



# 基于深度学习的古籍文字检测与识别算法

## (古籍文字识别技术报告)

组名：彼采AI兮

成员：郭文君，梁兆翔，方草青，徐伯辰

# 团队简介——彼采AI兮

团队来自北京理工大学自动化学院，均为研二在读。



郭文君

分工

技术报告撰写  
算法开发  
数据集采集与分析



梁兆翔

分工

OCR识别  
整体算法建立  
代码运行文档



徐伯辰

分工

数据集制作



芳草青

分工

数据集制作

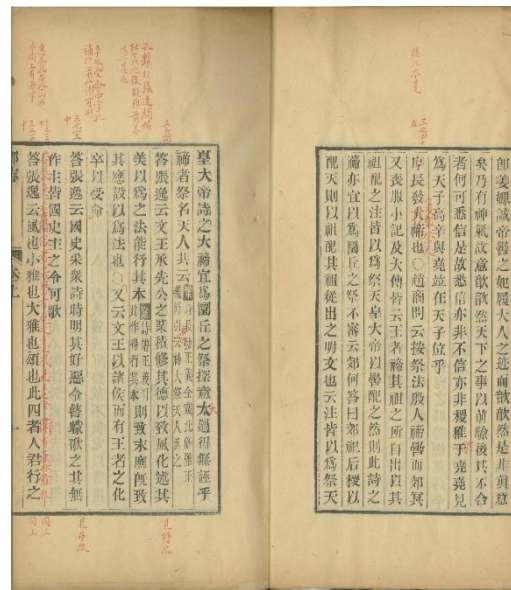
## 赛题重述

对跨越不同朝代、不同字体、不同存储格式、不同旋转方向的古籍电子图片进行检测、识别与语言逻辑恢复。

尽可能完整地提供图像信息，可包括但不限于正文文本、双行小字、朱批、文字坐标等古籍图像中所包含的信息。

## 问题理解

寻求以OCR算法为核心，对大规模复杂古籍数据的鲁棒解决方案。



p0000017.JPG\_1第1张图片识别结果:  
九劉燭雲校  
節髮緝誠帝嘗之妃履大人之迹而欲飲然非真意  
矣乃有口氣故意飲然天下之事以前驗後其不合  
者何可悉信是故悉信亦非不信亦非稷雅乎堯見  
為天子高辛與堯竝在天子位乎  
序長發跳禘也○趙商問云按祭法殷人禘嘗而郊冥  
又喪服小記及頭傳皆云王者禘其社之所自出以其  
[且/祖]配之注皆以為祭天皇帝以魯配之然射此詩之  
禘亦宜以為圓丘之祭不審云郊何答曰郊禘后稷以  
配則以祖配其祖從出之以文也云注皆以為祭天  
皇大帝詩之大禘宜為圓丘之祭探意太過得無誣乎  
禘者祭名天人共云[途長發正義全載止條繼正/義節云禘大祭天人異之]  
答張逸云文王承先公之業積修其德以致風化述其  
美以為之法能行其本[案詩譜正義訂/此作得行其本]則致未應既致  
其應設以為法也○又云文王以謂侯而有王者之化  
卒以履命  
答張逸云國史采衆時明其好惡命誓謗歌之其無  
作主皆國史主之令可歌  
答張逸云風也小雅也犬雅也頌也此四者人行之  
下黍

## 解决方案流程

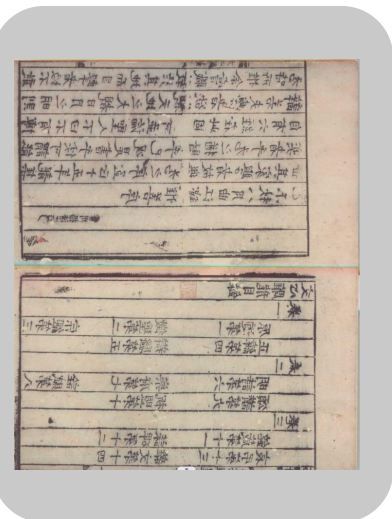
输入：古籍图片

数据预处理

OCR模型识别

原文语言逻辑恢复

输出：古籍全文txt



格式统一与方向统一

数据增强

基于Paddle OCR的  
原始模型

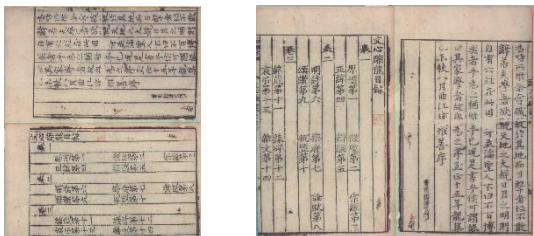
基于数据中高频古字  
词的模型微调

基于坐标，对**碎片化的  
OCR结果**进行鲁棒重组  
与原文语言结构恢复

0000017.JPG\_1第1张图片识别结果：  
九劉翔雲校  
節髮頻誠帝儀之紀履大人之迹而欲欲然非真意  
矣乃有口氣故寧敢欲然天下之事以前驗後其不合  
者何可悉信是故悉信亦非不信亦非稷推乎義見  
為天子高辛與禹在天子位乎  
序長發跳禘也○趙商問云按祭法般人禘而郊冥  
又喪服小記及頭傳皆云王者禘其社之所自出以其  
[且/祖]配之注皆以為祭天皇大帝以饗配之然此詩之  
禘亦宜以為禹丘之祭不審云郊何答曰郊禘後稷以  
配則以祖配其祖從出之以文也云注皆以為祭天  
皇大帝詩之大禘宜為禹丘之祭探意太過得無誤乎  
禘者祭名天人共云[途長發正義全載止條雖正/義節云禘大祭天人異之]  
答張逸云文王承先公之業積修其德以致風化述其  
美以屬之法能行其本[案詩譜正義訂]此作得行其本則致未應既致  
其應設以為法也○又云文王以讓侯而有王者之化  
卒以康命  
答張逸云國史采衆事明其好惡命醫禘歌之其無  
作主皆國史主之令可歌  
答張逸云風也小雅也大雅也頌也此四者人君行之  
下萎

