

# *Research and Application of Health Code Recognition Based on Paddle OCR under the Background of Epidemic Prevention and Control*

Dan Zhang, Yunjie Li\*

*ETI Center, Shenzhen Polytechnic, Shenzhen, Guangdong, China*

*\*Corresponding author: yxdbczd@163.com*

**Keywords:** Epidemic prevention and control, paddleocr, health code recognition, target text location algorithm

**Abstract:** Normalization of epidemic prevention and control, in order to solve the problems of low efficiency and high error probability of manual review of health code and trip card pictures by epidemic prevention personnel, an automatic health code recognition method based on paddleocr is proposed. This method uses the open source paddleocr technology for image recognition, and uses the target text location algorithm to output the required text. Practice shows that this method has good recognition effect and high accuracy. Applying this method to the campus epidemic prevention and control system can greatly improve the audit efficiency.

In recent years, the epidemic prevention and control work has become normal, but the epidemic prevention and control situation is still severe and complex. Universities and colleges have become key and difficult units in epidemic prevention and control because of their concentrated personnel, large scale and large mobility<sup>[1]</sup>. In order to implement the requirements of governments at all levels on the prevention and control of COVID, universities are working hard to implement the prevention and control work. Some universities have made use of information technology to build epidemic prevention and control systems<sup>[2]</sup>. In order to ensure the health and safety of teachers and students, health management, travel management and enrollment management should be strengthened. Health codes and travel cards have played an important role in epidemic prevention and control<sup>[3]</sup>. Teachers and students need to check the health code and travel card to enter the school. Those who are at risk are not allowed to enter the school. In order to strengthen the management of all kinds of staff and avoid waiting for teachers and students to gather at the school gate, our school has developed an epidemic prevention and control system. Before entering the school, teachers and students should submit the health code and travel card screenshot in the system, and pass the examination of the department's epidemic prevention specialist (each college or department arranges a special person responsible for epidemic prevention related matters, hereinafter referred to as epidemic prevention Specialist). Only after passing the examination can enter the campus. Therefore, the epidemic prevention specialist needs to review a large number of health code/travel card admission applications every day, and needs to open the health code and travel card pictures bit by bit to verify the picture information. People who are not green code, nucleic acid test is not in the specified time,

travel card route risk areas and other people who do not meet the requirements for admission are not allowed to enter the school.

Manual audit of health code, travel card picture, not only inefficient, easy to make mistakes, but also greatly affect the efficiency of business process audit, resulting in a very bad user experience. Each epidemic prevention specialist has a large number of staff, especially for students, who are engaged in tedious mechanical audit work for a long time, which brings not only physical fatigue, but also great psychological burden. It is urgent to seek technical means to change this situation. OCR technology based on deep learning can realize picture and text recognition, which can solve the above business pain points and improve the audit efficiency of epidemic prevention commissioners.

## 1. Introduction to OCR

### 1.1 Introduction to OCR technology

OCR(Optical Character Recognition) refers to the automatic recognition of text content in images, and is one of the important branches of computer vision<sup>[4]</sup>. Traditional OCR adopts pattern recognition technology, which has many drawbacks such as too many processing links, long process, poor recognition flexibility and difficult maintenance<sup>[5]</sup>. In recent years, with the rapid development of deep learning technology, the traditional OCR technology framework has been broken. OCR technology based on deep learning has gradually become a research hotspot, and is widely used in many fields such as electronic bill recognition, certificate recognition, license plate recognition, natural scene text recognition and so on<sup>[5]</sup>. Deep learning has strong image feature learning ability, which optimizes the traditional technical framework to a certain extent and improves the recognition effect and speed of OCR. Currently, the typical OCR recognition process includes three stages: preprocessing, text detection and character recognition<sup>[6]</sup>. As shown in Figure 1, the technical bottlenecks affecting the recognition accuracy lie in text detection and character recognition, which are the top priority of OCR technology research.



Figure 1: OCR recognition process

**Preprocessing:** usually it is to correct image imaging problems, including image noise reduction, tilt correction, blur removal, etc. At present, CNN based neural network model enhancement feature extraction is widely used to solve the above problems<sup>[7][8]</sup>.

