform of a
81  Gaussian is another Gaussian, applying a Gaussian blur has the effect of reducing the image's
82  high-frequency components; a Gaussian blur is thus a low pass filter.The equation of a
83  Gaussian function in one dimension is

$$G(r) = \frac{1}{\sqrt{2\pi\sigma^2}^N} e^{-r^2/(2\sigma^2)}$$

84

85  **Median Blur**: The median filter is a nonlinear digital filtering technique, often used to
86  remove noise from an image or signal. Such noise reduction is a typical pre-processing step
87  to improve the results of later processing. Median filtering is very widely used in
88  digital image processing because, under certain conditions, it preserves edges while
89  removing noise, also having applications in signal processing.

90  Below is an example of applying the above two kinds of smoothing method.



91  Figure 4: effect of image smoothing (original, Gaussian blur and median blur from left to
92  right)

93  We can see, median blur has a better performance on removing noise, so we chose it for
94  further exploration.

95

96

97

98

99

100

101

102

103

104

105
106

## 3      Experiment

We built a 4-digit captcha image recognition model to verify the feasibility of our design of the system. Also, the model is implemented with TensorFlow and ran on UCSD DSMLP.

### 3.1      Dataset

We searched the opening datasets on the Internet, but we cannot find a suitable dataset (one at least should have sufficient captcha images and corresponding labels to train our network). But luckily enough, there is a Python library called *captcha* (https://pypi.org/project/captcha/), we can use the library to build our own dataset (including training set and test set). Below is a demo of the generated dataset.



Figure 5: Running demo for Python captcha

We can see for each captcha image, we set their file name to be the content on the image, which serve as the supervised information(labels) during the training process. Also, since the dataset is generated and can be controlled strictly on the image format and dataset size, which is very convenient to verify our ideas and model. In this project, we generated 10000 images for each model and used 500 for testing and the rest for training.

### 3.2      Evaluation Protocol

For image captcha recognition, what we concern about mostly is the prediction accuracy. So here, we use prediction accuracy on each of the 4 digits as the only metrics to evaluate out models.

More importantly, let's consider another important problem for our model – robustness from the following two aspects:

● **Complexity of the code**: here, we explored different levels of complexity of the generated captcha code from numbers only (char set size: 10) to numbers + lower case characters (char set size: 36) and lastly the most complex case with numbers + lower and upper case characters (char set size: 62).
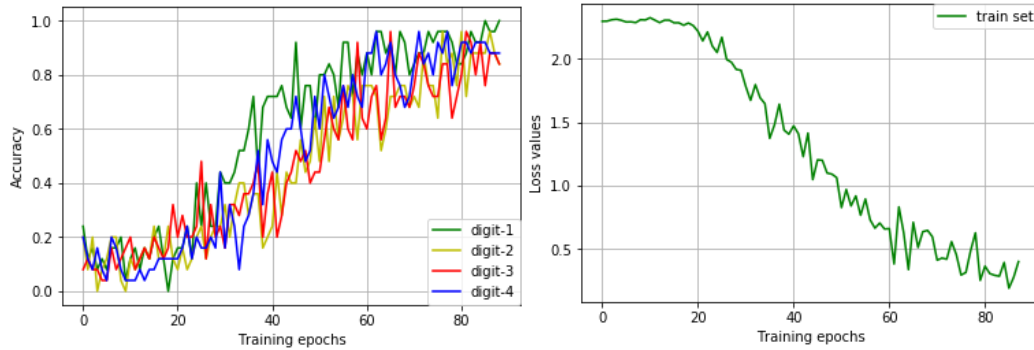
● **Noise**: As an improvement to our system, here instead of inputting the generated captcha images directly into the model, we borrowed some technology from traditional computer vision and pre-process the images with some noise cancellation operations like Gaussian blue and median blur. We compared the effect of image smoothing on our model.

