

OCR及领域优化

孟灿 孙雨豪 史宇辰 丁欣 董冠辰

目录

01 OCR简介、历史及应用

02 OCR相关技术

03 OCR领域现存问题及优化

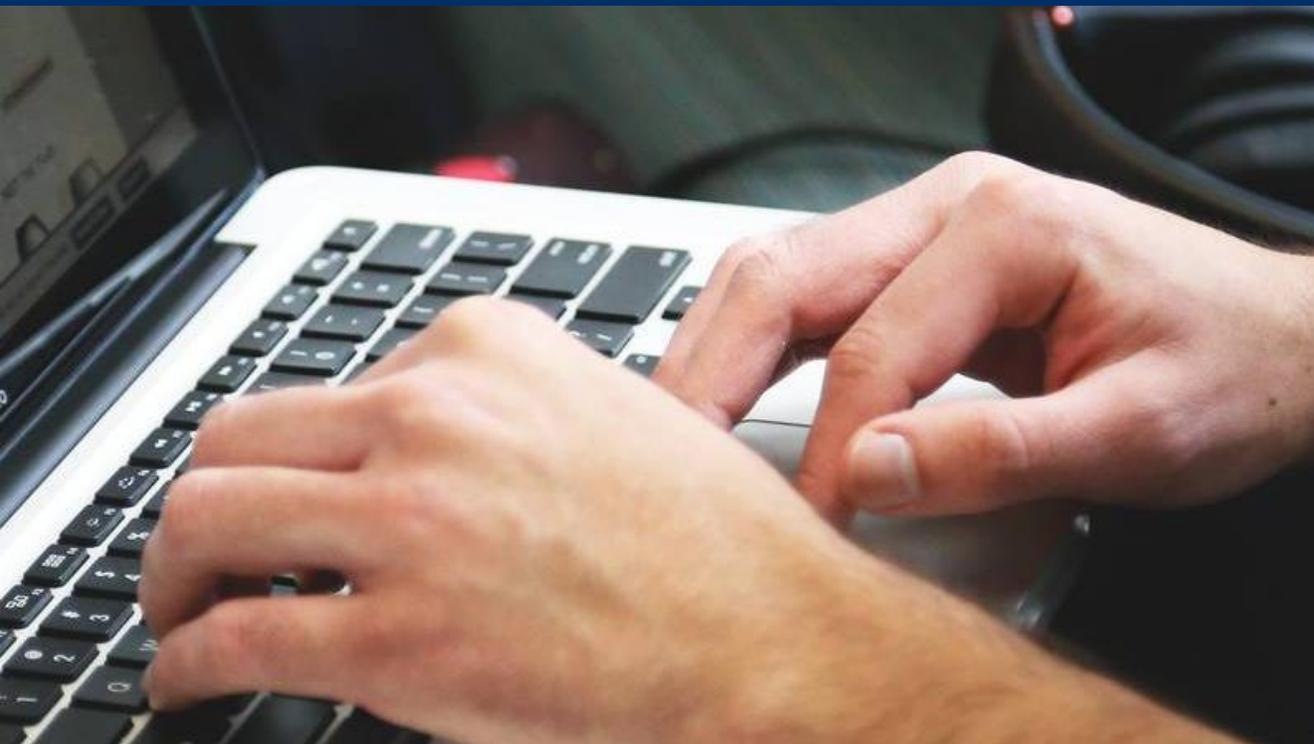
04 OCR未来发展方向

05 Demo展示



1 OCR简介、历史及应用

1.1 OCR简介



光学字符识别 (Optical Character Recognition) 是将文本的图像资料经过扫描图像机械对图像进行识别转换为机器编码的文本的过程。将其视为转换模拟数据，数字化的过程。

OCR历史

1929年

光学字符识别的概念产生

1957年

第一个OCR软件
一家基于光电技术公司Photron的ERA
(Electric Reading Automation)

1965年

第一个真正意义上的OCR产品是
IBM公司推出
IBM418

1966年

IBM公司的Casey和Nagy
发表了第一篇关于中文汉字识别的论文

1968年

东芝公司研制出了世界上第一个用于信函自动分拣的手写邮政编码识别系统

OCR历史

20世纪
70年代初

我国开始对数字、字母及符号等识别技术进行研究

20世纪
70年代末

开始对开始进行汉字字符识别相关的研究工作

1986年

国家“863”
计划

1989年

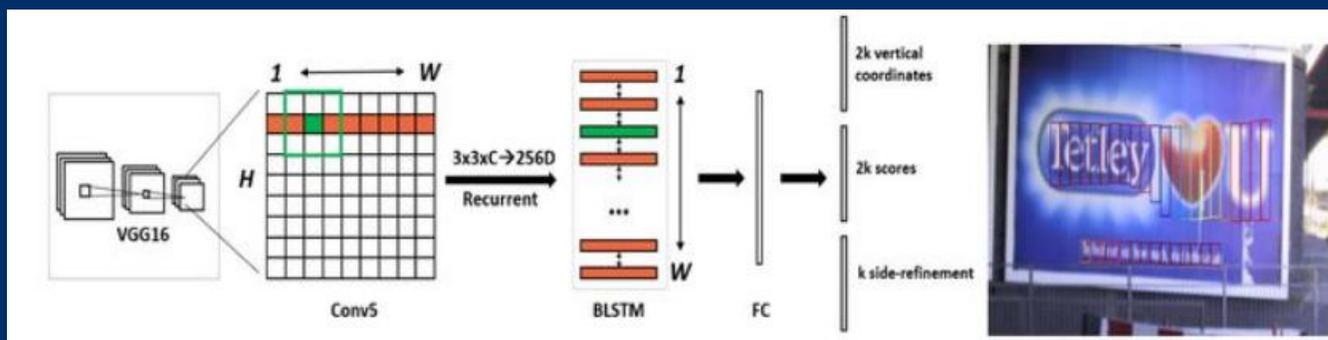
推出了我国首个中文OCR识别系统-清华文通TH-OCR1.0版。

1998年

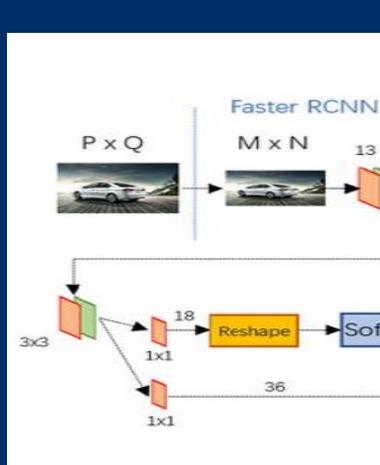
成立北京汉王科技有限公司
汉王制造手写识别技术成功授权微软

OCR发展过程

• 1. CTPN (Connection Text Proposal Network) 模型



• 2. Faster RCNN 模型——Ross B. Girshick



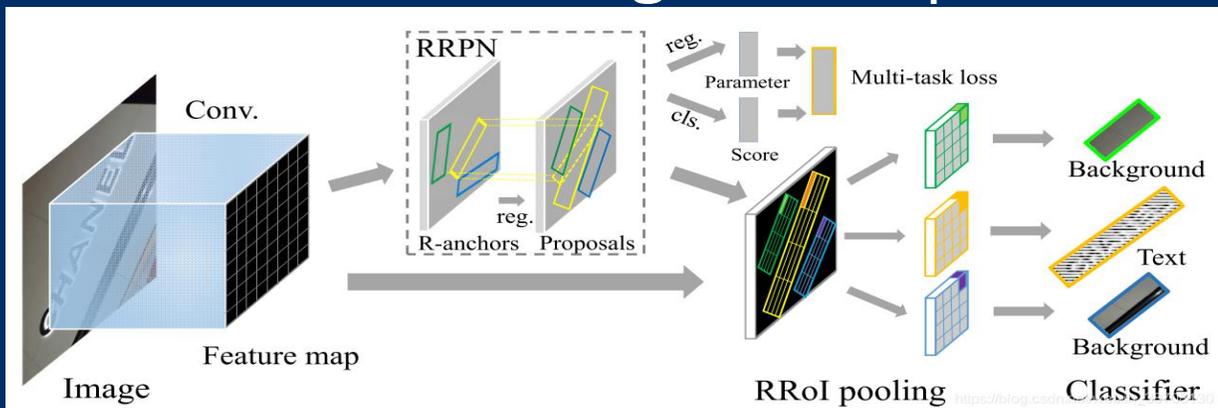
	使用方法	缺点	改进
R-CNN (Region-based Convolutional Neural Networks)	1、SS提取RP; 2、CNN提取特征; 3、SVM分类; 4、BB盒回归。	1、训练步骤繁琐(微调网络+训练SVM+训练bbox); 2、训练、测试均速度慢; 3、训练占空间	1、从DPM HSC的34.3%直接提升到了66% (mAP); 2、引入RP+CNN
Fast R-CNN (Fast Region-based Convolutional Neural Networks)	1、SS提取RP; 2、CNN提取特征; 3、softmax分类; 4、多任务损失函数边框回归。	1、依旧用SS提取RP(耗时2-3s, 特征提取耗时0.32s); 2、无法满足实时应用, 没有真正实现端到端训练测试; 3、利用了GPU, 但是区域建议方法是在CPU上实现的。	1、由66.9%提升到70%; 2、每张图像耗时约为3s。
Faster R-CNN (Fast Region-based Convolutional Neural Networks)	1、RPN提取RP; 2、CNN提取特征; 3、softmax分类; 4、多任务损失函数边框回归。	1、还是无法达到实时检测目标; 2、获取region proposal, 再对每个proposal分类计算量还是比较大。	1、提高了检测精度和速度; 2、真正实现端到端的目标检测框架; 3、生成建议框仅需约10ms。