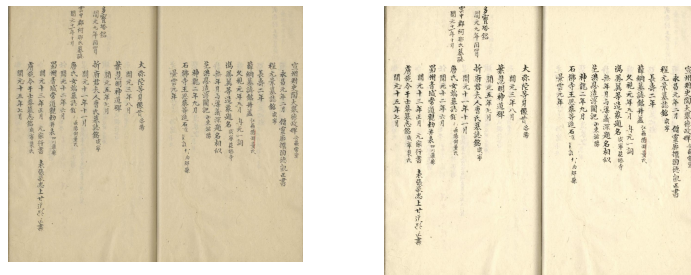
之后的解决方案PPT将按照该流程展开。

## ① 格式与旋转方向统一



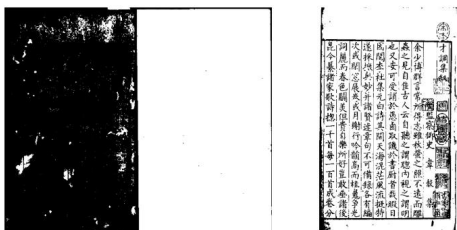
将图片统一转化为JPG格式  
检测图片尺寸长宽比，自动将  
图片统一旋转至正常阅读方向

## ② 对比度增强



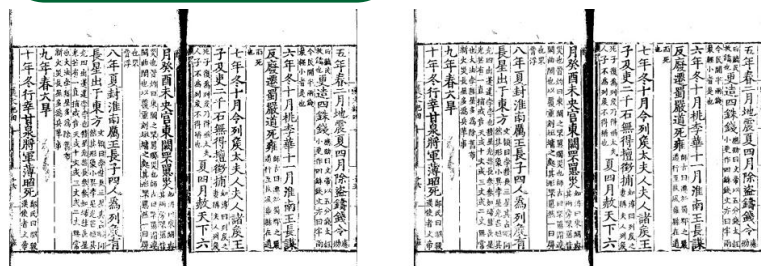
突出文本轮廓边界，提高识别  
准确性

## ③ 二值化

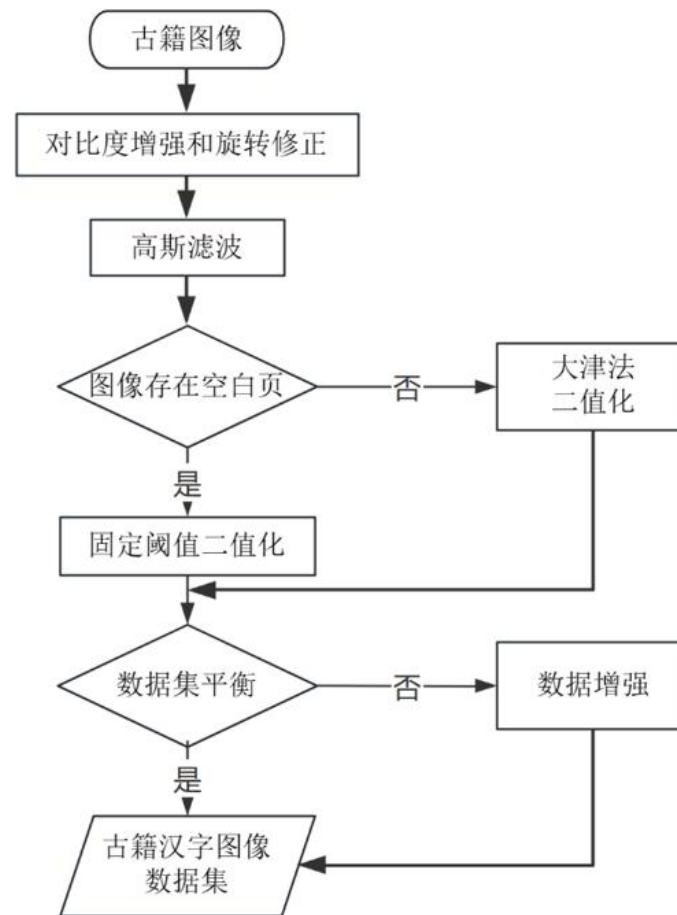


根据图像是否存在空白页，  
使用大津法或固定阈值二值化

## ④ 高斯滤波



选取合适的高斯核大小，有助于平滑图像  
中的瑕疵，比如折痕和褶皱



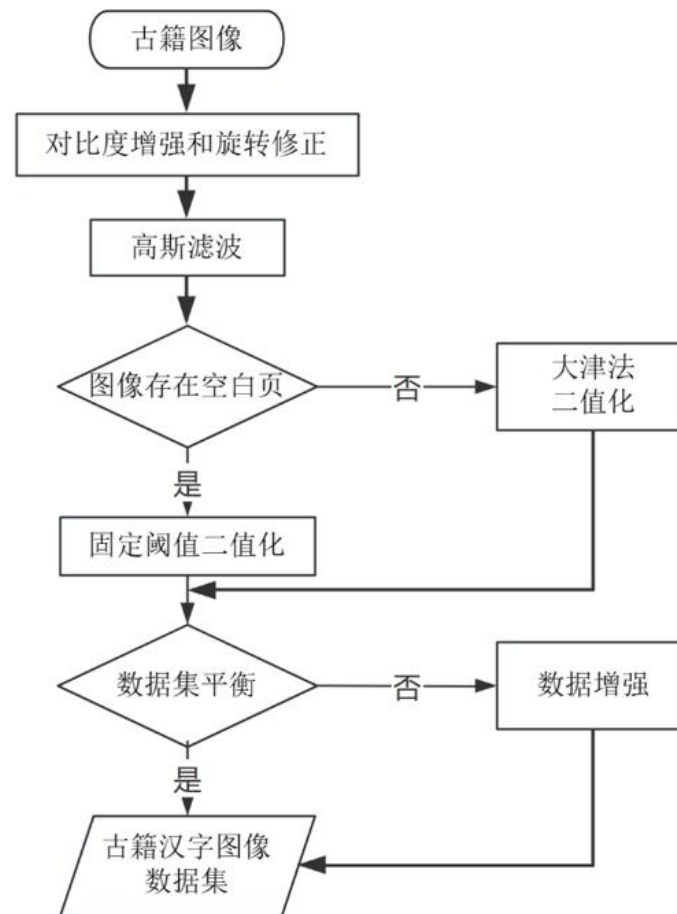




古籍原始图像（左）、灰度图像（中）和二值化图像（右）



古籍原始图像和数据预处理后的图像



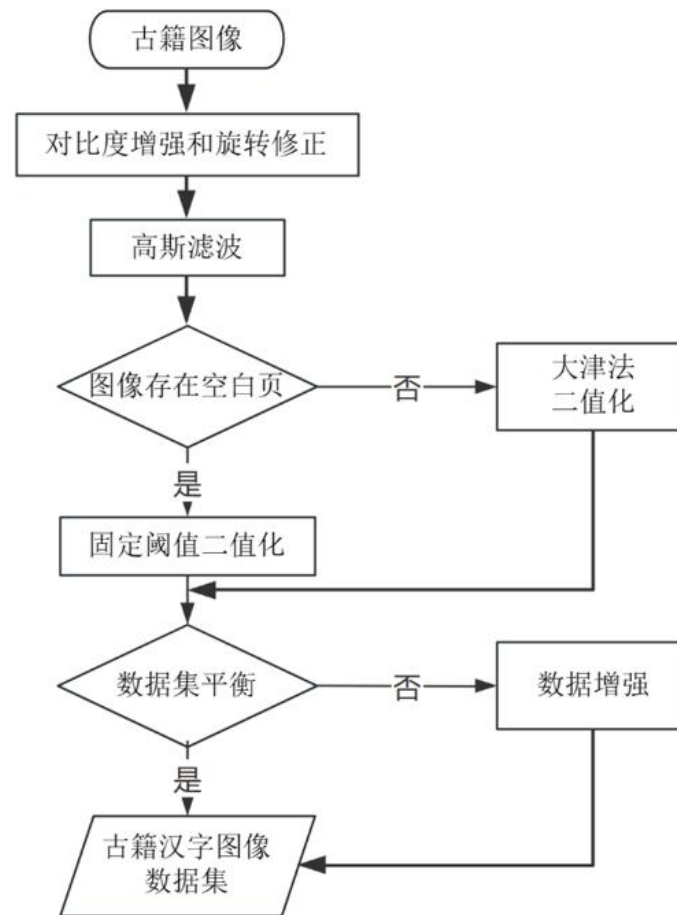
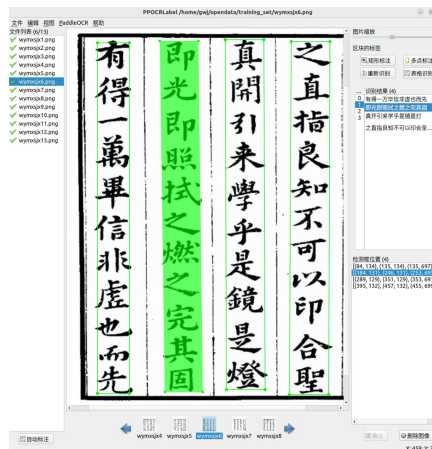
## ① 数据集分析

数据集中存在数十万张原始图片，为了得到**更高的性价比**，使我们能够以更小的标注工作量获得更好的识别结果，我们需要挑选出数据集中具有**代表性**的字词进行额外标注。