## 3.3    Results

In this part, we will present results we get on our model with different char set complexity: numbers only, numbers + lower case characters, number + lower + upper case characters and a comparison between the results of with or without smoothing operation. Typically, recognition accuracy on each digit can reach a level higher than 90%.
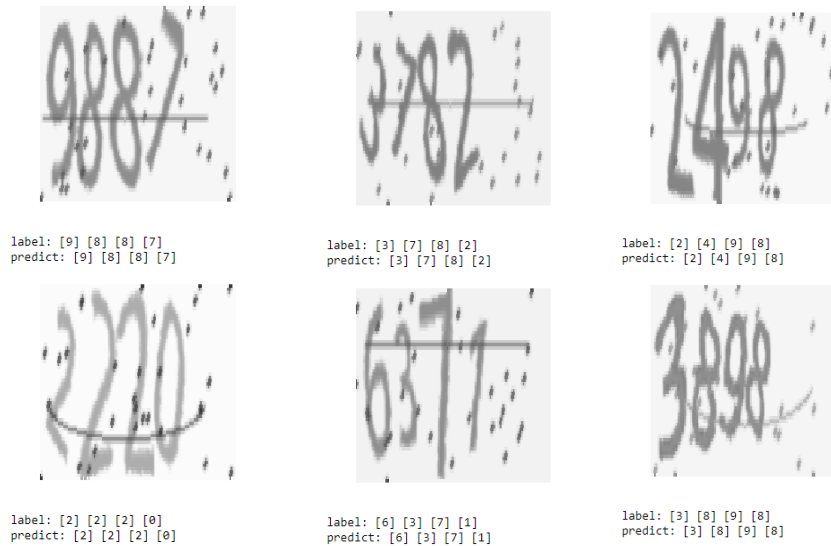
### 3.3.1 numbers only

We only used digits from 0 to 9 to compose a 4-digits captcha image. Below are the accuracies on test set and mean loss value on training set we get during the training process.



Figure 6: Accuracies and loss value during the training process (numbers only)
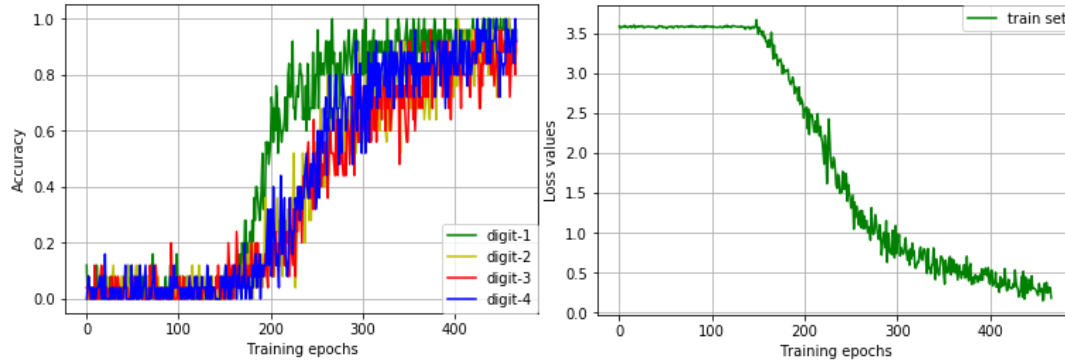
(sampling rate： 20 iterations/sample)