生成的anchor, 有位置关系以及一定的长宽比例, 整个训练过程通过调节anchor的坐标来达到与ground truth的最大IOU。如果anchor长宽比例设置不合理的话, 对一个长宽比例失调的物体就会很难找到一个将它整个包围的bounding box。

Faster-rcnn采用了nms, 这样会导致互相有遮盖的物体不能同时被检查出来

OCR发展过程

• 3. RRPN (Rotation Region Proposal Network) 模型

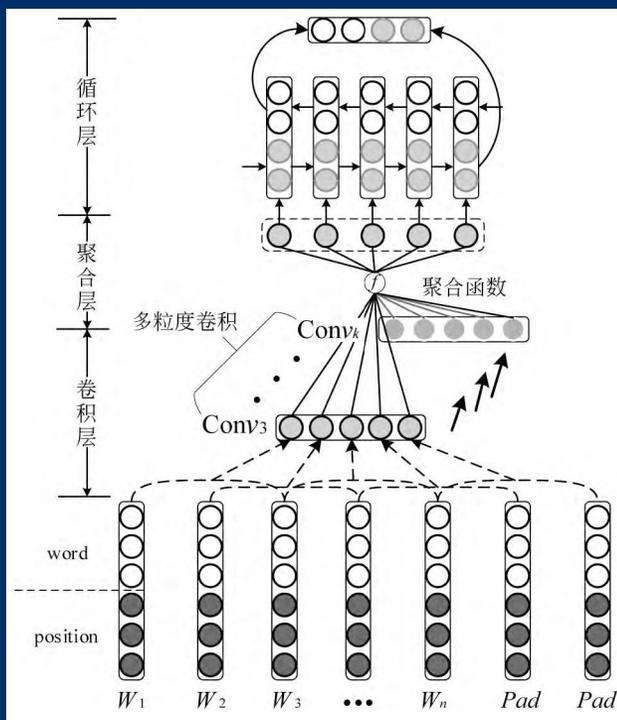


• 4. SegLink 模型

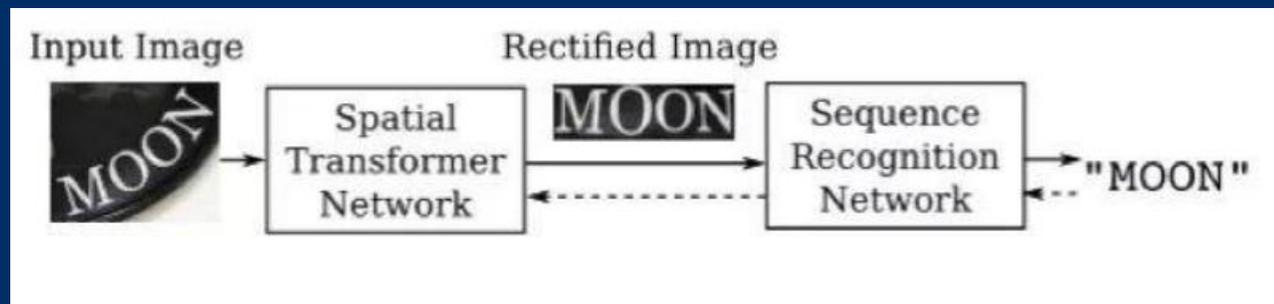


OCR发展过程

• 1. CRNN 模型



2. RARE 模型



OCR应用

1. 数字文档类

学院新闻

- 北理工计算机学院组织开展“北理工精神校友谈”纪念活动
- 北理工计算机学院2021级研究生新生开学典礼暨入学教育顺利召开
- 北理工计算机学院(唐山研究院)2021级研究生顺利开学
- 北理工计算机学院顺利完成2021级研究生新生迎新工作
- 北理工王国仁、刘驰团队获CCF-A类顶会KDD 2021最佳论文

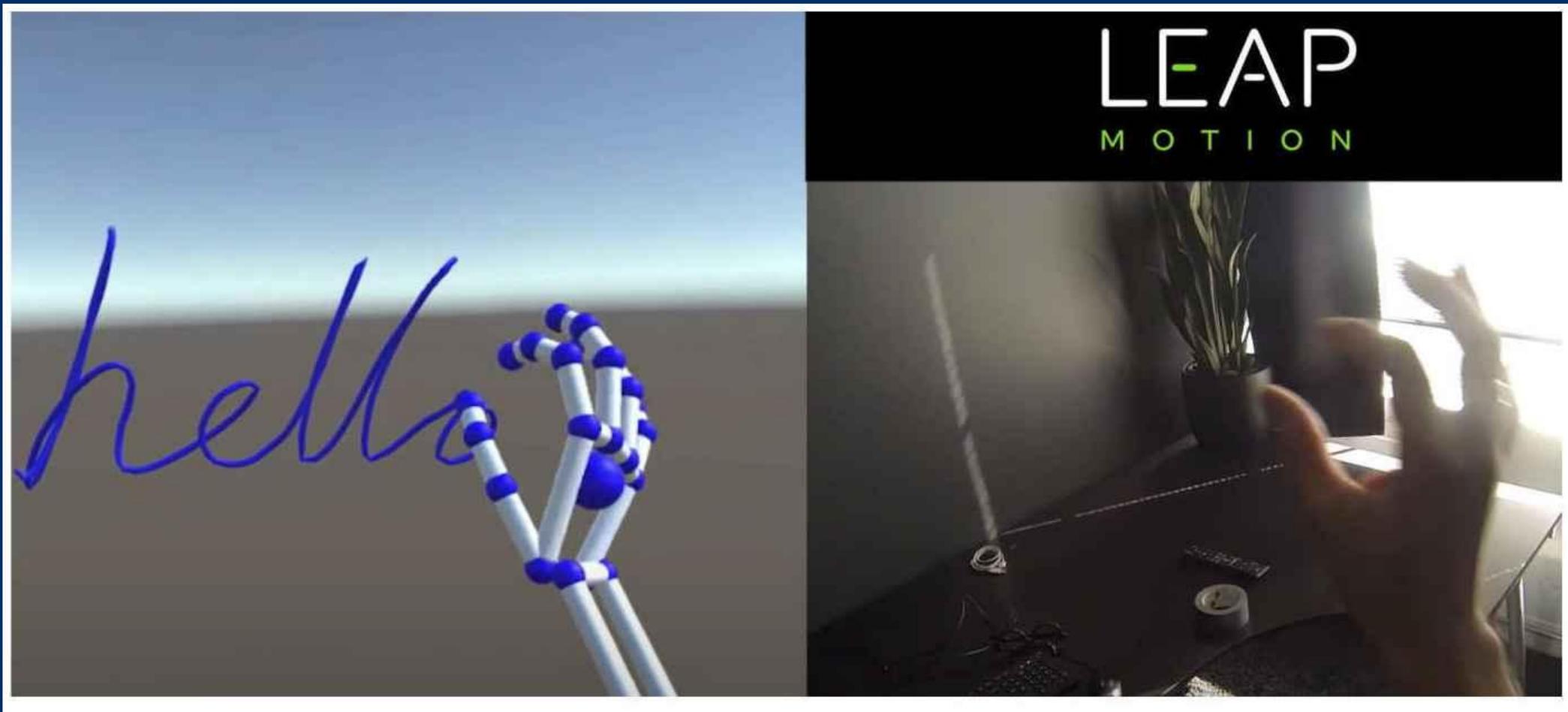
通知公告

- 北京理工大学计算机学院2021年研究生学业奖学金评审细则
- 计算机学院2021年研究生国家奖学金拟推荐人选公示通知
- 北京理工大学计算机学院2022年接收优秀应届本科毕业生推荐免试
- 北京理工大学计算机学院2021年研究生国家奖学金评审细则
- 北京理工大学计算机学院2022年推荐优秀应届本科毕业生免试