根据《古籍汉字字频统计》以及数据集中**朝代分布特点**，我们选择出了**低频分量**（即出现频次高）的字词作为代表，从古代字帖、书籍和石碑中提取了低频字样本加入到数据集中，以**减轻数据不平衡**问题。

## ② 数据集标注

我们使用百度飞桨推出的半自动文字标注工具**PPLabel**进行数据集标注。这帮助我们获得更高的标注效率。



## ① OCR base模型选择

经过对比与调研，我们选择百度飞桨 **PaddleOCR模型**作为流程中的**base算法**。PaddleOCR支持中英文数字组合识别、竖排文本识别、长文本识别。同时支持多种文本检测、文本识别的训练算法。我们使用飞桨提供的高精度预训练权重模型作为基础base。

由于base算法并不主要面对古籍模型，在古体字上的检测与识别精度并不理想。我们需要对其进行**调整**。



性别: 男	住院号:		
年龄: 34岁	科室: 内科		
代号	项目	结果	参考值
GLU	葡萄糖	6.271	3.9--6.11 mmol/L
TG	甘油三酯	2.121	0.4--1.88mmol/L
TCHO1	血清总胆固醇	6.361	2.8--5.7mmol/L
HDL	高密度胆固醇	1.80	0.83--1.91mmol/L
LDL	低密度胆固醇	4.32	0--3.15mmol/L
APOA1	载脂蛋白A1	1.28	1--1.5g/L
APOB	载脂蛋白B	0.93	0.6--1.14g/L