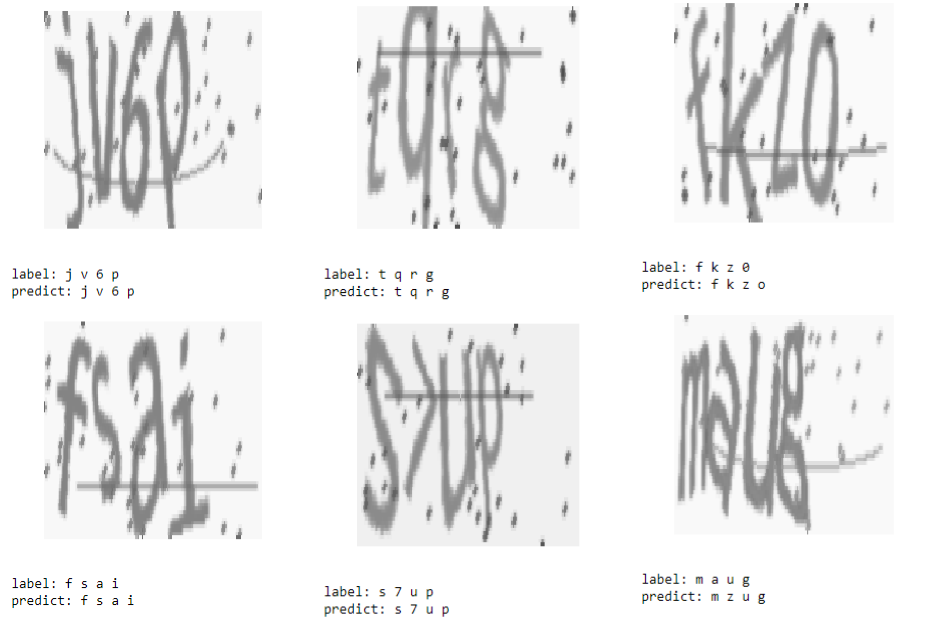Figure 7: Test demos for image captcha recognition (numbers only)

161 **3.3.2 numbers + lower case characters**

162 Here, we used digits from 0 to 9 and lower-case characters from 'a' to 'z' to compose a
163 4-digits captcha image. Below are the accuracies on test set and mean loss value on training
164 set we get during the training process.



165 Figure 8: Accuracies and loss value during the training process (numbers + lower case)

166 (sampling rate： 20 iterations/sample)

167



label: j v 6 p
predict: j v 6 p

label: t q r g
predict: t q r g

label: f k z 0
predict: f k z o

label: f s a i
predict: f s a i

label: s 7 u p
predict: s 7 u p

label: m a u g
predict: m z u g

168

169 Figure 9: Test demos for image captcha recognition (numbers + lower case)

170

171

172

173

174
175
176

### 3.3.3 numbers + lower and upper-case characters

Here, we used digits from 0 to 9 and lower-case characters from 'a' to 'z' and upper-case characters from 'A' to 'Z' to compose a 4-digits captcha image. Below are the accuracies on test set and mean loss value on training set we get during the training process.
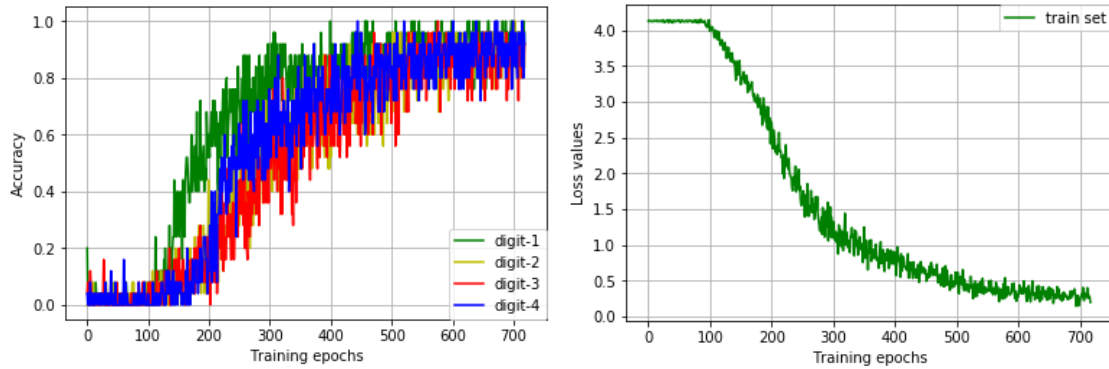