- 1: 学院新闻 0.992
- 2: 《北理工计算机学院组织开展“北理工精神校友谈”纪念活动》 0.956
- 3: 《北理工计算机学院2021级研究生新生开学典礼暨入学教育顺利召开》 0.968
- 4: 《北理工计算机学院(唐山研究院)2021级研究生顺利开学》 0.974
- 5: 《北理工计算机学院顺利完成2021级研究生新生迎新工作》 0.94
- 6: 《北理工王国仁、刘驰团队获CCF-A类顶会KDD 2021最佳论文》 0.935
- 7: 通知公告 0.999
- 8: 《北京理工大学计算机学院2021年研究生学业奖学金评审细则》 0.976
- 9: 《计算机学院2021年研究生国家奖学金拟推荐人选公示通知》 0.972
- 10: 《北京理工大学计算机学院2022年接收优秀应届本科毕业生推荐免试》 0.923
- 11: 北京理工大学计算机学院2021年研究生国家奖学金评审细则 0.9
- 12: 《北京理工大学计算机学院2022年推荐优秀应届本科毕业生免试》 0.900

OCR应用

- 2. 手写识别类



OCR应用

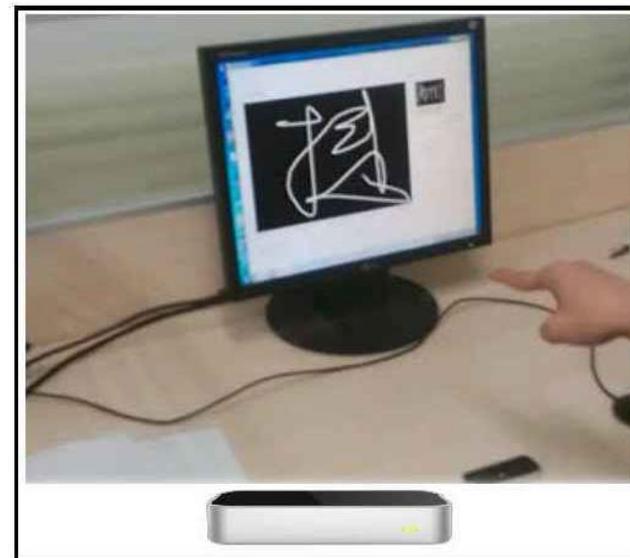
- 2. 手写识别类



(a) 基于穿戴式设备的空中手写



(b) 基于Kinect的空中手写



(c) 基于LeapMotion的空中手写

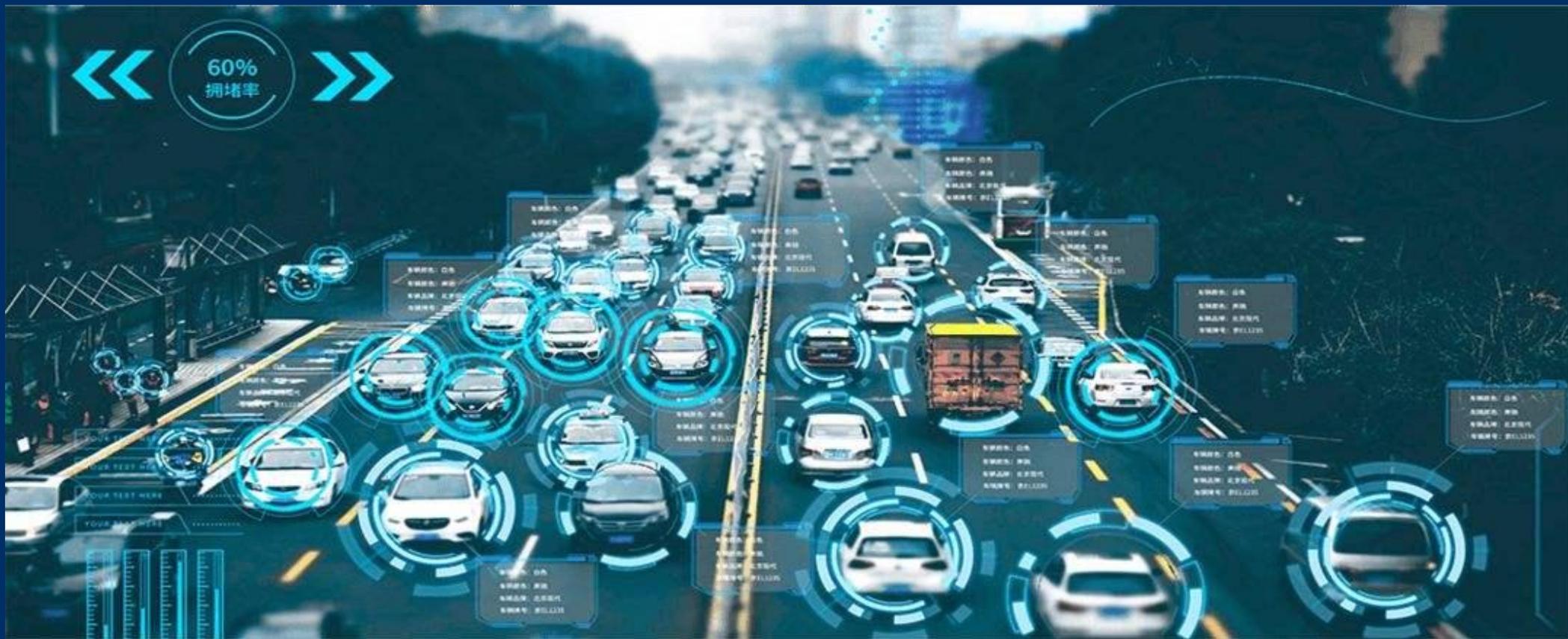
OCR应用

• 3. 自然场景类



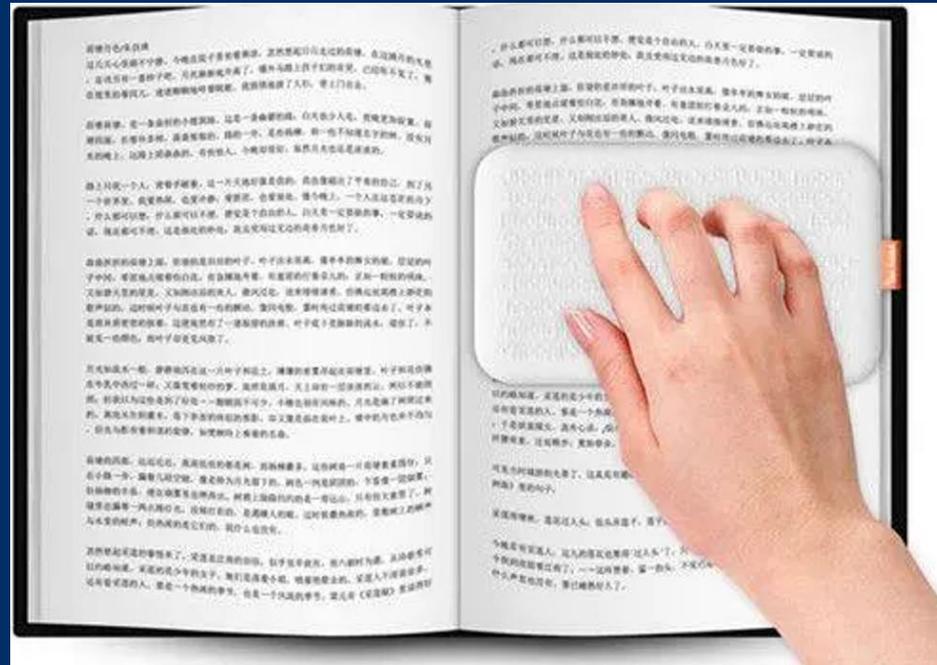
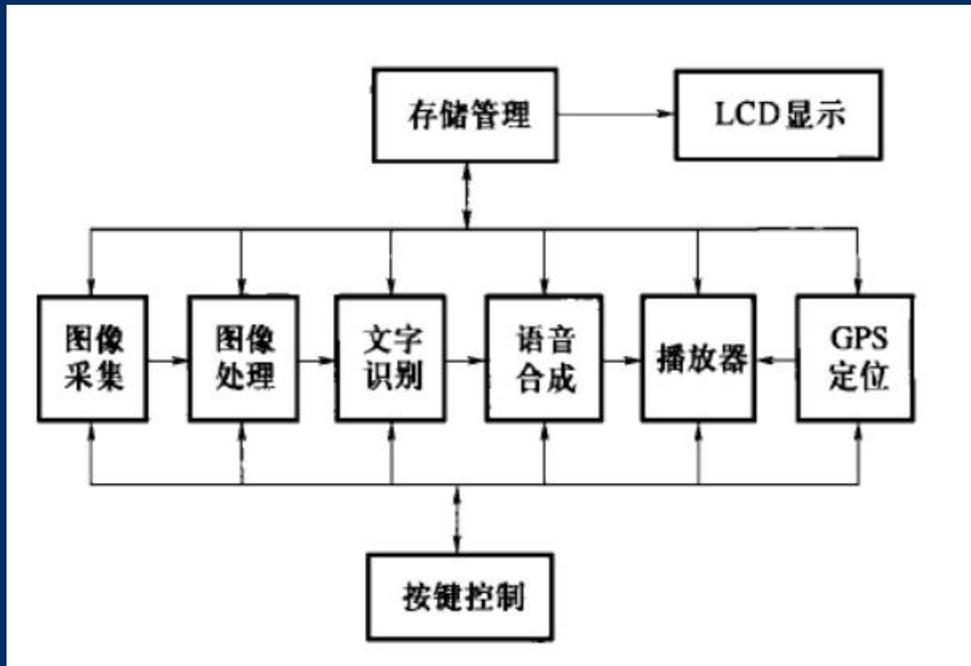
OCR应用

• 3. 自然场景类



OCR应用

• 4. 视觉识别类





2 OCR相关技术