## ② 模型调整

我们使用的高精度的预训练权重模型所使用的文本数据库非常庞大，我们的自定义数据集对来说占比非常小。因此，在训练过程中我们非常警惕模型的**过拟合情况**。为此，我们主要采取了**2个措施**：

- 1.在自定义数据集中加入大量的以预训练模型识别结果的的数据为标注的**非古籍数据**。
- 2.我们控制训练轮数在100个epoch以内寻找表现较好的训练结果，以防止在毁坏原生模型的情况下，对其**微调**以在古籍表现上更好。



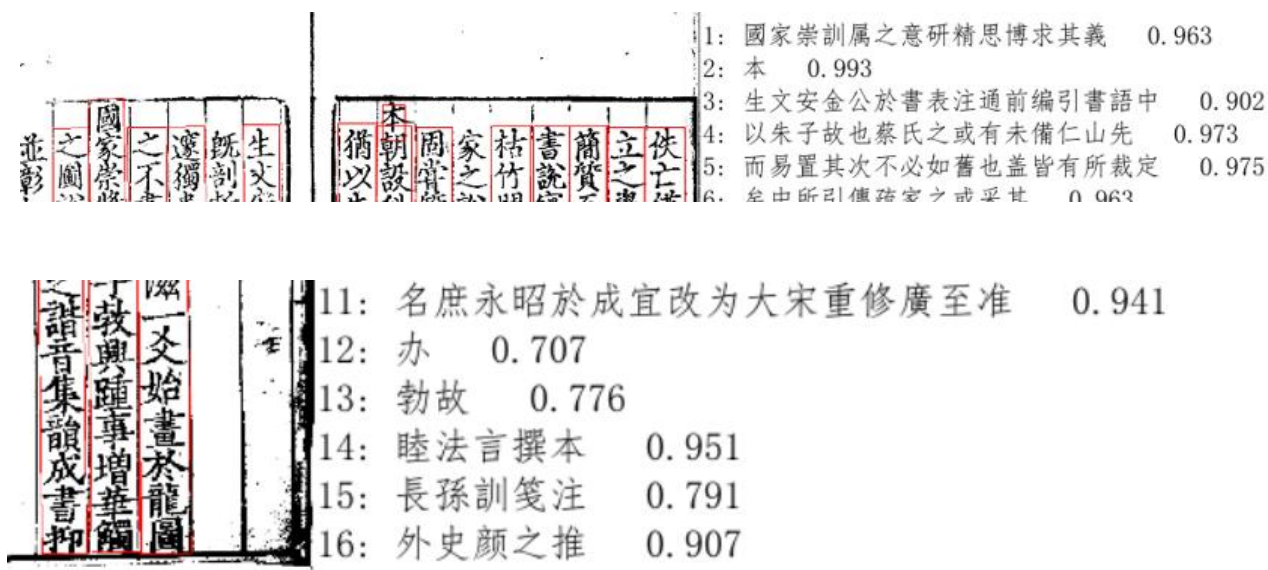
性别: 男	住院号:		
年龄: 34岁	科室: 内科		
代号	项目	结果	参考值
GLU	葡萄糖	6.271	3.9--6.11 mmol/L
TG	甘油三酯	2.121	0.4--1.88 mmol/L
TCHO1	血清总胆固醇	6.361	2.8--5.7 mmol/L
HDL	高密度胆固醇	1.80	0.83--1.91 mmol/L
LDL	低密度胆固醇	4.32	0--3.15 mmol/L
APOA1	载脂蛋白A1	1.28	1--1.5 g/L
APOB	载脂蛋白B	0.93	0.6--1.14 g/L



## ① OCR微调模型识别

使用完成调整后的模型对经过预处理的数据集进行OCR检测与识别，输出识别到的文字段与对应矩形检测框坐标信息。

然而，由于存在部分文字检测不到，识别结果是**碎片化**的若干个**文字段**。并且由于古籍的语言顺序是**右上至左下**，识别结果**顺序**也与期望不符。同时，有时还会有部分识别错误，在古籍中识别结果为英文、阿拉伯数字、罗马数字。



使用完成调整后的模型对经过预处理的数据集的识别结果为左图，可以看到，由于“国家”、“本朝”在**偏上**位置，使得默认的输出结果将其错误的放在了文首。

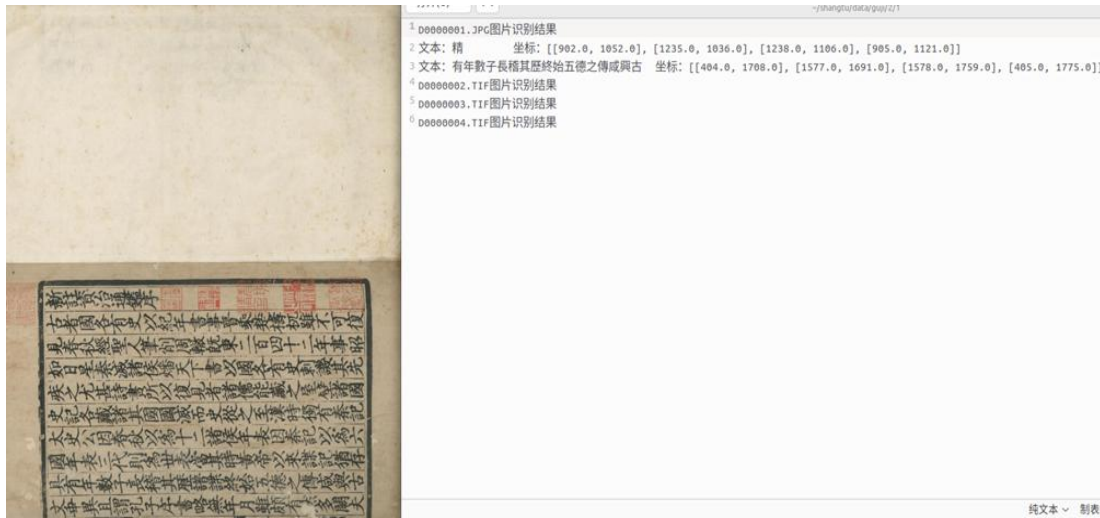
同时，由于个别字检测难度高，会导致检测框的碎片化，而**框与框之间的高低关系**也将导致错误的语言顺序。