**Text detection:** locate the text area in the image. Target recognition algorithm based on convolutional neural network is applied to text detection. Scholars constantly propose improved algorithms and adjust models, such as TextBoxes, CTPN, SEGLINK, EAST and other algorithms, to make the detection results more accurate<sup>[9]</sup>.

**Character recognition:** identify the text content in the image. On the basis of traditional single character recognition, the context sequence information is introduced, and the neural network model CRNN, which depends on the temporal relationship, can be used to improve the accuracy of character recognition<sup>[10]</sup>.

### 1.2 PaddleOCR

PaddleOCR is a set of rich, advanced and practical OCR tool library realized by Baidu based on deep learning technology, which is open source. It is an ultra-lightweight OCR system. PaddleOCR

has been continuously optimized, improved and updated since its release on May 14, 2020. PaddleOCR open-source a variety of language recognition models, the most commonly used PP-OCR based ultra-lightweight Chinese and English OCR model, universal Chinese and English OCR model, while also providing Japanese, Korean, French, German and other more than 80 language recognition models<sup>[11]</sup>. PaddleOCR provides a variety of text detection training algorithms (EAST, DB) and a variety of character recognition training algorithms (Rosetta, CRNN, STAR-Net, RARE), supports user-defined training, provides rich predictive reasoning deployment schemes, supports PIP quick installation, It can run on Windows, Linux, MacOS and other systems<sup>[12] [13]</sup>.

In view of the above excellent features of PaddleOCR, this paper proposes a health code image recognition method based on PaddleOCR. This method is composed of PaddleOCR image recognition and target text location algorithm. The specific identification process is as follows: First, PaddleOCR technology is adopted to identify the image text, and then the text required by the target text positioning algorithm is used to output the text required by the target, so as to realize the automatic identification of health code/travel card images. Finally, the research results are applied to the epidemic prevention and control system of our school to realize the batch audit and automatic audit of business processes and greatly improve the audit efficiency of epidemic prevention specialists.

## **2. PaddleOCR based health code recognition method**

### **2.1 PaddleOCR recognition**

#### **2.1.1 Setting up the environment**

- (1) Install the python environment;
- (2) Install PaddlePaddle;
- (3) Download PaddleOCR repository and install PaddleOCR dependency package;
- (4) Install CUDA (recognize using GPU) environment.

The detailed installation steps will not be covered here, but can be found in the PaddleOCR open source community.

#### **2.1.2 Health Code recognition**

After the environment is ready, a health code image (Yuekang code image is taken as an example in this paper) is sent to PaddleOCR for recognition.

PaddleOCR can output all text box coordinates recognized, text content and text confidence. The result of paddleocr is unstructured and uncorrelated data printed by line and text box, much of which is not what we care about. Our concerns for health code audits are as follows:

- (1) Picture time, the screenshot submitted according to the requirements should not be 12 hours earlier than the submission time;
- (2) Health code color, whether it is green code;
- (3) Nucleic acid test information, including nucleic acid test results, test time;
- (4) Vaccination information, including vaccination injection information and the last vaccination time.

Obtain the above target text and output it in a structured way to assist the epidemic prevention specialist in health code review.

## 2.2 Target text location algorithm

### 2.2.1 Introduce coordinates

The ocr method of PaddleOCR can output the coordinates of all the recognized text boxes. Each text box has four coordinates, which are respectively  $(x1,y1)$ ,  $(x2,y2)$ ,  $(x3,y3)$ , and  $(x4,y4)$ . Each coordinate is marked clockwise with the top left coordinate as the vertex, as shown in Figure 2.

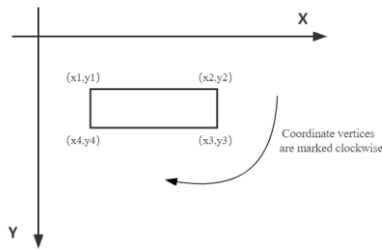


Figure 2: Text box coordinates

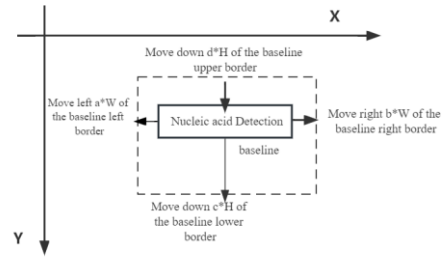


Figure 3: Target identification box

### 2.2.2 Select a benchmark

Take the benchmark in the picture, using the target text "Nucleic acid Detection" as an example. Take the first text box that contains "nucleic acid test" as the baseline. Get the width and height of each Chinese character in the benchmark by using the formula as follows:

$$W = ((x2 + x3)/2 - (x1 + x4)/2) / L$$

$$H = (y3 + y4)/2 - (y1 + y2)/2$$

Where  $W$  represents the average width of each character in the datum;  $L$  represents the number of characters in the benchmark;  $H$  represents the height of each character in the baseline.

### 2.2.3 Determine the target identification box

Take the coordinates of the current reference box as the reference, expand up, down, left and right according to the target content you want to output to get the target identification box, as shown in Figure 3.

The target identification box is based on  $a*W$  left shift of the baseline left border,  $b*W$  right shift of the baseline right border,  $c*H$  down of the baseline lower border, and  $d*H$  down of the baseline upper border. According to the actual layout of the picture and the position of the target box, the offset coefficient can be adjusted flexibly. For example, with "nucleic acid detection" as the reference frame, the target frame is based on the reference left frame moved left  $1.4*W$ , right frame moved right  $0.8*W$ , lower frame moved down  $5.2H$ , upper frame moved down  $1.2*H$ ; With "Member Management" as the reference frame, the target frame is moved left  $14*W$  based on the left frame of the reference frame, right frame is moved right  $-2*W$ , bottom frame is moved down  $20*H$ , and top frame is moved down  $1.5*H$ .

### 2.2.4 Result output

Based on the above algorithm, "nucleic acid detection", "COVID vaccine" and "member management" are selected as the benchmarks to output the target text content.

Get the target identification box part implementation code as follows:

```

#key: benchmark value
BMK = {'nucleic acid detection':[1.4,0.8,5.2,1.2],'COVID vaccine':[3,3,5.5,1],' member
management':[14,-2,20,1.5]}
def getTargetByBase(ocrResult,key):
    targetBox = None
    for line in ocrResult:
        box = line[0]
        txt = line[1][0]

        if key in txt:
            x1,y1 = box[0]
            x2,y2 = box[1]
            x3,y3 = box[2]
            x4,y4 = box[3]
            W = round(((x2+x3)/2-(x1+x4)/2)/len(txt),2) # W=((x2+x3)/2-(x1+x4)/2)/L
            H = round((y3+y4)/2 - (y1+y2)/2,2) # H=(y3+y4)/2-(y1+y2)/2

            targetBox1 = [round((x1-BMK[key][0]*W),2), y1+BMK[key][3]*H]
            targetBox2 = [round((x2+BMK[key][1]*W),2), y2+BMK[key][3]*H]
            targetBox3 = [round((x3+BMK[key][1]*W),2), round((y3+BMK[key][2]*H),2)]
            targetBox4 = [round((x4-BMK[key][0]*W),2), round((y4+BMK[key][2]*H),2)]
            targetBox = [targetBox1,targetBox2,targetBox3,targetBox4]
    return targetBox
return targetBox

```

## 2.3 Results and analysis

In order to evaluate the effectiveness of the algorithm in this paper, 852 health code pictures were collected for experimental analysis, including some special fonts and pictures presented on behalf of relatives. The identification statistics of key fields of health code are shown in Table 1:

Table 1: Key field identification statistics of health code

Total number of pictures	Keywords	Identify accurate	Identify inaccurate	Accuracy rate	Average accuracy rate
852	nucleic acid detection	849	3	99.65%	97.34%
	COVID vaccine	851	1	99.88%	
	member management	788	64	92.49%	

It can be seen from Table 1 that the identification accuracy of the algorithm in this paper for the keywords "nucleic acid detection" and "COVID vaccine" is over 99%, while the accuracy of "member management" is relatively low at 92.49%. According to the analysis, it is mainly caused by the large span of the target box of "member management" and the special font. The average accuracy rate of the algorithm in this paper is 97.34%. Generally speaking, the algorithm in this paper achieves good results in health code recognition. However, if the target text positioning algorithm is used to obtain the text content in the target box, the span of the target box should not be too large, the more compact the effect is better.