OCR技术实现文字识别的过程



图像输入：读取不同图像格式文件；

图像预处理：将每一个文字图像分检出来交给识别模块识别。主要包括图像二值化，噪声去除，倾斜校正等；

文字检测：将图片中的文字区域位置检测出来，以便于进行后面的文字识别，只有找到了文本所在区域，才能对其内容进行识别。

文本识别：在文字检测的基础上，对文本内容进行识别，主要解决的问题是每个文字是什么。

2.1.1 二值化



图像二值化 (Image Binarization) 就是将图像上的像素点的灰度值设置为0或255, 也就是将整个图像呈现出明显的黑白效果的过程。二值图像也常常用作原始图像的掩模, 它就像一张部分镂空的纸, 把我们对不感兴趣的区域遮掉。进行二值化有多种方式, 其中最常用的就是采用阈值法进行二值化。

2.1.2 噪声去除



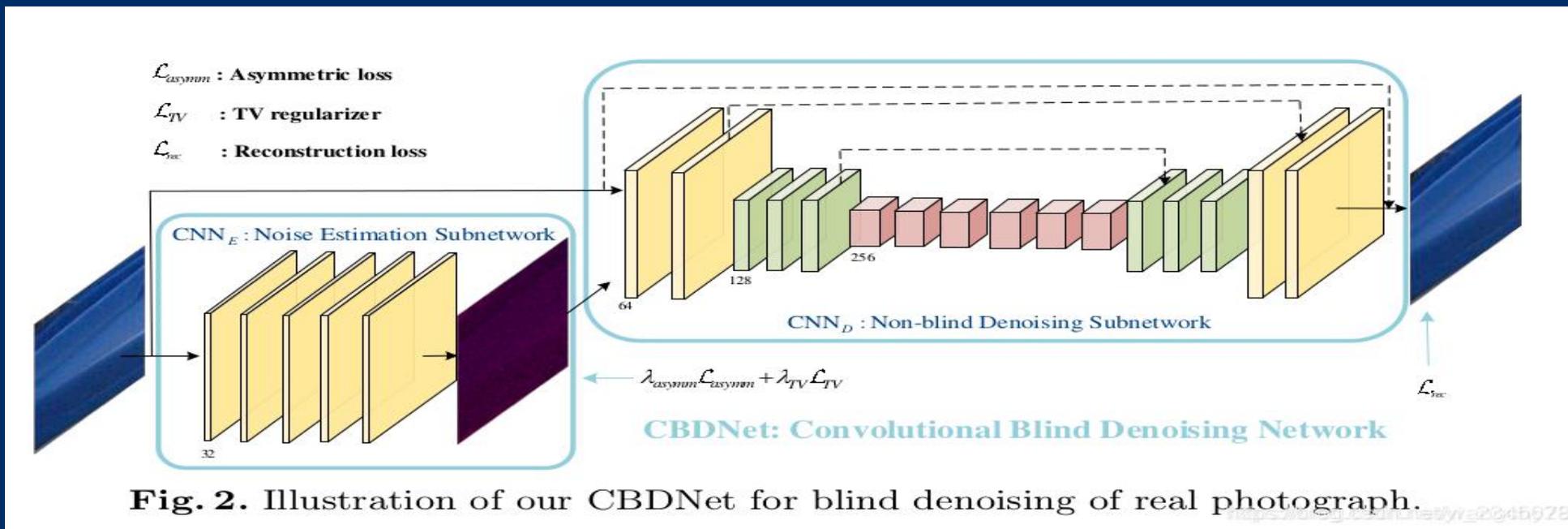
- 平滑噪声的目的：
 1. 模糊
- 比如在拍摄一下大的物体，中间有些小细缝或者小窟窿、小孔；我们希望把这些去掉，跟周围的物体同一个颜色，我们就用平滑的方法，在提取较大的物体前，去除太小的细节；在图片中的小间断的地方，将这些间断链接起来；
 2. 消除噪声
- 改善图像质量，降低干扰；平滑滤波是对低频进行增量的增强，将高频滤除掉；相当于底图滤波器用来消除图像中的随机噪声，起到平滑的作用；

CBDNet

提出了一个更加真实的噪声模型，其考虑了信号依赖噪声和ISP流程对噪声的影响，展示了图像噪声模型在真实噪声图像中起着关键作用。

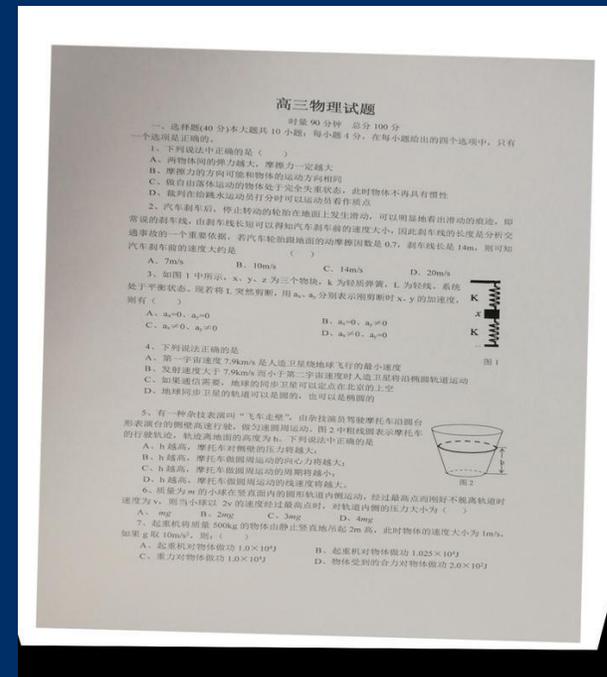
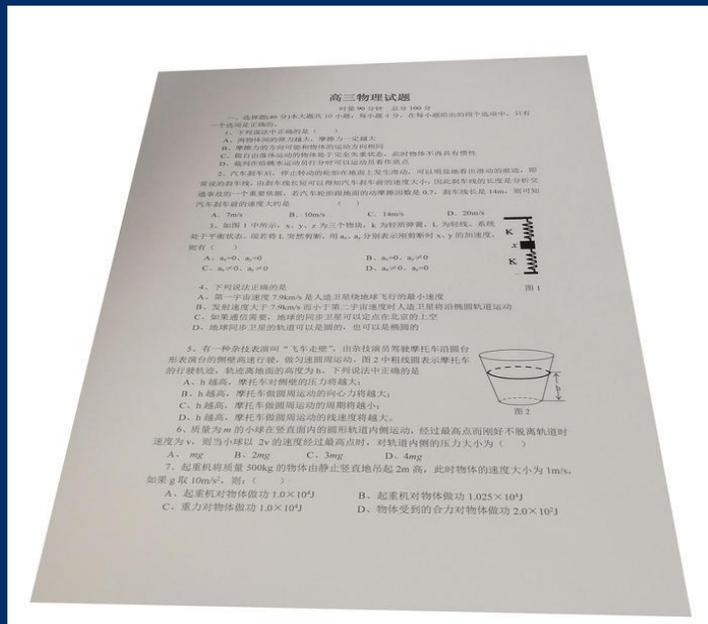
提出了CBDNet模型，其包括了一个噪声估计子网络和一个非盲去噪子网络，可以实现图像的盲去噪（即未知噪声水平）。

