

# OCR

## 前言

HOG 特征：全称 Histogram of Oriented Gradient（方向梯度直方图），由图像的局部区域梯度方向直方图构成特征；LBP 特征：全称 Local Binary Pattern（局部二值模式），通过比较中心与邻域像素灰度值构成图像局部纹理特征；Haar 特征：描述图像的灰度变化，由各模块的像素差值构成特征；核函数（Kernels）：从低维空间到高维空间的映射，把低维空间中线性不可分的两类点变成线性可分的；

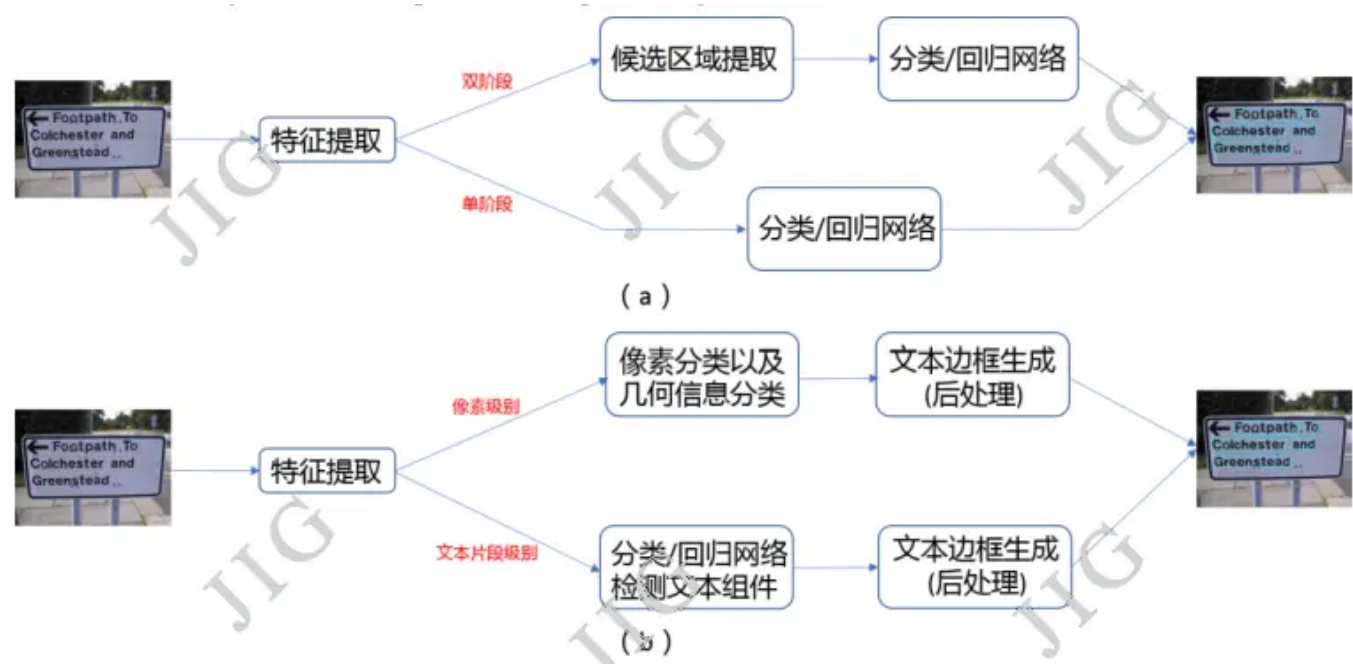
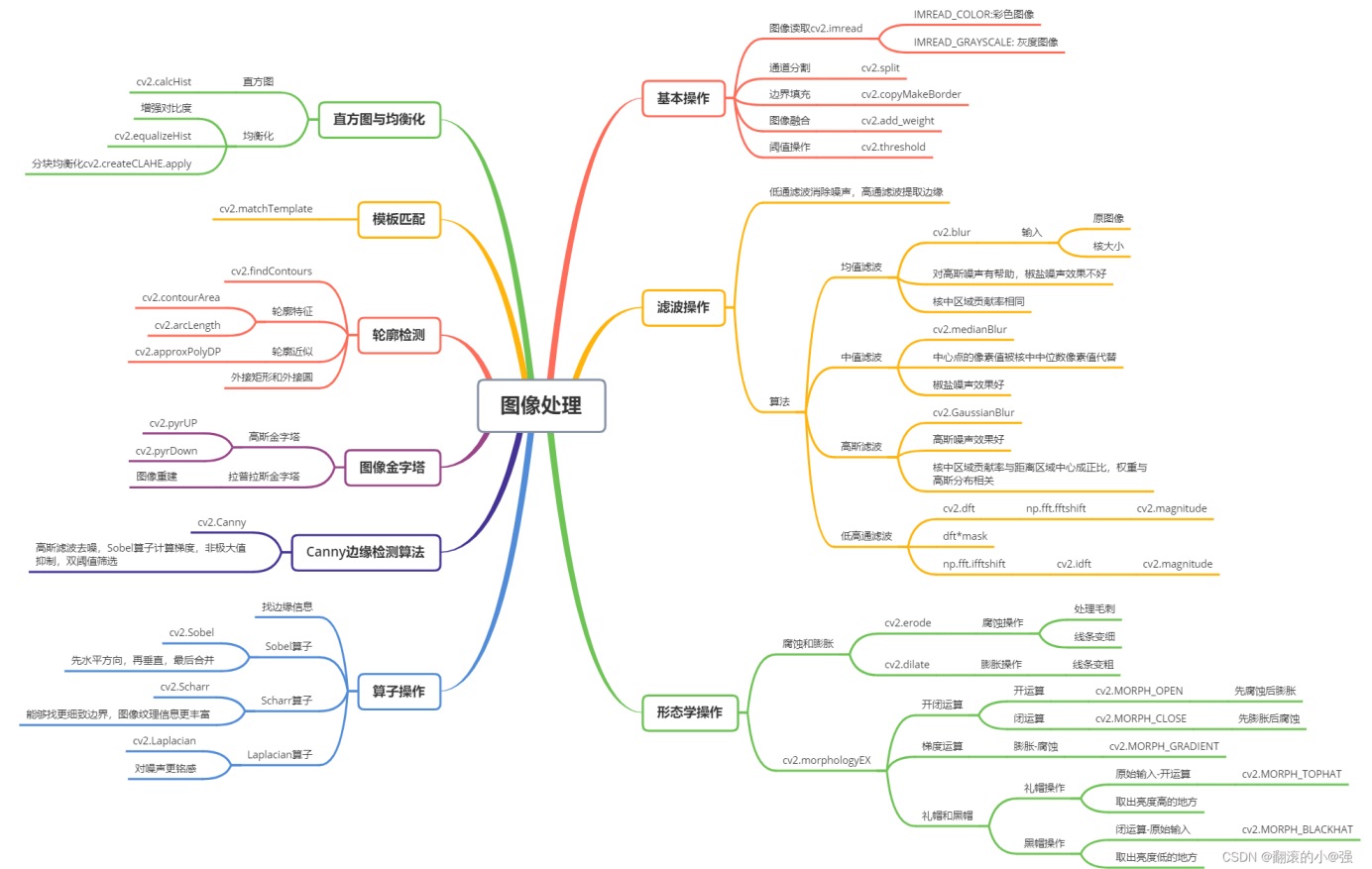