## ② OCR识别结果重组

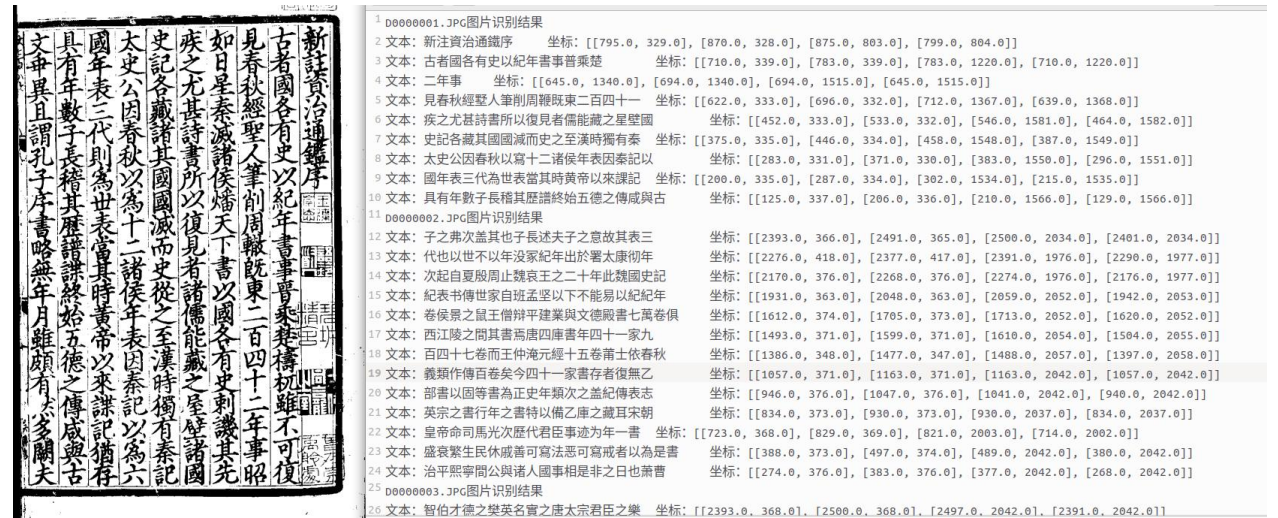
首先，我们的算法将识别结果中的明显错误，如：数字、英文、特殊字符删除。

我们估计数据集中每一列平均宽度的像素值，依据检测结果中不同文字段的**检测框位置信息**，按照先从右到左、后从上到下的顺序进行**重组**。我们将文字段恢复为**每列**的形式，并按照从右到左的顺序输出。

经过实验，我们还设置了一定的**像素阈值**，因此我们的算法能够应对矩形框的检测误差，并给出正确的结果。这使得我们的算法拥有较强的**鲁棒性**。



原始模型对古籍图像原始数据进行文字识别的结果



采用本文算法预处理和调整模型的文字识别结果

## ① 部分结果

今獨盤鼎在耳  
新城道中  
山村五絕  
湖上夜歸  
寒食未明至湖上太守未來兩縣  
令先在  
次韻孫莘老見贈時莘老移處州  
因以別之

飲湖上初晴後雨  
往富陽新城李節推先行三日留  
風水洞見待  
風水洞和李節推  
獨游富陽普照寺  
自普照游二庵  
富陽妙庭觀董雙成故宅發地得  
丹鼎覆以銅盤承以琉璃盆  
盆既碎丹亦為人爭奪持去

- 1: 飲湖上初晴雨 0.944
- 2: 往富陽新城李節推先行三日 0.903
- 3: 風水洞見待 0.899
- 4: 風水洞和李節推 0.907
- 5: 獨游富陽普照寺 0.991
- 6: 自普照游二庵 0.996
- 7: 富陽妙庭觀董雙成故宅發地得 0.987
- 8: 丹鼎覆以銅盤承以琉璃盆 0.964
- 9: 盆既碎丹亦為人爭奪持去 0.906
- 10: 今獨盤鼎在耳 0.863
- 11: 新城道中 0.994
- 12: 山村五色 0.921
- 13: 湖上夜歸 0.898
- 14: 寒食未明至湖上太守未來兩縣 0.952
- 15: 令先在 0.939
- 16: 次孫莘老見贈時莘老移處州 0.731
- 17: 因以別之 0.997
- 18: 別 0.996



