後數位化時代的 OCR算法與應用創新

釋賢超 北京市海淀區龍泉寺藏經辦公室

https://www.gj.cool/author

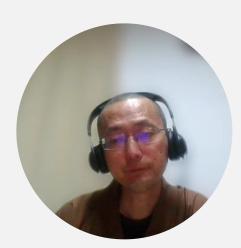


内容大纲

OCR: 发展前景和融合框架

古籍酷: 应用平台

籍智达:事业拓展



OCR与大模型

- · 多模态大模型可以实现OCR的大部分能力
- · OCR也能够辅助大模型获得更好表现
- OCR相对大模型的优势——性价比; 劣势——通用性
- •中文古籍的特殊性:数据偏少,使用面窄,难以吸引资本和技术

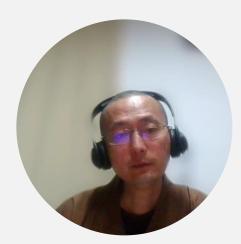
中文古籍OCR框架:多模型融合

• 版式特征: 双行夹注

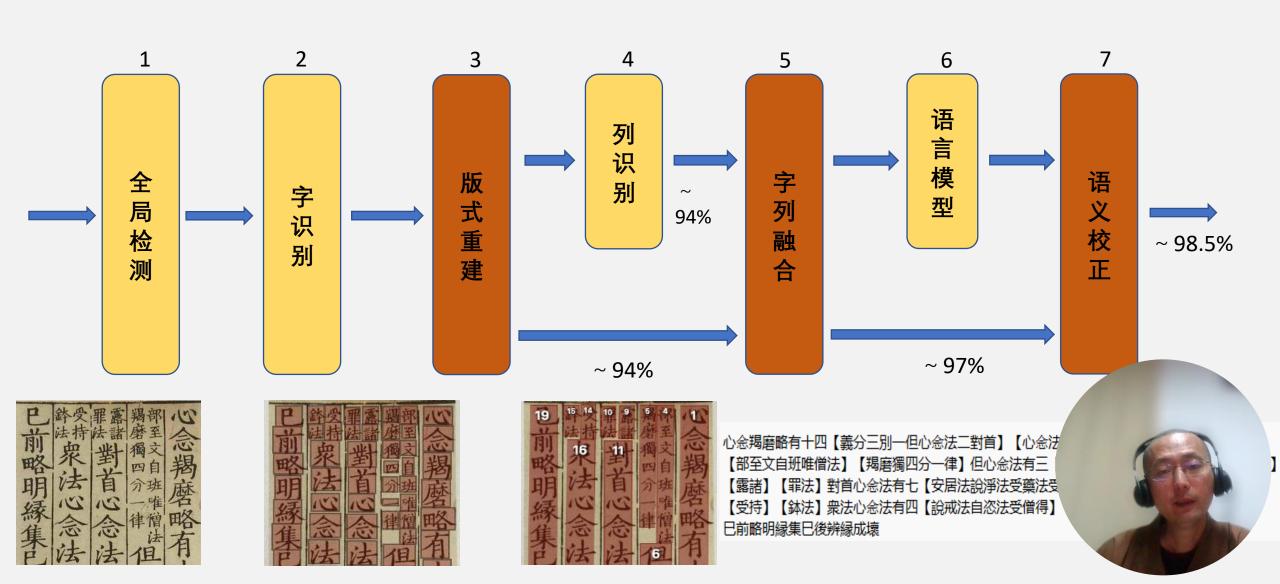
• 开发时间: 2021年5月-12月

• 知识准备: 2020年字符串比对

• 算法创新: 针对古籍版式的字符级版式重建



Multi-Model Ensemble Framework



古籍酷: 应用平台



X

● 准备就绪,请上传/拍照



問梅詩社詩鈔卷一 後附東籬會詩 癸未 探梅結社紀事

九興詩

我過北郭里每懷沈尚書詩社問遺老僉曰久荒蕪中 吳盛文藪大雅輪誰扶黄子豪于詩夙以山水娛激我 出西郊尋詩訪林逋歸田有彭子欣然命籃輿春風來 吹送一路三人俱行行詣古寺冷香襲襟裾巡簷佇清 影啜茗曲茄蕊面游小雪柄