图 2 自然场景文本检测方法流程图。(a) 自顶向下；(b) 自底向上

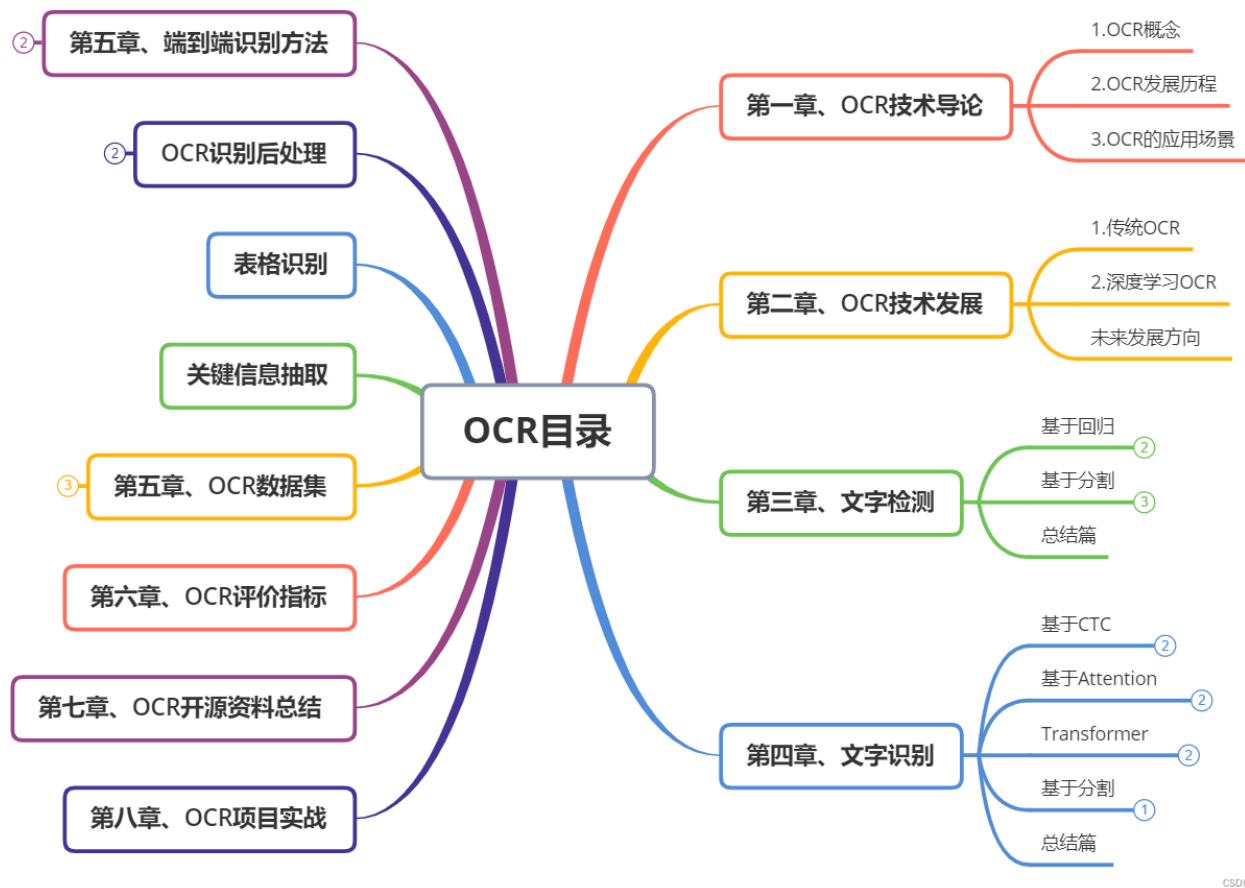
## 图像处理导图



Opencv 思维导图



目录



# OCR 入门教程系列（一）：OCR 基础导论

## OCR 发展



文档文字识别：可以将图书馆、报社、博物馆、档案馆等的纸质版图书、报纸、杂志、历史文献档案资料等进行电子化管理，实现精准地保存文献资料。

自然场景文字识别：识别自然场景图像中的文字信息如车牌、广告干词、路牌等信息。对车辆进行识别可以实现停车场收费管理、交通流量控制指标测量、车辆定位、防盗、高速公路超速自动化监管等功能。

票据文字识别：可以对增值税发票、报销单、车票等不同格式的票据进行文字识别，可以避免财务人员手动输入大量票据信息，如今已广泛应用于财务管理、银行、金融等众多领域。。

证件识别：可以快速识别身份证、银行卡、驾驶证等卡证类信息，将证件文字信息直接转换为可编辑文本，可以大大提高工作效率、减少人工成本、还可以实时进行相关人员的身份核验，以便安全管理。