The algorithm in this paper is not only suitable for the recognition of health code, but also suitable for the recognition of travel card. Select "Please accept the green travel card" and "You arrived or passed through within the previous 14 days" as the benchmark, adjust the benchmark parameters, and obtain the mobile phone number, update time and route city information of the travel card.

### 3. Health code identification application

#### 3.1 Interface Encapsulation

Flask, a python lightweight Web application framework that is flexible, portable and safe, was adopted to develop the application interface for the encapsulation of the interface. Images were imported by POST method. The text required by the target was identified by the above algorithm, and then cleaned by regular expression matching (if there are multi-recognized characters or text, it is filtered out). Finally, the identification result in JSON format is returned<sup>[14]</sup>.

For example, call the health code identification interface, pass a health code screenshot, the interface returns the result:

```
{
  "data": {
    "Health Code":{
      "Health Code Results": "Green Code",
      "Health code time": "04-02 16:35:49"
    },
    "Nucleic acid detection":{
      "Nucleic acid test result": "48 hours negative",
      "Nucleic acid detection time": "2022-04-02 03:31"
    },
    "COVID Vaccine":{
      "COVID Vaccine Results": "The whole process of vaccination has been completed",
      "Vaccination time": "2021-12-10"
    }
  }
}
```

Call the travel card identification interface, pass in a travel card screenshot, interface returns the result:

```
{
  "data": {
    "Mobile number": "181 * * * 9406",
    "Update time": "2022.03.01 10:13:53",
    "Cities by Way": ["Dongguan City, Guangdong Province *", "Dongguan City, Guangdong Province *"],
    "Whether it passes through medium and high risk areas": "Yes"
  }
}
```

#### 3.2 Application of epidemic prevention and control system

Since 2022, the situation of epidemic prevention and control has been severe and complex. According to the latest requirements of epidemic prevention and control work, all units must

strengthen personnel management, strictly guard the campus entrance, and avoid loopholes in prevention and control. The school staff must fill in the daily nucleic acid and daily health report, and submit the health code/travel card. The epidemic prevention specialist of the department shall review the health code/travel card in strict accordance with the latest requirements of the school, and the green code for returning to school will be obtained after passing the review. All personnel shall enter and leave the campus with the green code. The specific verification process of health code/travel card is shown in Figure 4.

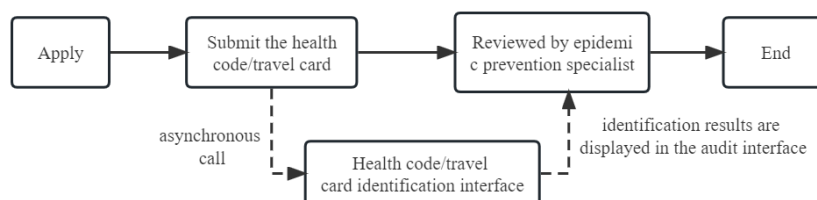


Figure 4: Health code/travel card review process

The on-campus staff submits the screenshot of the health code/travel card in the epidemic prevention and control system to the epidemic prevention specialist for review. The system asynchronously calls the identification interface of the health code/travel card, and the identification result is displayed on the review interface as auxiliary identification information. Health code identification results mainly display: nucleic acid detection results, detection time, health code time, health code color and other information; Travel card identification results mainly display: mobile phone number, update time, way city and other information. According to the requirements of the school's latest audit rules (constantly adjusted according to the epidemic prevention and control situation), if there are problems in the identification results, such as nucleic acid test results within 48 hours, nucleic acid test identification results within 72 hours, the red bold font will indicate "abnormal health code identification results", if all the results meet the requirements, the green font will indicate "no abnormal identification results"; If the identification matches the risk city of the travel card route, the red eye-catching font will prompt "abnormal identification result of the travel card", otherwise, "no abnormal identification result". If the epidemic prevention specialist enters the list to be reviewed, he or she can screen those who are not abnormal and review them in batches. For those who are abnormal, the process can be returned or terminated in batches.