## ① 部分结果

麗人行一首	曲江一首	飲中八仙歌一首	送孔巢父謝病歸游江東兼呈李白一首	九日寄岑參一首	示從孫濟一首	同諸公登慈恩寺塔一首	苦雨奉寄西公兼呈王街君一首	贈衛八處士一首	醉歌行一首	醉時歌一首	秋雨數三首	白絲行一首	天育驃騎歌一首	高都護驄馬行一首	兵車行一首	貧交行一首	今夕行一首	玄都壇歌一首
-------	------	---------	------------------	---------	--------	------------	---------------	---------	-------	-------	-------	-------	---------	----------	-------	-------	-------	--------

- 1: 立都增歌一首 0.990
- 2: 今夕行一首 0.997
- 3: 貧交行一首 0.978
- 4: 兵車行一首 0.926
- 5: 高都護馬行一首 0.901
- 6: 天育驃暗歌一首 0.887
- 7: 白絲行一首 0.979
- 8: 秋雨三首 0.998
- 9: 欺亭前甘菊花一首 0.913
- 10: 醉時歌一首 0.988
- 11: 醉歌行一首 0.995
- 12: 贈衛八處三首 0.677
- 13: 苦雨奉寄西公兼呈王街君一首 0.962
- 14: 同公登慈恩寺塔一首 0.999
- 15: 示孫濟一首 0.997
- 16: 九日寄岑条一首 0.926

- 17: 送孔巢父病歸游江東兼呈李白一首 0.970
- 18: 飲中八仙歌一首 0.997
- 19: 曲江首 1.000
- 20: 麗人行一首 0.968



### ① 部分结果

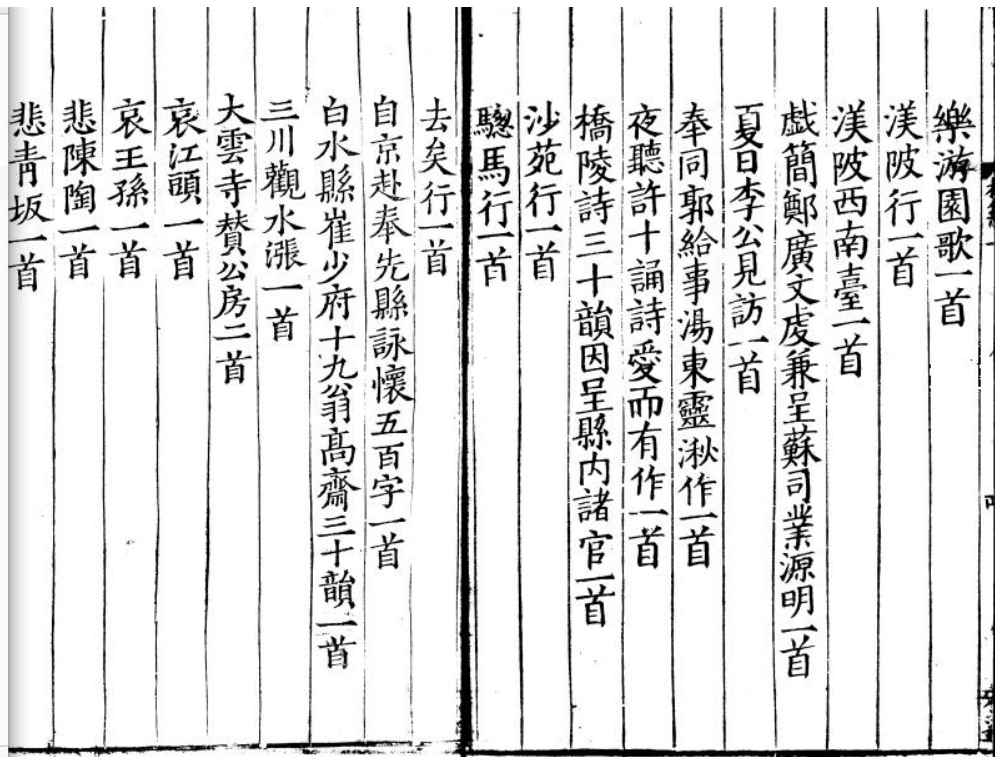


- 1: 办故煤 0.705
- 2: 一又准大中祥符元年六月五日 0.993
- 3: 办道有形器之適物有象數之滋一父始畫於龍圖 0.957
- 4: 八體遂生於鳥書契是造文字教興事增華 0.852
- 5: 類浸長載以發本尚律之音集成書抑 0.918
- 6: 亦久矣朕聿遵 0.762
- 7: 先志遵揚素風設教崇文科取士考程進實 0.919
- 8: 用焉而舊本既學者多必豕魚之革乃朱紫 0.893
- 9: 以洞分爰儒臣叶宣精力校增正刊綜 0.951
- 10: 其網條灼然叙列伸之基刻垂于將來仍特換於新 0.862
- 11: 名庶永昭於成宜改為大宋重修廣至准 0.941
- 12: 办 0.707
- 13: 勃故 0.776
- 14: 睦法言撰本 0.951
- 15: 長孫訓箋注 0.791
- 16: 外史顏之推 0.907

- 17: 儀同三司劉臻 0.939
- 18: 著作郎魏淵 0.954
- 19: 武陽太守思道 0.992
- 20: 散常侍李若 0.996
- 21: 國子博士 0.985
- 22: 蜀王议參軍辛德源 0.908
- 23: 吏部侍郎薛道衡已上八人同撰集 0.948
- 24: 郭知玄拾遺正更以朱箋三百字 0.946
- 25: 關亮增加字 0.941
- 26: 薛增加字 0.993