OCR 生态



OCR 的技术路线

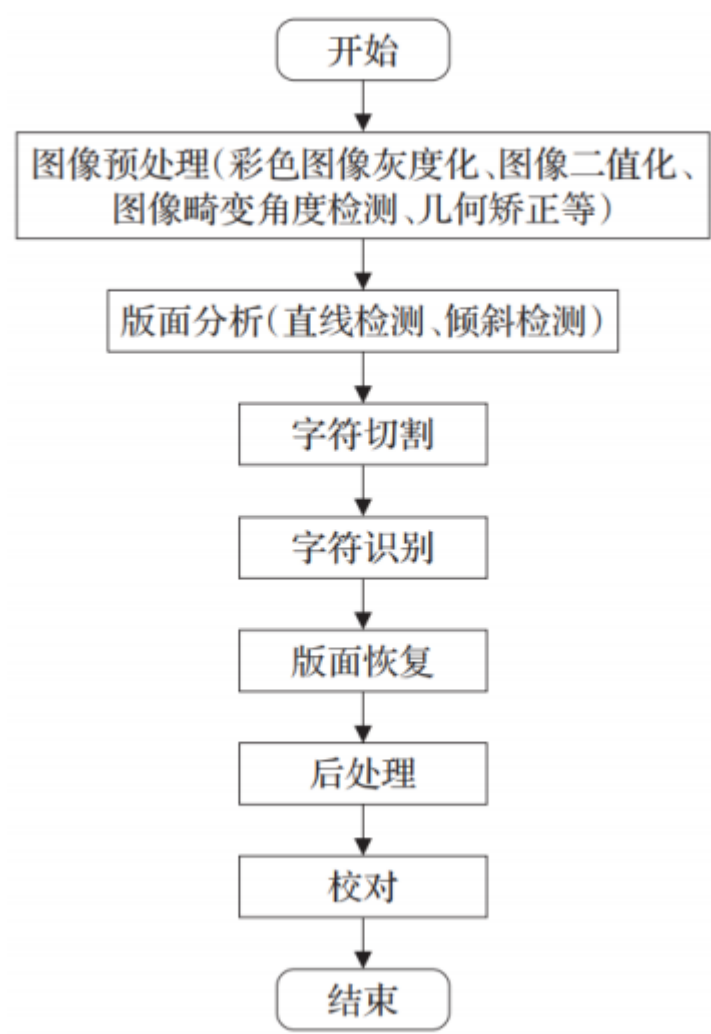


其中 OCR 识别的关键路径在于文字检测和文本识别部分，这也是深度学习技术可以充分发挥功效的地方。

传统 OCR 技术流程

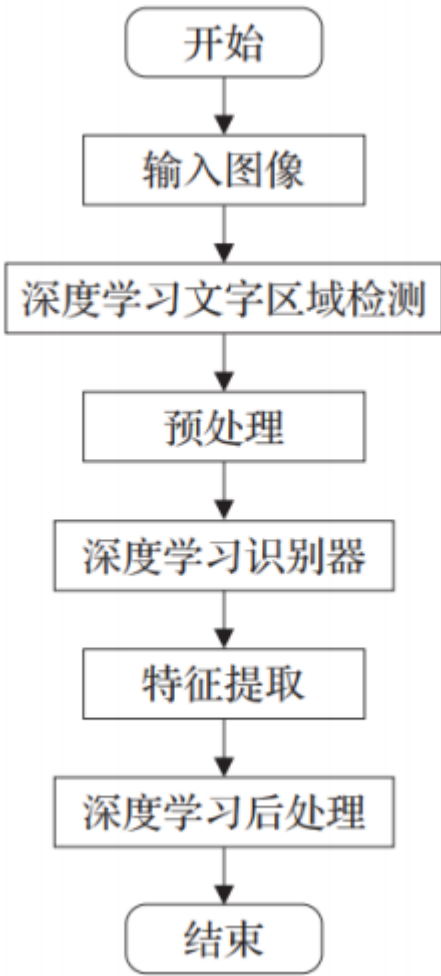
- 1、水平投影垂直投影 2、模板匹配 3、查找轮廓 findcontours

传统的光学字符识别过程为：图像预处理（彩色图像灰度化、二值化处理、图像变化角度检测、矫正处理等）、版面划分（直线检测、倾斜检测）、字符定位切分、字符识别、版面恢复、后处理、校对等。

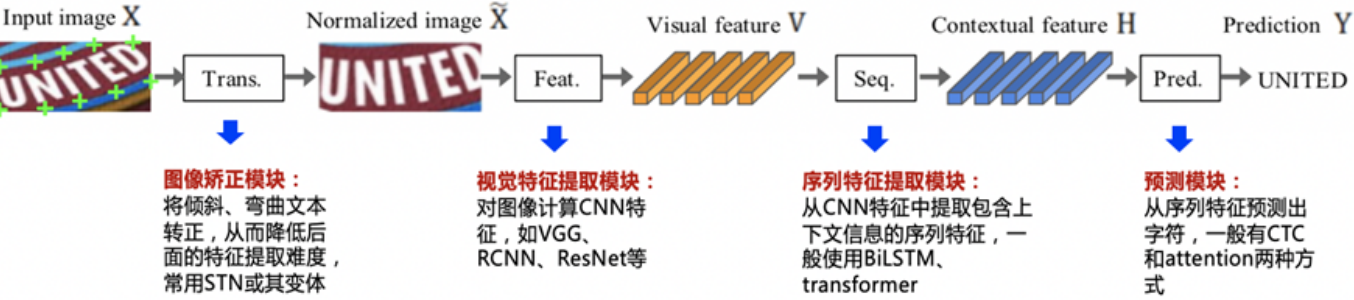


深度学习 OCR 技术流程

深度学习图像文字识别流程包括：输入图像、深度学习文字区域检测、预处理、特征提取、深度学习识别器、



深度学习后处理等。  
现有多数深度学习识别算法具体流程包括图像校正、特征提取、序列预测等模块，流程如图所示：



常用文字检测框架介绍

DBNet ( 早期 14-16)