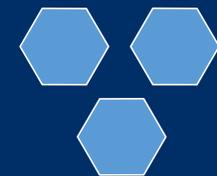
提出了非对称学习的损失函数，并允许用户交互式调整去噪结果，增强了去噪结果的鲁棒性。将合成噪声图像与真实噪声图像一起用于网络的训练，提升网络的去噪效果和泛化能力。



2.1.3 倾斜校正

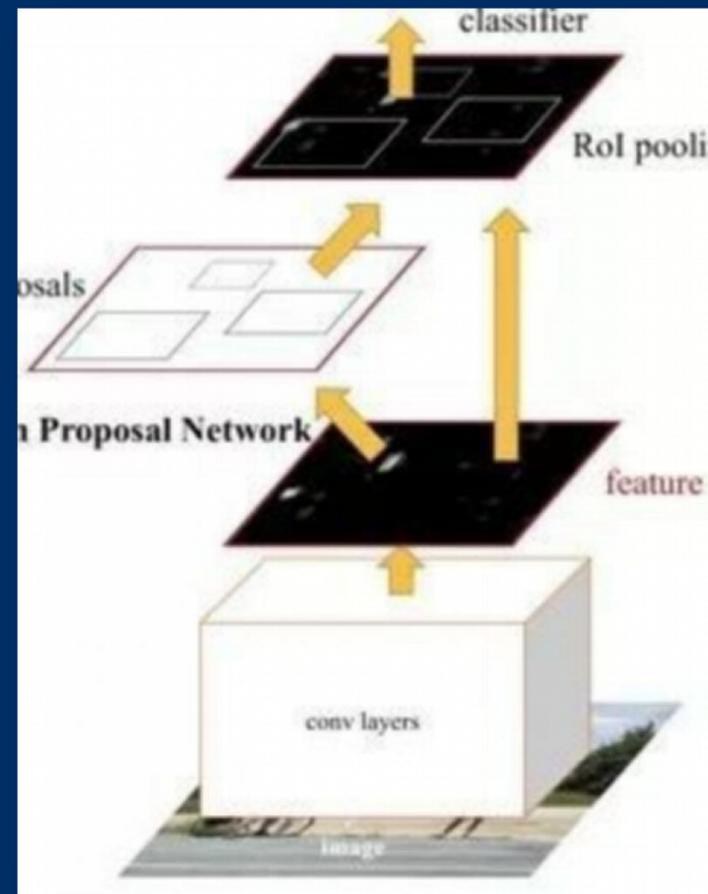
图像倾斜可以分为两种情况，一种是平面倾斜，这种情况下拍照设备与试卷平行，拍出来的图像只需要进行旋转即可完成矫正；另一种是Z轴倾斜，这种情况下拍照设备与试卷存在一定的角度，拍出来的图像要先进行透视变换，然后再进行旋转等操作才可以完成矫正。



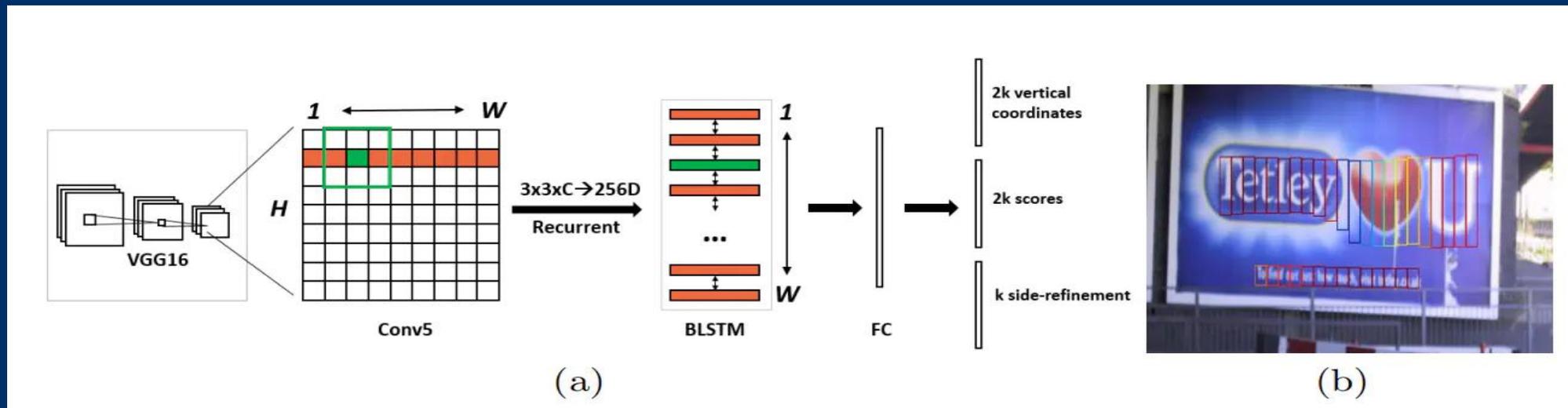


2.2.1 Faster R-CNN

Faster R-CNN采用辅助生成样本的RPN网络，将算法结构分为两个部分，先由RPN网络判断候选框是否为目标，再经分类定位的多任务损失判断目标类型，整个网络流程都能共享卷积神经网络提取的特征信息，节约计算成本，且解决Fast R-CNN算法生成正负样本候选框速度慢的问题，同时避免候选框提取过多导致算法准确率下降。对于受限场景的文字检测，Faster R-CNN的表现较为出色。可以通过多次检测确定不同粒度的文本区域。