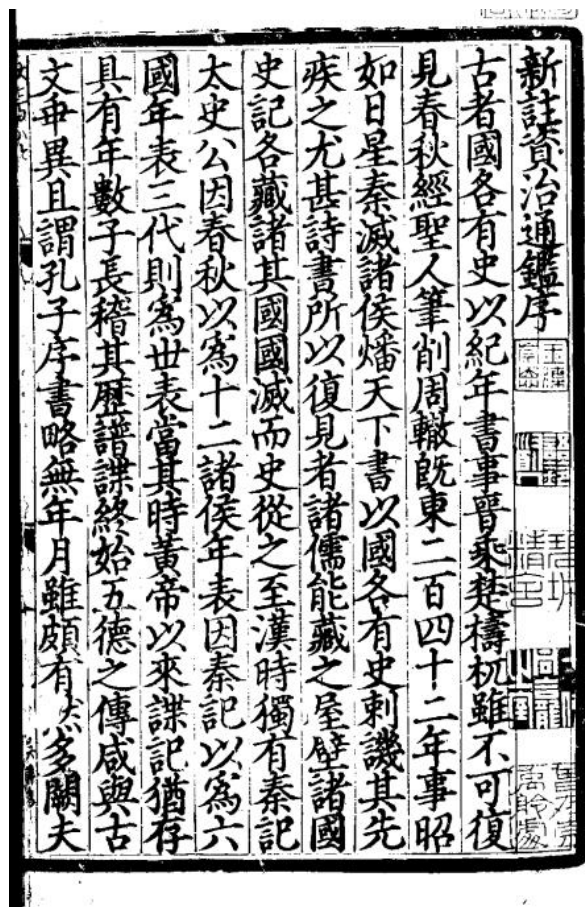
## ② 导出结果

1	文本: 樂游園歌一首	坐标: [[3508.0, 481.0], [3672.0, 481.0], [3672.0, 1281.0], [3508.0, 1281.0]]
2	文本: 漢陂行一首	坐标: [[3352.0, 489.0], [3488.0, 489.0], [3488.0, 1141.0], [3352.0, 1141.0]]
3	文本: 漢陂西南臺一首	坐标: [[3164.0, 485.0], [3324.0, 485.0], [3324.0, 1448.0], [3164.0, 1448.0]]
4	文本: 簡鄭廣文虔兼呈蘇司業源明一首	坐标: [[2996.0, 489.0], [3141.0, 489.0], [3141.0, 2462.0], [2996.0, 2462.0]]
5	文本: 夏日李公見訪一首	坐标: [[2832.0, 509.0], [2953.0, 509.0], [2953.0, 1499.0], [2832.0, 1499.0]]
6	文本: 奉同郭給事湯東靈湫作一首	坐标: [[2640.0, 494.0], [2785.0, 493.0], [2797.0, 2095.0], [2653.0, 2096.0]]
7	文本: 夜許十誦詩愛而有作一首	坐标: [[2457.0, 490.0], [2617.0, 489.0], [2625.0, 2099.0], [2465.0, 2100.0]]
8	文本: 陵詩三十因呈縣內官一首	坐标: [[2289.0, 497.0], [2449.0, 497.0], [2449.0, 2295.0], [2289.0, 2295.0]]
9	文本: 沙苑行一首	坐标: [[2112.0, 492.0], [2260.0, 488.0], [2279.0, 1127.0], [2130.0, 1131.0]]
10	文本: 馬行一首	坐标: [[1957.0, 489.0], [2121.0, 489.0], [2121.0, 1134.0], [1957.0, 1134.0]]
11	文本: 去矣行一首	坐标: [[1768.0, 484.0], [1890.0, 480.0], [1907.0, 1107.0], [1786.0, 1111.0]]
12	文本: 自京赴奉先縣懷五百字一首	坐标: [[1593.0, 466.0], [1738.0, 465.0], [1750.0, 2111.0], [1606.0, 2112.0]]
13	文本: 白水縣崔少府十九翁高齋三十一首	坐标: [[1394.0, 451.0], [1578.0, 449.0], [1598.0, 2445.0], [1415.0, 2447.0]]
14	文本: 三川觀水漲一首	坐标: [[1226.0, 467.0], [1374.0, 464.0], [1392.0, 1354.0], [1243.0, 1357.0]]
15	文本: 大雲寺贊公房二首	坐标: [[1070.0, 469.0], [1215.0, 469.0], [1215.0, 1476.0], [1070.0, 1476.0]]
16	文本: 哀江頭一首	坐标: [[898.0, 477.0], [1051.0, 477.0], [1051.0, 1118.0], [898.0, 1118.0]]
17	文本: 二首	坐标: [[741.0, 836.0], [854.0, 846.0], [833.0, 1105.0], [720.0, 1095.0]]
18	文本: 哀王孫	坐标: [[711.0, 481.0], [863.0, 481.0], [863.0, 911.0], [711.0, 911.0]]
19	文本: 一首	坐标: [[559.0, 891.0], [676.0, 891.0], [676.0, 1114.0], [559.0, 1114.0]]
20	文本: 悲陳陶一	坐标: [[547.0, 481.0], [699.0, 481.0], [699.0, 947.0], [547.0, 947.0]]
21	文本: 悲青坂一首	坐标: [[379.0, 493.0], [516.0, 493.0], [516.0, 1149.0], [379.0, 1149.0]]