首先，DB 是一种基于分割的文本检测算法。在各种文本检测算法中，基于分割的检测算法可以更好地处理弯曲等不规则形状文本，因此往往能取得更好的检测效果。但分割法后处理步骤中将分割结果转化为检测框的流程复杂，耗时严重。因此作者提出一个可微的二值化模块 ( Differentiable Binarization，简称 DB )，将二值化阈值加入训练中学习，可以获得更准确的检测边界，从而简化后处理流程。DB 算法最终在 5 个数据集上达到了 state-of-art 的效果和性能。

- 主要思想：先获取图像中的文本区域，再利用 **opencv、polygon** 等后处理得到文本区域的最小包围曲线；



- DB 提出可微分阈值，通过一个近似于阶跃函数的二值化函数使得分割网络在训练时学习文本分割的动态阈值，使模型提升精度，简化后处理；
- DB 的 backbone 是典型的 FCN 结构，由多层上采样和下采样的特征图 concat 完成。

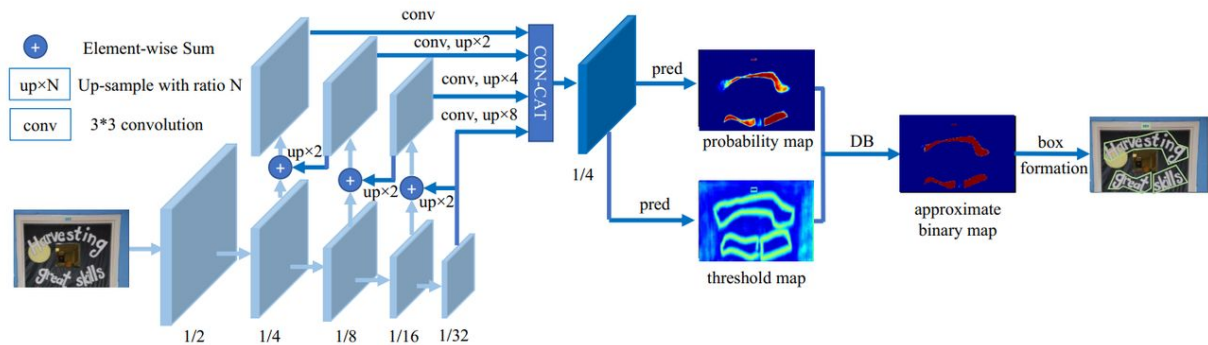
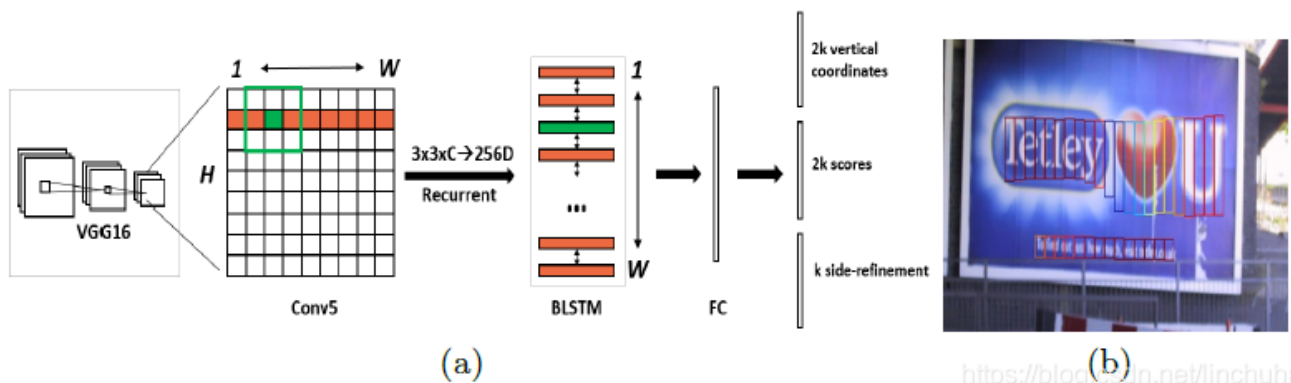


Figure 3: Architecture of our proposed method, where “pred” consists of a  $3 \times 3$  convolutional operator and two de-convolutional operators with stride 2. The “1/2”, “1/4”, ... and “1/32” indicate the scale ratio compared to the input image.