2.2.2 文本检测网络CTPN



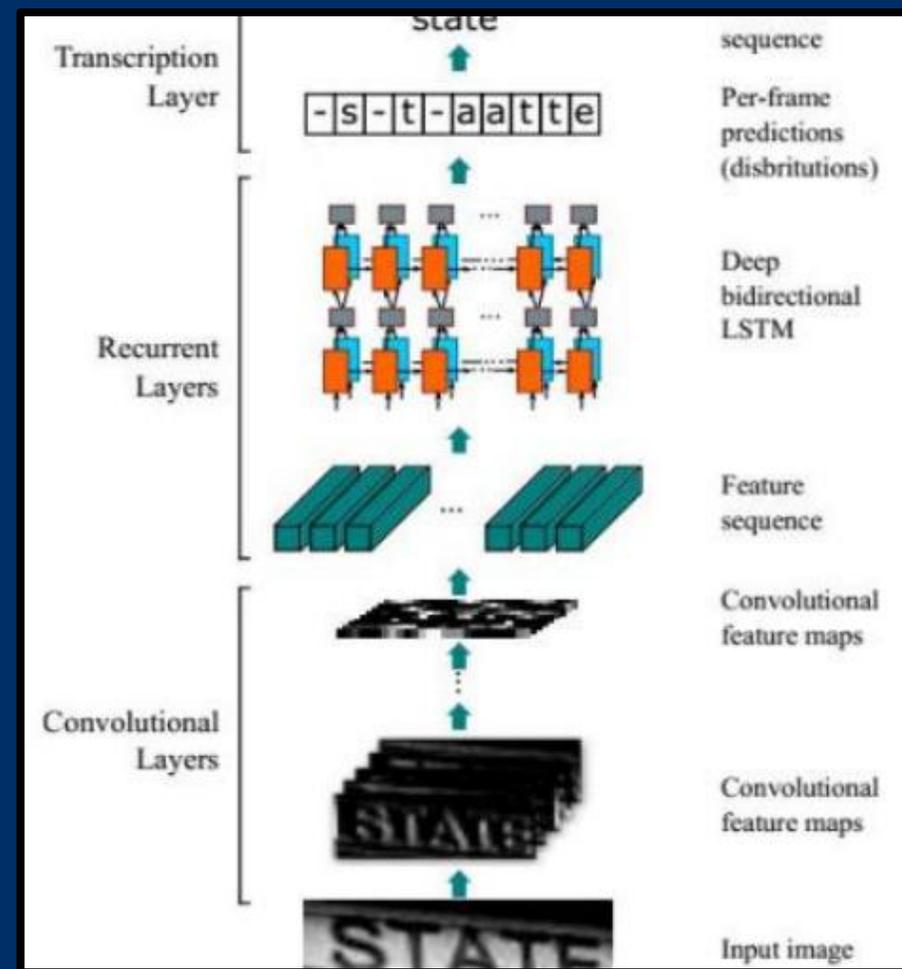
1. 采用VGG16网络提取空间特征，下采样16倍后的卷积特征图Conv5层，输出维度为 $N \times H \times W \times C$ ；
2. 采用3x3大小的滑动窗口在feature map上滑窗，每个窗口得到3x3xC的特征向量，每个点都融合了周围9个点的信息，输出维度为tf: $N \times H \times W \times C$ (或caffe: $N \times H \times W \times 9C$)；
3. 把feature map reshape为 $(NH) \times W \times C$ ，输入BLSTM中提取每一列的特征，输出维度为 $(NH) \times W \times 256$ ，然后reshape为 $N \times 256 \times H \times W$ ；
4. 输入FC全卷积层，输出维度为 $N \times H \times W \times 512$ ；
5. 输入RPN网络，输出三个分支，从上到下分别为：
 - (1) 维度 $N \times H \times W \times 2k$ ，k表示每个像素位置有k个anchor，每个anchor有2个坐标，分别为中心y坐标 v_c 和高度 v_h ；
 - (2) 维度为 $N \times H \times W \times 2k$ ，k表示每个像素位置有k个anchor，每个anchor有2个前背景得分；
 - (3) 维度为 $N \times H \times W \times k$ ，k表示每个像素位置有k个anchor，每个anchor的水平精修side-refinement比例 o 。
6. 得到很多text proposal，使用nms来过滤掉多余的box；
7. 使用基于图的文本行构造算法，将得到的一个一个的box合并成本行。

2.3.1 文本识别网络CRNN



- 整个CRNN网络结构包含三部分，从下到上依次为：
- CNN（卷积层），使用深度CNN，对输入图像提取特征，得到特征图；
- RNN（循环层），使用双向RNN（BLSTM）对特征序列进行预测，对序列中的每个特征向量进行学习，并输出预测标签（真实值）分布；
- CTC loss（转录层），使用 CTC 损失，把从循环层获取的一系列标签分布转换成最终的标签序列。

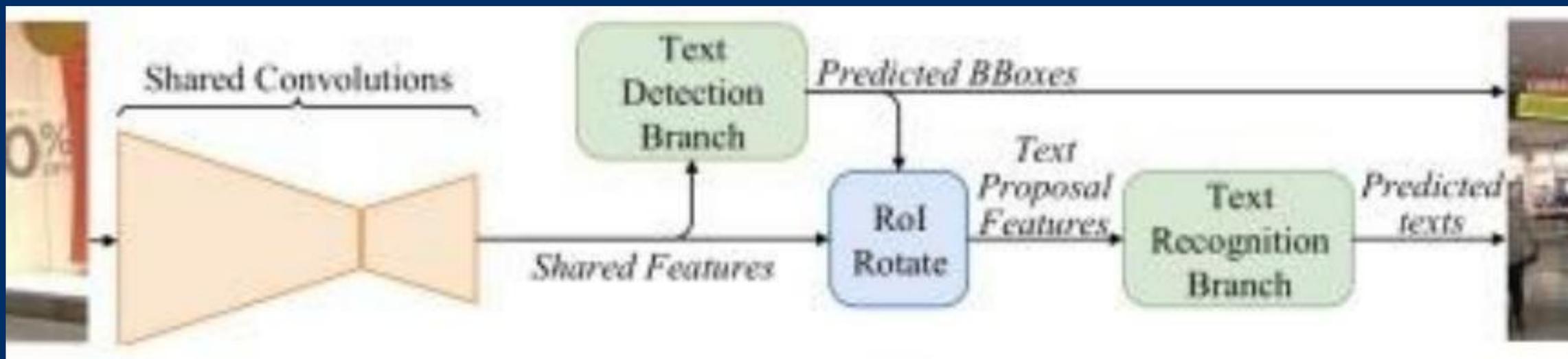
CRNN主要用于端到端地对不定长的文本序列进行识别，不用先对单个文字进行切割，而是将文本识别转化为时序依赖的序列学习问题，就是基于图像的序列识别。



2.3.2 端到端的OCR

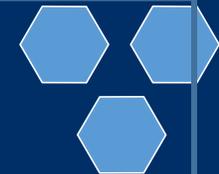


与检测-识别的多阶段OCR不同，深度学习使端到端的OCR成为可能，将文本的检测和识别统一到同一个工作流程中。目前比较受到瞩目的一种端到端框架叫做FOTS(Fast Oriented Text Spotting)。FOTS的检测任务和识别任务共享卷积特征图。一方面利用卷积特征进行检测，另一方面引入了RoIRotate，一种用于提取定向文本区域的算符。得到文本候选特征后，将其输入到RNN编码器和CTC解码器中进行识别。同时，由于所有算符都是可微的，因此端到端的网络训练成为可能。由于简化了 workflow，网络可以在极低运算开销下进行验证，达到实时速度。



3.OCR领域现存问题及优化

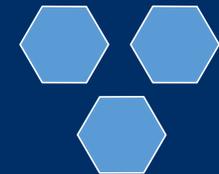




Ocr 此前遇到的问题与解决方法

由于人类识字的机理及过程并不清楚，汉字识别的研究还只能停留在一般模式识别问题的研究上。1966年IBM公司的Casey和Nagy首次提出了一个识别1000汉字的识别方案。至今，中文OCR技术取得了令人瞩目的重大进展。其研究工作，大致分为以下三个大的阶段。

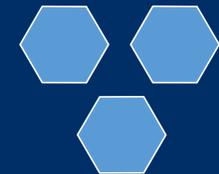




第一阶段

从70年代末到80年代末期。这一阶段主要研究的是汉字识别的算法和方案探索。

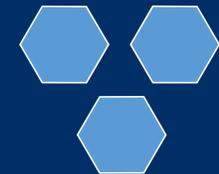




第二阶段

从90年代初期开始到90年代中期, 汉字OCR技术进入了一个更为关键和重要的阶段—将实验室的研究成果推向市场。

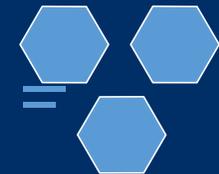




第三阶段

从90年代中期开始, 致力于中文OCR技术和中文OCR系统性能的提高, 特别是汉字识别率的提高和稳健性的增强。



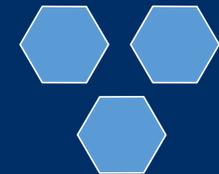


现阶段面临的问题

在OCR技术的研究发展上,特别是汉字OCR技术的真正实用化和被广大用户方便使用方面,还有许多问题需要进一步解决,主要包括三个方面。

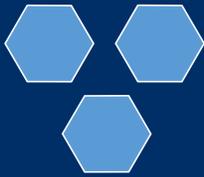
- 一. 系统总体性能的进一步提高。
- 二. 系统固化以及系统各部分的质量和性能的稳定提高。
- 三. 扩大OCR核心技术的应用范围。