上传拍照 @

标点⊕

翻译 🕓

清理宣

電記 暇日與文文起 朱德升過此

今亦三人來後先若合符黃子大歡喜

自动标点

砂 問梅詩社詩鈔卷一後附。東籬會詩。癸未探梅結社紀 事。九興詩。我過北郭里,每懷沈尚書。詩社問遺 老, 愈日久荒蕪。中吳盛文藪, 大雅輪誰扶。黄子豪 于詩, 夙以山水娛。邀我出西郊, 尋詩訪林逋。歸田 有彭子, 欣然命籃輿。春風來吹送, 一路三人俱。行 行詣古寺,冷香襲襟裾。巡簷佇清影,啜茗兼茹蔬。 再游小雲栖周忠介公書額。不辭幽徑紆,仰瞻忠介 書。凜然生眉鬚,相傳過此地。公嘗偕文朱公作是庵 記云: 暇日與文文起、朱德升過此, 今亦三人來。後 先若合符, 黄子大歡喜。

上传拍照 @

标点 🗩 📗 翻译 🚱 📗 清理 🗑

文言文翻译

X

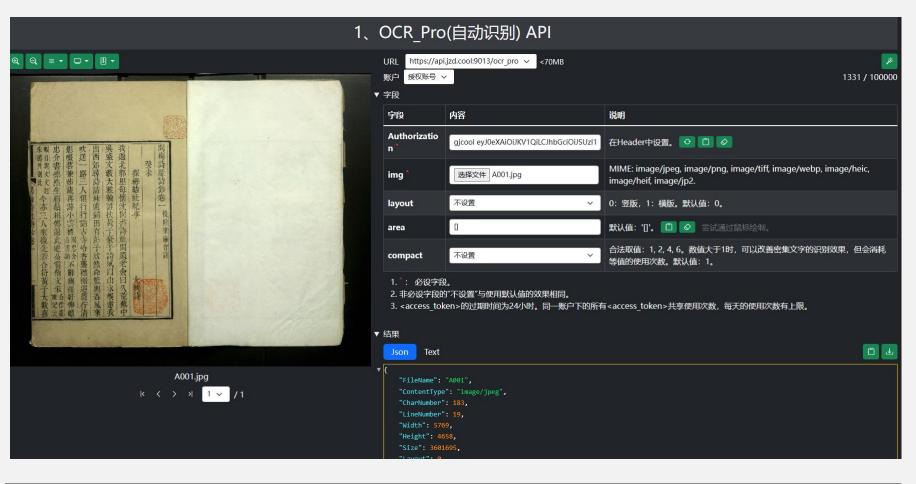
❷ 请他读《梅诗》和《诗经》,然后才将之附在后面。 东门诗是这样。二十六日,探索梅子结社纪事。九兴 诗。我路过北郭外城,常常思念着沉于尚书。诗社问 遣老,都说长期荒芜。中吴盛有文才,大雅有谁扶 持? 黄子豪富于诗,早以山水为娱乐。邀我去西郊, 寻访林逋。归田前有彭子,高兴兴地坐在车上。春风 吹来送走了我,一路上三人都来了。他走到古寺,脱 下了清香,穿上了衣襟。我只有等待清晨的影子,吃 茶叶也吃蔬菜。第二次游历小云,是周朝忠诚耿介的 公文。不辞幽隐曲折,仰瞻忠节书籍。严生凛然眉毛 长满胡须, 一路前来经过这里。不辞幽隐曲折, 仰瞻 忠介书。谢安曾经和文朱公一起作这样的记载:闲暇 的时候与文起、朱德升一起经过这里,如今也有三个 人来。不辞幽隐曲径,仰瞻忠介书,皇后若合符,黄 子大喜。