### CTPN (中期 17-18)

CTPN 模型主要包括三个部分，分别是卷积层、Bi-LSTM 层、全连接层，其结构如下图所示。



在卷积层部分，CTPN 选取 **VGG16** 模型前 5 个 **conv5= stage** 得到 **feature maps** 作为图像最后的特征，假设此时 **feature maps** 的尺寸为  $W \times H \times C$ ；由于文本之间存在序列关系，因此，作者引入了递归神经网络，采用的是一层 Bi-LSTM 层，作者发现引入了递归神经网络对文本检测的效果有一个很大的提升，如下图所示，第一行是不采用递归神经网络的效果，第二行是采用了 Bi-LSTM 后的效果。具体的做法是采用一个的滑动窗口，提取 **feature maps** 上每个点周围的区域作为该点的特征向量表示，此时，图像的尺度变为，然后将每一行作为序列的长度，高度作为 **batch\_size**，传入一个 128 维的 Bi-LSTM，得到 Bi-LSTM 层的输出为；将 Bi-LSTM 的输出接入全连接层，在这一部分，作者引入了 **anchor** 的机制，即对每一个点用  $k$  个 **anchor** 进行预测，每个 **anchor** 就是一个盒子，其高度由  $[273, 390, \dots, 11]$  逐渐递减，每次除以 0.7，总共有 10 个。作者采用的是三个全连接层分支。

### OCR 评价指标

OCR 评价指标包括字段粒度和字符粒度的识别效果评价指标。

以字段为单位的统计和分析，适用于卡证类、票据类等结构化程度较高的 OCR 应用评测。以字符（文字和标点符号）为单位的统计和分析，适用于通用印刷体、手写体类非结构化数据的 OCR 应用评测。具体指标包括以下几个：



**a 字段召回率**，指被完全正确识别字段（测试输出结果与字段的所有字符完全匹配）数量与总字段数比值。

**b 字段准确率**，指被完全正确识别字段（测试输出结果与字段的所有字符完全匹配）数量与测试返回识别结果的字段数量比值。

**c 字符召回率**，指被完全正确识别字符数量与真实字符总数的比值，可以反应识别错和漏识别的情况。

**d 字符准确率**，指被完全正确识别字符数量与测试返回的字符数的比值，可以反应识别错和多识别的情况。

**e  $F_\beta$  - Score**，可以综合反映字符识别召回效果和字符识别准确效果，计算公式如下：

$$F_\beta - Score = (1 + \beta^2) * \frac{\text{字符召回率} * \text{字符准确率}}{\beta^2 * (\text{字符召回率} + \text{字符准确率})}$$

**f 最小编辑距离**，表示测试结果要与标注结果一致需要修改的字符数，忽略引擎返回行的顺序与原图标注的顺序，适用于电商广告、手机截图等样本版式不规范的 OCR 应用评测。

**g 全图编辑距离**，表示测试返回结果要与标注结果一致需要修改的字符数，要求引擎返回的每一行文字顺序要和标注顺序一致，适用于文档、表格等样本版式较为规范的 OCR 应用评测。

编辑距离：

编辑距离是针对二个字符串（例如英文字）的差异程度的量化量测，量测方式是看至少需要多少次的处理才能将一个字符串变成另一个字符串。在莱文斯坦距离中，可以删除、加入、替换字符串中的任何一个字元，也是较常用的编辑距离定义，常常提到编辑距离时，指的就是莱文斯坦距离。

测试指标说明：

平均识别率： $[1 - (\text{编辑距离} / \max(1, \text{groundtruth 字符数}, \text{predict 字符数}))] * 100.0\%$  的平均值；平均编辑距离：编辑距离，用来评估整体的检测和识别模型；平均替换错误：编辑距离计算时的替换操作，用于评估识别模型对相似字符的区分能力；平均多字错误：编辑距离计算时的删除操作，用来评估检测模型的误检和识别模

型的多字错误；平均漏字错误：编辑距离计算时的插入操作，用来评估检测模型的漏检和识别模型的少字错误；