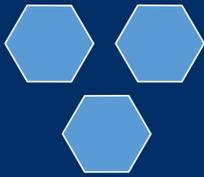
一种基于文字位置信息的 OCR优化方法





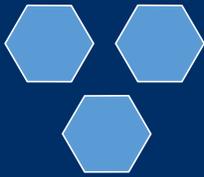
1. 设定需要提取的约束信息项及提取范围；
2. 基于文字位置信息的OCR优化方法，对待识别图像进行预处理





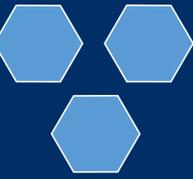
3. 设定需要提取的约束信息项及提取范围，
4. 通过临近连通区域搜索方法，得到图片所有文字及其坐标位置信息。





5. 计算所有文字信息的平均行间距
6. 基于所述间距与平均行间距的大小，判断当前文字与下一行文字是否需要合并。





7. 提取每一个待识别的约束信息项，对提取到的信息项文字数据进行格式化输出。
8. 模板设计模块，用于设定需要提取的约束信息项及提取范围。





4.OCR未来发展方向

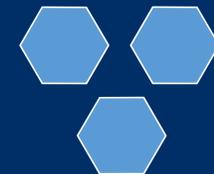
OCR技术发展三大趋势

一体化的端到端OCR模型
兼具高性能高效率的OCR
从感知到认知的智能OCR



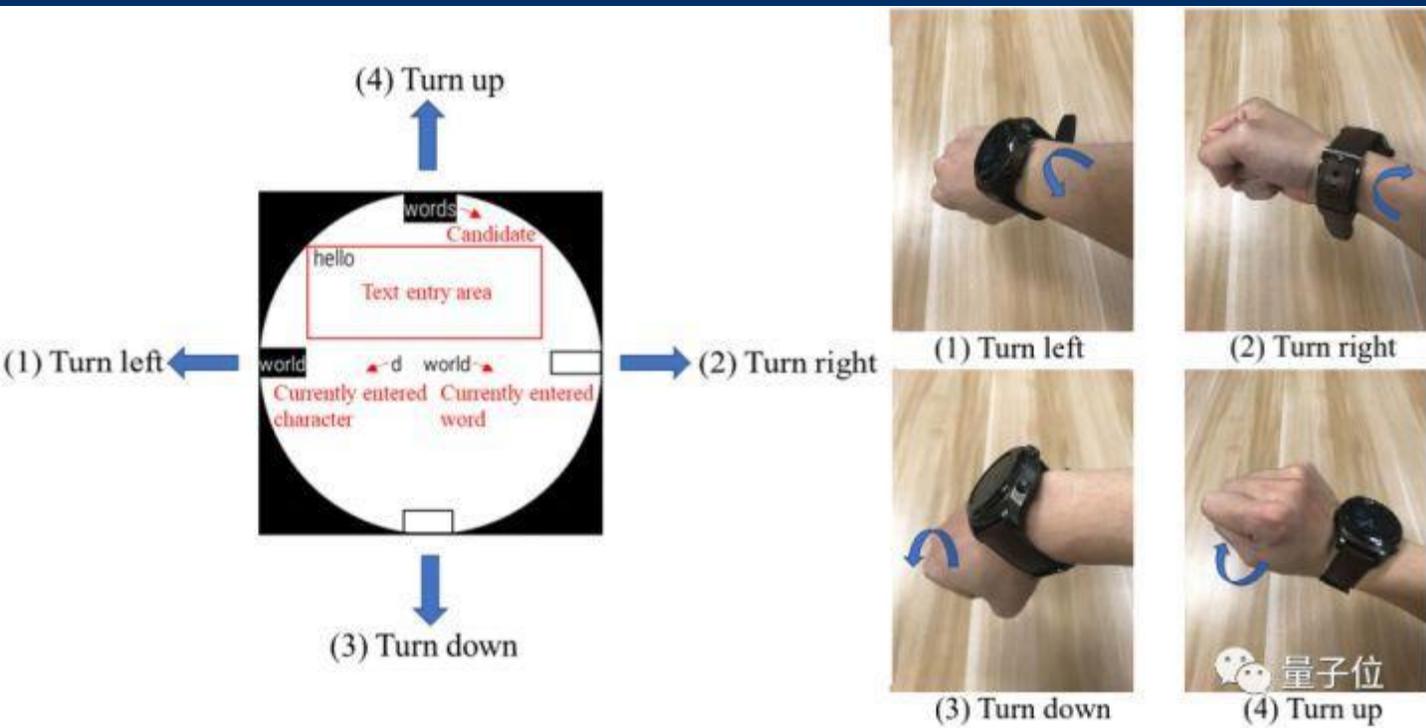
在美国权威机构GRAND VIEW RESEARCH发布的《全球OCR (Optical Character Recognition) 市场预测以及趋势分析》中指出2018年全球OCR市场规模已经达到了52.7亿美元，并预测全球OCR市场将以13.7%的复合年增长稳健发展，至2025年全球OCR市场规模将达到133.81亿美元。





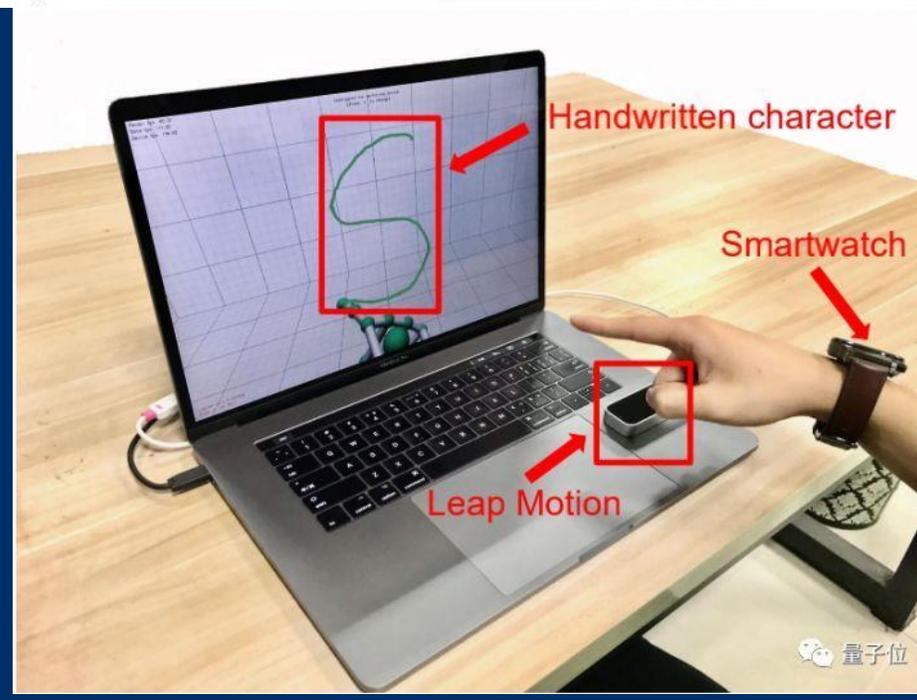
OCR产业生态图

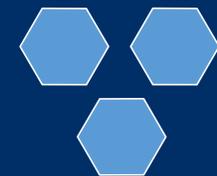




character	a	A	b	B
scenario				
on-surface (non-rotating-wrist)				
on-surface (rotating-wrist)				
in the air (non-rotating-wrist)				
in the air (rotating-wrist)				

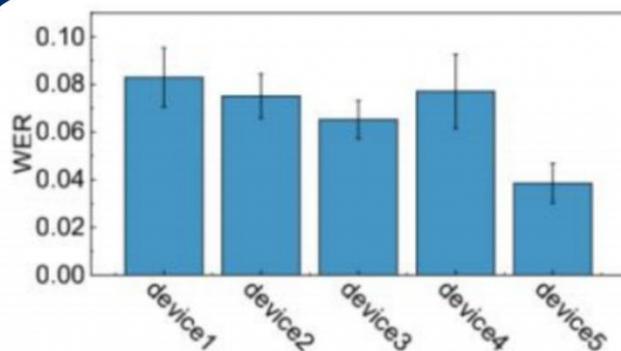
空气输入法！
浙大最新研究：空中动动手指就能打字



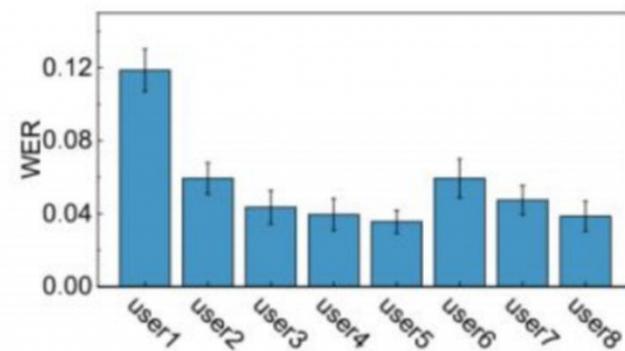


词错误率

- 衡量指标为WER（词错误率，Word Error Rate），计算方法就是将错误字/识别正确的字。
- 其中错误字涉及三种类型：漏字（用I表示，即校对成正确的拼写时需要再插入的字的数量）、多字（用D表示，即需要删除的字数）、错字（用S表示，即需要替换的字数）。
- 左图中，同一个用户使用5种不同智能手表测试AirText获得的准确率得分分别为：8.3%、7.5%、6.5%、7.7%和3.9%。
- 右图中，8位不同用户使用同一手表获得了11.2%、5.9%、4.3%、4.0%、3.6%、5.9%、4.7%和3.9%的WER。
- 从中我们可以看出，与不同的设备相比，不同的用户对准确性的影响更大。