ſ°l

上传拍照 @

标点 🖭 📗 翻译 🖗









详尽的文档说明















₩ -

66 🔻

 \leq

品 固定

₩ -

无限缩放, 毫厘可见





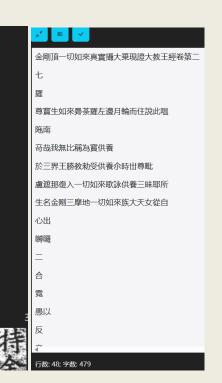
辅助编辑: 多多益善















時世尊毗

奈<mark>☆</mark> 時時

所所

24

袋 供 供











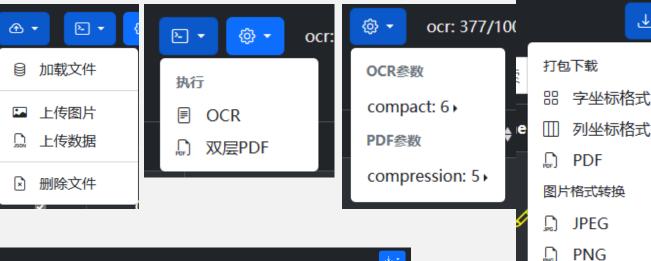








批量处理



TIFF

WEBP

AVIF

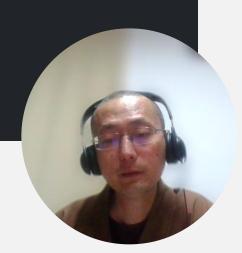
HEIC

JPEG 2000



隐私政策

- 1、本网站<mark>不搜集、不存储用户数据</mark>(包括但不限于:上传图片、OCR结果以及用户编辑内容,上传文本与自动标点结果),也不对用户数据进行"<mark>合理利用</mark>"(例如 改进模型和相关算法等),仅对用户最近一次的上传图片和OCR结果进行临时性存储。
- 2、本网站不存储用户密码明文,并对用户密码进行了高强度加密,理论上可以保证除用户本人之外的任何人无法在合理时间内还原密码明文。
- 3、本网站不对用户使用OCR、自动标点所产生的数据主张权利或授予权利,不对用户合理使用数据的行为附加任何条件或限制,但同时也不承担与之相关的任何直 接或间接的法律责任和道德义务。
- 4、出于维护网站正常运行的需要,除了用户注册时提供的信息之外,本着"<mark>最低程度"、"不触及用户内容</mark>"的原则,用户以下的使用行为和数据将被酌情记录:
 - (1) 提交ocr或自动标点单次请求的时间、IP地址;
 - (2) 上传图片的体积;
 - (3) 自动标点原文的字数;
 - (4) 持续长时间API请求所提交的图片或文本。



更新历史

版本: RC 5.24

一般更新:

1、API演示。优化jp2文件加载,增加下载进度提示。

— 2024.3.29

版本: RC 5.23

一般更新:

- 1、首页快速体验。增加上传进度提示。
- 2、API演示。增加演示账号的当前用量的实时反馈。

— 2024.3.2*&*

版本: RC 5.22

重要更新:

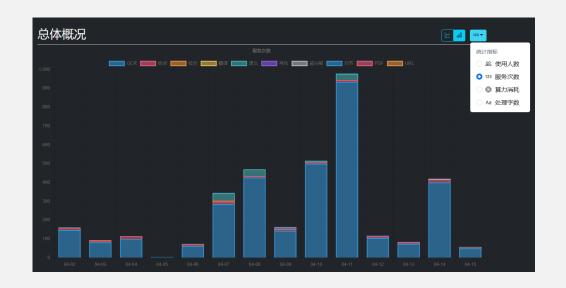
- 1、双层PDF API参数调整。compression取值范围调整为0-5的整数,数值越大则图像压缩程度越大,生成的PDF文件体积越小。
- 2、标注平台。优化图片上传按钮,新增文件信息展示。批量处理界面的图片上传、数据上传分为两个独立按钮,增加当前使用量的实时反馈。

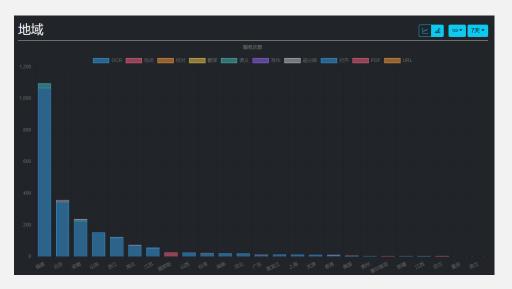
版本: RC 5.21

一般更新:

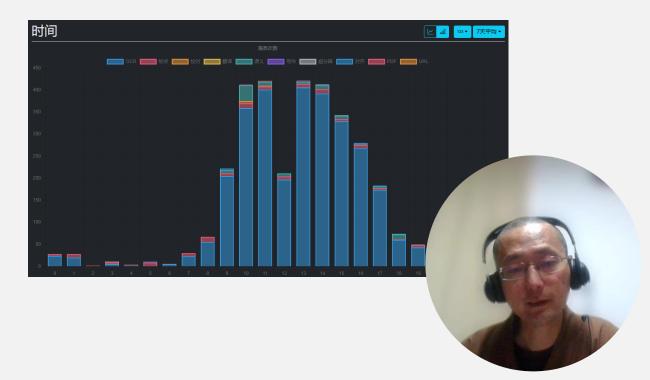
1、导航栏调整。







服务类型						□ OCR •	
日期	¢	使用人数	¢	服务次数 💠	算力消耗 ◆	处理 <u>字数</u> 服务类型	
2024-04-15		11		50	88	12325	
2024-04-14		35		400	568	93069	
2024-04-13		20		73	136	28733 🔾 🚇 🛔	翻译
2024-04-12		18		106	250	46711 G 🖮 🗎	
2024-04-11		54		934	1389	81346	
2024-04-10		26		500	794	33158 ⊝ ≡ 🛪	型分辨 Ht字
2024-04-09		31		143	374	44263 A	
2024-04-08		36		426	562	74659 ○ ∞ ∪	JRL
2024-04-07		47		285	584	93242	
2024-04-06		15		62	106	24020	
2024-04-05		2		4	7	2159	
2024-04-04		20		99	177	38991	
2024-04-03		12		82	209	32214	
2024-04-02		24		148	252	42209	



北京某大学和某互联网公司合作开发的某典平台

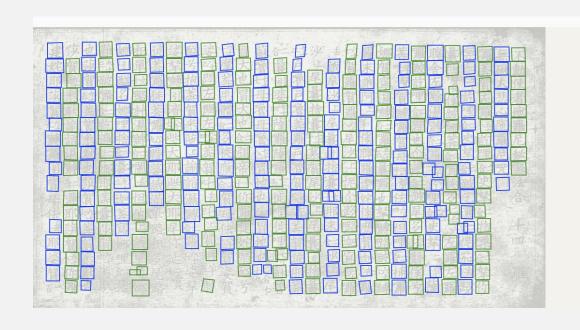
文字识别校对帮助文档

环节目标: 切分正确、顺序正确、文字正确

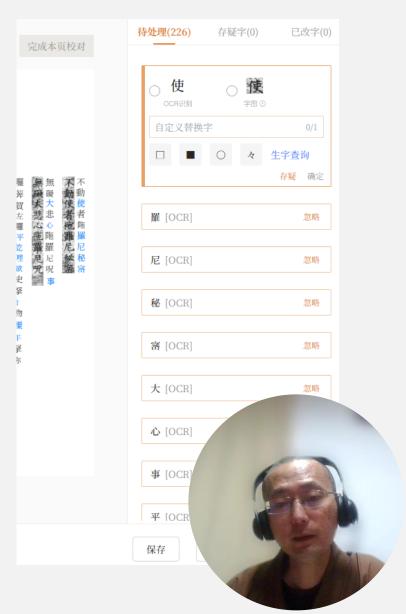
1. 调整切分: 保证每个应被识别的字, 都有字框; 并且字框归属于正确的列框。

2. 调整顺序: 文字间的阅读顺序正确。

3. 校对文字: 每个字都识别为正确的写法。



不動使者陁羅尼秘密 無礙大悲心陁羅尼呪事 **噻辨賀左囉平訖哩欲史拏哈物嘌年拏你** 禀加望訖哩合史拏合跛乞灑云合抳茄去跢 義之呼傒跛娜舞合賀徙多二合左羅左羅 聴舎左輪濕嚩合囉至乞三哩丑合史拏合薩 **哪**跛訖哩二蟬也人論合跛尾多三夲翳後兮 **养浮聲賀崎囉賀母理聲佉此怛哩苔補囉** 娜賀寧濕嚩合囉十鬟囉也拏嚩路跛碎嚩 羅末誠阿上唎後聹羅建姹溪麼賀迦法羅 三賀羅賀羅本尾沙注抳尒跢路迦寫摩囉 識尾沙上尾鬟引上捨曩七七那未二少字尾 尾囊捨襲十慕賀尾沙上尾襲引捨襲 法戸噜戸噜莾羅戸嚕賀曚切三莾賀跛那莾 曩引婆呌薩囉薩囉五徙哩徙哩七際蘇曽 蘇嚕切母嚕母嚕注冒地也能冒地二也七 冒大也冒大也合弭帝子你羅建婬翳醯 十座葬思體合多徙廐合賀母餘智 娑斗悶左悶左十四莾賀吒去吒上賀無文 醯兮抱莾賀悉陁上諭詣濕嚩合羅 娑拏嚩引濟十七些大耶譽大耶尾 徒恭囉徙莾囉碎瞻婆識滿 枳單协路計濕嚩合隬去怛他上識單九一娜 娜引醯名娜哩捨合鱗九迦莾寫 雖子時關能約擺汀縣耶瑟屬恭智