## ② 导出结果



- 1 D0000001.JPG图片识别结果
- 2 文本: 新注資治通鑑序 坐标: [[795.0, 329.0], [870.0, 328.0], [875.0, 803.0], [799.0, 804.0]]
- 3 文本: 古者國各有史以紀年書事普乘楚 坐标: [[710.0, 339.0], [783.0, 339.0], [783.0, 1220.0], [710.0, 1220.0]]
- 4 文本: 二年事 坐标: [[645.0, 1340.0], [694.0, 1340.0], [694.0, 1515.0], [645.0, 1515.0]]
- 5 文本: 見春秋經聖人筆削周鞭既東二百四十一 坐标: [[622.0, 333.0], [696.0, 332.0], [712.0, 1367.0], [639.0, 1368.0]]
- 6 文本: 疾之尤甚詩書所以復見者儒能藏之星壁國 坐标: [[452.0, 333.0], [533.0, 332.0], [546.0, 1581.0], [464.0, 1582.0]]
- 7 文本: 史記各藏其國國滅而史之至漢時獨有秦 坐标: [[375.0, 335.0], [446.0, 334.0], [458.0, 1548.0], [387.0, 1549.0]]
- 8 文本: 太史公因春秋以寫十二諸侯年表因秦記以 坐标: [[283.0, 331.0], [371.0, 330.0], [383.0, 1550.0], [296.0, 1551.0]]
- 9 文本: 國年表三代為世表當其時黃帝以來課記 坐标: [[200.0, 335.0], [287.0, 334.0], [302.0, 1534.0], [215.0, 1535.0]]
- 10 文本: 具有年數子長稽其歷譜終始五德之傳咸與古 坐标: [[125.0, 337.0], [206.0, 336.0], [210.0, 1566.0], [129.0, 1566.0]]
- 11 D0000002.JPG图片识别结果
- 12 文本: 子之弗次盖其也子長述夫子之意故其表三 坐标: [[2393.0, 366.0], [2491.0, 365.0], [2500.0, 2034.0], [2401.0, 2034.0]]
- 13 文本: 代也以世不以年没豕紀年出於署太康初年 坐标: [[2276.0, 418.0], [2377.0, 417.0], [2391.0, 1976.0], [2290.0, 1977.0]]
- 14 文本: 次起自夏殷周止魏哀王之二十年此魏國史記 坐标: [[2170.0, 376.0], [2268.0, 376.0], [2274.0, 1976.0], [2176.0, 1977.0]]
- 15 文本: 紀表書傳世家自班孟堅以下不能易以紀紀年 坐标: [[1931.0, 363.0], [2048.0, 363.0], [2059.0, 2052.0], [1942.0, 2053.0]]
- 16 文本: 卷侯景之鼠王僧辯平建業與文德殿書七萬卷俱 坐标: [[1612.0, 374.0], [1705.0, 373.0], [1713.0, 2052.0], [1620.0, 2052.0]]
- 17 文本: 西江陵之間其書焉唐四庫書年四十一家九 坐标: [[1493.0, 371.0], [1599.0, 371.0], [1610.0, 2054.0], [1504.0, 2055.0]]
- 18 文本: 百四十七卷而王仲淹元經十五卷莆士依春秋 坐标: [[1386.0, 348.0], [1477.0, 347.0], [1488.0, 2057.0], [1397.0, 2058.0]]
- 19 文本: 義類作傳百卷矣今四十一家書存者復無乙 坐标: [[1057.0, 371.0], [1163.0, 371.0], [1163.0, 2042.0], [1057.0, 2042.0]]
- 20 文本: 部書以固等書為正史年類次之蓋紀傳表志 坐标: [[946.0, 376.0], [1047.0, 376.0], [1041.0, 2042.0], [940.0, 2042.0]]
- 21 文本: 英宗之書行年之書特以備乙庫之藏耳宋朝 坐标: [[834.0, 373.0], [930.0, 373.0], [930.0, 2037.0], [834.0, 2037.0]]
- 22 文本: 皇帝命司馬光次歷代君臣事迹為年一書 坐标: [[723.0, 368.0], [829.0, 369.0], [821.0, 2003.0], [714.0, 2002.0]]
- 23 文本: 盛衰繁生民休戚善可寫法惡可寫戒者以為是書 坐标: [[388.0, 373.0], [497.0, 374.0], [489.0, 2042.0], [380.0, 2042.0]]
- 24 文本: 治平熙寧間公與諸人國事相是非之日也蕭曹 坐标: [[274.0, 376.0], [383.0, 376.0], [377.0, 2042.0], [268.0, 2042.0]]
- 25 D0000003.JPG图片识别结果
- 26 文本: 智伯才德之焚英名實之唐太宗君臣之樂 坐标: [[2393.0, 368.0], [2500.0, 368.0], [2497.0, 2042.0], [2391.0, 2042.0]]

## ③ 效果统计

模型	数据	检测准确率	识别准确率
原始模型	原始图像	19.57%	8.621%
原始模型	预处理后图像	59.38%	57.47%
本文算法训练后模型	预处理后图像	63.22%	60.15%

各种情况的检测准确率和识别准确率

采用本文数据预处理方法以及基于深度学习的古籍文字检测与识别算法训练后的模型和原始模型相比，**检测准确率和识别准确率均明显提高。**



原始模型与原始数据(上)与本文算法流程(下)效果对比



## ① 技术学习

在开发初期团队成员共同学习相关的技术基础，查找相关资料，并定期组织团队内部的学习分享会议，分享各自的学习内容、想法和方案。

## ② 明确分工

在得到初步的解决方案思路后进行分工，明确每个成员的具体任务，并定期开线上会议进行交流，不仅需要分享各自的进展，并且需要交流遇到的困难，这是非常重要的。

## ③ 积极沟通

在整个参赛过程中，我们强调团队合作和开放的沟通氛围。每个成员都可以提出想法并分享解决问题的方法，积极且高频率地沟通，互相学习。



北京理工大学  
BEIJING INSTITUTE OF TECHNOLOGY

谢谢各位评委老师!