The application of health code/travel card identification interface has greatly improved the audit efficiency of epidemic prevention specialists and reduced the working pressure of epidemic prevention specialists. The follow-up can even consider automatic audit according to the identification results, completely abandoning the mechanical manual audit mode, freeing the manpower, shortening the audit time, and improving the sense of gain and happiness brought by information technology to teachers and students.

#### 4. Summary

In the current situation of normal epidemic prevention and control, this paper proposed a PaddleOCR health code recognition method, using PaddleOCR technology for picture text recognition, using target text location algorithm to output the text required by the target, and applied the research results in the campus epidemic prevention and control system, to assist the epidemic prevention commissioner to batch review the health code/travel card. It greatly reduces the pressure of epidemic prevention commissioners, improves the information experience of teachers and students, and obtains good social benefits. It is a successful case that OCR technology based on deep learning

is applied to the daily life of teachers and students. In recent years, the rapid development of artificial intelligence technology, how to apply the new technology in the construction of campus information, it is worth our thinking. We need to constantly deepen and refine the demand for information business, combine the existing information system, make full use of new technology, create a new innovative service model, solve the pain points and difficulties in the information construction of colleges and universities, and constantly provide higher quality information services for teachers and students.

## References

- [1] Cen J. *Research on the judicial application of the crime of impairing the prevention and control of infectious diseases under the background of novel coronavirus epidemic*[C]// 2020 International Conference on Public Health and Data Science (ICPHDS). 2020.
- [2] Lao Chen, Jinyang Liu, Yunhui Du, et al. *Design and development of the prevention and control system Model of COVID in colleges and universities* [J]. *Chinese Journal of ICT in Education*, 2021(05): 76-79.
- [3] Jian Chen, Hong-Lin Chen, Jia-ning Li. *Study on the development of health codes in China in the post-epidemic era* [C]. *Proceedings of 2021 Annual Science and Technology Conference of Chinese Society for Environmental Sciences (III)*. 2021: 508-515.
- [4] Hongtao Lu, Mukun Luo. *Survey on New Progresses of Deep Learning Based Computer Vision* [J]. *Journal of Data Acquisition and Processing*, 2022, 37(2): 247-278.
- [5] Rihua Wang. *Research on Key Technologies and Application of Intelligent OCR Recognition Based on Deep Learning* [J]. *Designing Techniques of Posts and Telecommunications*, 2021(08): 20-24.
- [6] Tan J, Ma Y, Men K, et al. *Tuberculosis in Epidemic Prevention and Control Based on Big Data Technology*[J]. *Journal of Physics: Conference Series*, 2021, 1881(4):042036-.
- [7] Thi Tuyet Hai Nguyen, Adam Jatowt, Mickael Coustaty, Antoine Doucet. *Survey of Post-OCR Processing Approaches* [J]. *ACM COMPUTING SURVEYS*, 2021, 54(06): 1-11.
- [8] Xiaopei Hou, Ying Gao. *Application of Convolutional Neural network CNN algorithm to text classification* [J]. *Science and Technology & Innovation*, 2019, 0(4): 158-159.
- [9] Zhi Tian, Weilin Huang, Tong He, Pan He, Yu Qiao 0001. *Detecting Text in Natural Image with Connectionist Text Proposal Network*. [J]. *Journal of CoRR*, 2016, abs / 1609.03605.
- [10] Shi Baoguang, Bai Xiang, Yao Cong. *An End-to-End Trainable Neural Network for Image-Based Sequence Recognition and Its Application to Scene Text Recognition*. [J]. *IEEE transactions on pattern analysis and machine intelligence*, 2017, 39(11).
- [11] Baidu open source OCR introduction [EB/OL]. [2021-12-21]. <https://gitee.com/computer-vision/PaddleOCR#/>
- [12] Jianmin Qiu. *Improvement and Practice of Open Source PaddleOCR Technology in Enterprise Business License Recognition* [J]. *Modern Information Technology*, 2021, 5(09): 65-69+74.
- [13] Xing Wang, Yongfeng Zheng, Yongbing Yan, et al. *Research on Ticket Recognition Algorithm based on OCR Technology* [J]. *Intelligent Computers and Applications*, 2021, 11(11): 101-106.
- [14] Huaqiao University. *A method and system for Structural output of OCR recognition Results: CN201910145824.0* [P]. 2019-06-07.