中研院文字識別與校對平台





籍智达: 事业拓展

2021年, 注册成立。

2022年, API网络服务。校对平台网络免费体验。

2023年,服务器租用。发布"数字万舟"公益计划。

2024年,服务器出售。校对平台本地化部署。





合作单位

复旦大学中国语言文学系 浙江大学数字人文研究中心 山东大学数字人文实验室 武汉大学数字人文研究中心 对外经济贸易大学中国语言文学学院 江苏师范大学中华家文化研究院 商丘师范学院汉梁文化研究中心 香港中文大学图书馆



















"数字万舟"计划简介

中文古籍OCR工业版自2022年发布以来广受各界好评,得到了一批忠实用户的鼎力支持。他们的积极反馈和殷切期望,为我们注入了源源不断的前进动力。正所谓"士不可不弘毅,任重而道远。"一方面是尘封的古籍在手中日渐老去,另一方面是数字列车从身边飞驰而过。愿望与现实间的巨大落差,无时无刻不在触动每一位中华学人心中的软处。

有感于此,古籍酷正式对外发布"数字万舟"计划,面向社会征集优秀中文古籍数字化项目,为传承中华古籍文化贡献绵薄之力。符合申请条件的项目,将获赠中文古籍OCR工业版免费使用额度。

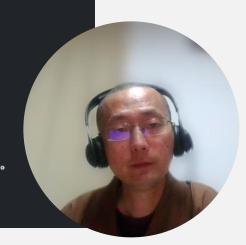
申请者

申请者应为以下三者之一:

- 1、学术或教育单位;
- 2、社会非营利机构;
- 3、个人。

申请条件

- 1、健康。符合申请者所在国家或地区以及服务对象所在国家或地区的法律和社会伦理,具有正面社会价值。
- 2、公益。不含任何盈利性成分。
- 3、透明。在网络主页及时公布项目的进展情况和成果信息,准确描述古籍酷的技术贡献,在描述中应包含"古籍酷"和"gi.cool"的明确字样。
- 4、授权古籍酷及其关联单位对该项目的公开内容进行转载或展示。



数字万舟的受赠者

• 单位:

- 复旦大学中国语言文学系
- 对外经济贸易大学中国语言文学学院
- 江苏师范大学中华家文化研究院
- 商丘师范学院汉梁文化研究中心
- 铜仁市图书馆

个人:

- 哈佛大学、埃默里大学、俄亥俄州立大学、匹兹堡大学、 马来西亚理科大学, ……
- 北京大学、复旦大学、浙江大学、北京师范大学、中央美术学院、华东师范大学、华中师范大学、华中科技大学、西北大学、湖北大学、广东外语外贸大学、河南财经政法大学、南京理工大学、厦门大学、香港中文大学、澳门理工大学、……

项目类型

- 说文,集韵,训诂小学
- 左传,清史稿,清宫档案
- 经籍志,省郡县府镇志,山水志,运河文献
- 文集,小说
- 类书
- 中医方剂
- 家谱,族谱,家训
- 科举文献
- 民间文献
- 石经、石刻
- 政区史、海洋史、农村史
- 专名提取、历史计量学
- 语料库
- 维基文库、中国哲学电子书计划





中文古籍的文字勘探與處理. Chinese Classic Text Mining and Processing

圖書館歡迎中大師生參加以下關於中文古籍識別及文字處理平台的工作坊,提升研究效率。

古籍酷智能平台與AI技術在中文古籍數字化中的實踐應用

講者:

賢超法師 北京市海淀區龍泉寺藏經辦公室主任

張敬女士 北京籍智達數字科技有限公司總經理

3月星 MARCH 期 23四 Thursday

香港中文大学图书馆中文古籍的文字勘探与处理工作坊。2023-3



杭州: 第二届东亚古籍数字人文研讨会。2023-10

籍智达与嘉兴精严讲寺签署战略合作协议。本人受聘嘉兴市嘉兴藏文化研究会顾问。2023-11







黄河科技学院:中华优秀传统文化与黄河文化学术论坛。2023-11



北京: 第二十届全国科技翻译研讨会。2023-11



曲阜: 全国家文化研究机构联席会议二次会议暨颜氏家训家

风与中华民族家文化研讨会。2023-11

展望