Figure 10: Accuracies and loss value during the training process

(numbers + lower, upper case)



label: d l r V
predict: d I c V

label: 9 u g G
predict: 9 u 8 G

label: F 7 h R
predict: F T h R

label: D D a i
predict: D a a i

label: i O a X
predict: i 0 a X

label: j A b N
predict: j A b V

Figure 11: Test demos for image captcha recognition (numbers + lower, upper case)

190 ### 3.3.4 Median blur

191 In this part, we before feeding captcha images into AlexNet, we preprocessed them with
192 median blur to remove noise like the lines and spots on the images. Below are the accuracies
193 on test set and mean loss value on training set we get during the training process. (We
194 compare with part 3.3.2, captcha built from numbers and lower case character)
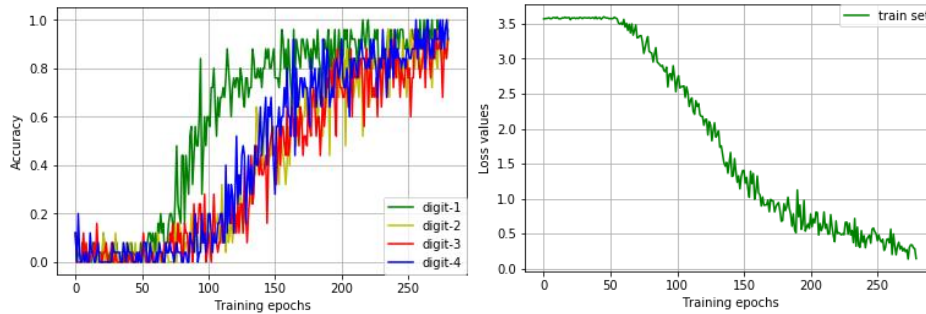


195 Figure 12: Accuracies and loss value during the training process (numbers + lower case)

196 (with image smoothing, sampling rate： 20 iterations/sample)

197

198 As we can see, after applying image smoothing, accuracy of the model doesn't get improved
199 much as we expected since we have got a high recognition accuracy in the previous part, but
200 our model can converge fasted compared with experiments in the previous part without
201 denoising (same learning rate). Here, it took 5620 iterations for convergence, but in part
202 3.3.2 without smoothing, it took nearly 9000 iterations for convergence.

203

204

205 ## 4 Conclusions and Discussions

206 Trough the whole project, first we proposed a novel idea of using convolutional neural
207 network to recognize captcha images. Here, different from traditional image classification
208 problems (since the simplest 4-digit number captcha has $10^4$ kinds of labels), we proposed a
209 joint training method to train the model and predict the character on each digit with a
210 SoftMax layer separately. Below are some interesting conclusions we get from the results
211 from the previous part.

212 ● **Accuracy**: Our system can reach a high accuracy of more than 95% on each digit on test
213 set generally.

214 ● **Training process**: During the training process, since we randomly initialize the parameters
215 in the SoftMax layer, we suffered a problem called "cold start", that is, during the initial
216 hundreds of iterations, accuracies and losses would stay flattened. And then accuracy on the
217 first digit would began to raise (the green line above) ahead of other digits, which is because
218 of the joint training method we use.

219 ● **Drawbacks**: From the running demos above we can see, although our model has reached a
220 high accuracy on test set, but there are still some mistakes like 0-o, 9-g, 7-T etc. This may
221 because the characters are not written in a clear and neat format, and even for a human, it is
222 still a little difficult to distinguish these pairs sometimes. Also, denoise operation help little
223 with this situation.

224

225

226

227

## 5    Reference

[1] Yann LeCun. (Proc. IEEE 1998) Gradient-Based Learning Applied to Document Recognition.

[2] Alex Krizhevsky, Ilya Sutskever, Geoffrey E. Hinton. (2012) ImageNet Classification with Deep Convolutional Neural Networks.