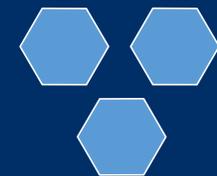


(a) one-user multi-devices

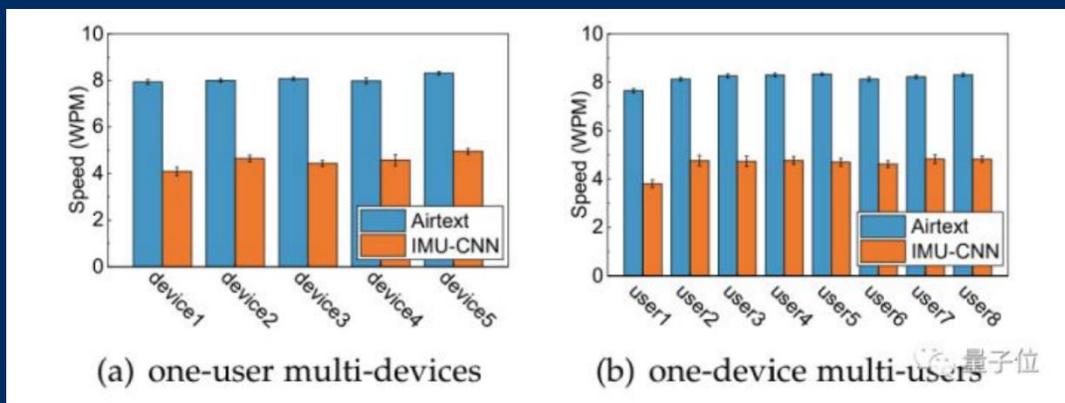


(b) one-device multi-users





速度测试



- 由于BLSTM基线的WER约为57%，错误率太高，研究人员只比较了AirText和IMU-CNN两者的速度。
- 衡量指标是WPM（单词/每分钟，Word Per Minute），其计算方式为用总体输入字数-错误字数/时间。
- 结果AirText的平均WPM为8.1，而IMU-CNN基线的WPM仅为4.6。
- 研究人员指出，此输入速度与一些基于双手触摸屏的文本输入方法相当（这些方法在实际应用中的WPM为9.1、9.8WPM）。
- 总体来看，AirText的准确率不错，但速度还需要进步。
- 慢的主要原因还是因为它每拼写一个字符就需停顿一会儿。
- 研究人员正在研究破解办法，并表示最终有兴趣将AirText商业化。





5.Demo展示

简介

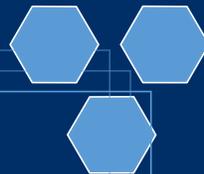
Python-tesseract是一款用于光学字符识别（OCR）的python工具，即从图片中识别出其中嵌入的文字。Python-tesseract是对Google Tesseract-OCR的一层封装。它也同时可以单独作为对tesseract引擎的调用脚本，支持使用PIL库（Python Imaging Library）读取的各种图片文件类型，包括jpeg、png、gif、bmp、tiff和其他格式，。作为脚本使用它将打印出识别出的文字而非写入到文件。所以安装pytesseract前要先安装PIL和tesseract-ocr这两依赖库





环境要求

- Python 3.6+
- PIL库
- 安装Google Tesseract OCR
- 系统：windows/mac/linux，我的系统是Windows10



demo展示

tochr2.py > ...

```
1  import pytesseract
2  from PIL import Image
3  text = pytesseract.image_to_string(Image.open(r"C:\Users\LENOVO\Desktop\py\OIP-C (1).jfif"))
4  print(text)
5
```



The Art of War

Sunzi

 China Intercontinental Press

```
PS C:\Users\LENOVO\Desktop\py> c  
her' '50643' '--' 'c:\Users\LENOV  
} China Intercontinental Press.
```

```
PS C:\Users\LENOVO\Desktop\py> 
```



```
1 import pytesseract
2 from PIL import Image
3 text = pytesseract.image_to_string(Image.open(r'C:\Users\LENOVO\Desktop\py\OIP-C (1).png'))
4 print(text)
```

```
PS C:\Users\LENOVO\Desktop\py> c:; cd 'c:\Users\LENOVO\Desktop\py'; & 'D:\python\python.exe' 'C:\Users\LENOVO\.vscode\extensions\ms-python.python-2022.2.1924087327\pythonFiles\lib\python\debugpy\launcher' -' 'c:\Users\LENOVO\Desktop\py\tocr2.py'
```

RUN

```
import pytesseract
from PIL import Image
```

```
text = pytesseract.image_to_string(Image.open(r'C: \Users\LENOVO\Desktop\py\OIP-C (1).png"))
print(text)
```

```
PS C:\Users\LENOVO\Desktop\py> █
```

PaddleOCR

- PaddleOCR是百度开源的一款基于深度学习的ocr识别库，对中文的识别精度相当不错，可以应付绝大多数的文字提取需求。

✦ paddleocr.py > ...

```
1  ocr = PaddleOCR(use_angle_cls=True, lang="ch")
2
3  img_path = r'C:\Users\LENOVO\Desktop\py\C%ZHP$07EKVRZ]LU6}0QB`E.png'
4
5  result = ocr.ocr(img_path, cls=True)
6  for line in result:
7      print(line)
8
9  from PIL import Image
10 image = Image.open(img_path).convert('RGB')
11 boxes = [line[0] for line in result]
12 txts = [line[1][0] for line in result]
13 scores = [line[1][1] for line in result]
14 im_show = draw_ocr(image, boxes, txts, scores)
15 im_show = Image.fromarray(im_show)
16 im_show.show()
17
```

新建PPTX 演示文稿.pptx

```
PS C:\Users\LENOVO\Desktop\py> c:; cd 'c:\Users\LENOVO\Desktop\py\paddleocr.py'
```

新建PPTX 演示文稿.pptx

```
PS C:\Users\LENOVO\Desktop\py> █
```



```
PS C:\Users\LENOVO\Desktop\py> c:; cd 'c:\Users\LENOVO\Desktop\py\paddleocr.py'
```

```
Kurfüstendamm
```

```
11-17
```

```
PS C:\Users\LENOVO\Desktop\py> █
```

dddocr

- dddd_ocr也是一个用于识别验证码的开源库，识别效果也是非常不错，对一些常规的数字、字母验证码识别有奇效。



```
D:\py38\python.exe C:/Users/75061/Desktop/paddleocr.py
```

欢迎使用dddocr，本项目专注带动行业内卷，个人博客：wenanzhe.com

训练数据支持来源于：<http://146.56.204.113:19199/preview>

爬虫框架feapder可快速一键接入，快速开启爬虫之旅：<https://github.com/Boris-code/feapder>

```
3n3d
```

```
Process finished with exit code 0
```



```
D:\py38\python.exe C:/Users/75061/Desktop/paddleocr.py  
欢迎使用ddddocr，本项目专注带动行业内卷，个人博客:wenanzhe.com  
训练数据支持来源于:http://146.56.204.113:19199/preview  
爬虫框架feapder可快速一键接入，快速开启爬虫之旅: https://github.com/BinaryHuk  
Huk
```

```
Process finished with exit code 0
```

感谢您的聆听与观看

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Lorem ipsum dolor sit amet, consectetur