# Data Extraction Using Documents RPA

**Raste Suchitra[1]\* Nalwande Nikita[2] Mr. Phule U. G.[3]**

[1,2] TY. Students, Department of Computer Engineering, Sahakar Maharshi Shankarrao Mohite Patil Institute of Technology and Research, Akluj, Solapur, Maharashtra, India

[3] Lecturer, Department of Computer Engineering, Sahakar Maharshi Shankarrao Mohite Patil Institute of Technology and Research, Akluj, Solapur, Maharashtra, India

*Abstract – Scanned receipts OCR and key facts extraction (SROIE) constitute the processeses of spotting textual content from scanned receipts and extracting key texts from them and shop the extracted checks to established documents. SROIE performs important roles for lots report evaluation packages and holds terrific business potentials, however little or no studies works and advances had been posted on this area. In popularity of the technical challenges, significance and massive business potentials of SROIE, we prepared the ICDAR 2019 opposition on SROIE. In this opposition, we installation 3 tasks, namely, Scanned Receipt Text Localization (Task 1), Scanned Receipt OCR (Task 2) and Key Information Extraction from Scanned Receipts (Task 3). A new dataset with one thousand entire scanned receipt pictures and annotations is created for the opposition. The opposition opened on tenth February, 2019 and closed on fifth May, 2019. There are 29, 24 and 18 valid submissions acquired for the 3 opposition tasks, respectively. In this document we are able to provide the motivation, opposition datasets, project definition, assessment protocol, submission statistics, performance of submitted strategies and consequences evaluation. According to the extensive hobbies won thru SROIE and the wholesome number of submissions from academic, studies institutes and enterprise over one of kind countries, we agree with the opposition SROIE is successful. And it's far thrilling to examine many new thoughts and approaches are proposed for the brand new opposition project set on key facts extraction. According to the overall performance of the submissions, we agree with there may be nevertheless a huge hole at the anticipated facts extraction overall performance.*

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - *x* - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

## INTRODUCTION

Scanned receipts OCR are a technique of spotting textual content from scanned dependent and semi-dependent receipts and in voices. On the alternative hand, extracting key texts from receipts and invoices and keep the texts to dependent files can serve many packages and services, which includes green archiving, rapid indexing and file analytics. Scanned receipts OCR and key statistics extraction (SROIE) play vital roles in streamlining file-extensive procedures and workplace automation in lots of financial, accounting and taxation areas. However, SROIE additionally faces huge challenges. With perform an greatly boosted with the aid of using current breakthroughs in deep getting to know technology in phrases of accuracy and processing speed, OCR is turning into mature for lots realistic duties (which include call Huang is with Shanghai Jiao tong University, China, huangzheng@sjtu.edu.cn. Jianhua He is with Aston University, UK,j.he7@aston.ac.uk. Kai Chen (kaichen@onlyou.com) is with Onlyou, China.Xiang Bai (xbai@hust.edu.cn) is with Huazhong University of Science and Technology, China. Dimosthenis Karatzas (dimos@cvc.uab.es), Universitat Autonoma de Barcelona, Spain. Shijian Lu (Shijian.Lu@ntu.edu.sgis with Nanyang Technological University, Singapore. C. V. Jawahar(jawahar@iiit.ac.in) is with IIIT Hyderabad, India.card reputation, registration code reputation and hand-written textual content reputation). However, receipts OCR has a lot better accuracy necessities than the overall OCR duties for lots industrial packages. And SROIE turns into extra hard whilst the scanned receipts have low quality. Therefore, within side the present SROIE structures, human assets are nonetheless closely utilized in SROIE. There is an pressing want to investigate and develop rapid, green and strong SROIE structures to lessen or even dispose of guide work. With the traits of OCR structures going to be extra intelligent and file analytics, SROIE holds unparalleled potentials and opportunities, which attracted big pastimes from big companies, which includes Google, Baidu and Alibaba. Surprisingly, there are little studies works posted within side the subject matter of SROIE. While strong reading, file format evaluation and named entity reputation are applicable to the SROIE, none of the existing studies and beyond ICDAR competitions absolutely address the issues confronted with the aid of using SROI. In reputation
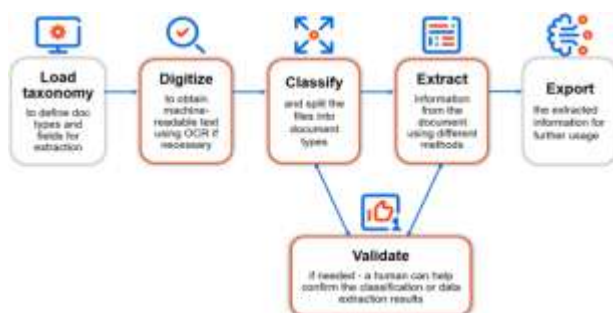
of the above challenge importance and big industrial potentials of SROIE, we prepared the ICDAR 2019 opposition on SROIE, aiming to attract interest from the network and sell studies and improvement efforts on SROIE. We agree with the opposition might be of pastimes to the ICDAR network from numerous aspects. First, to guide the opposition, a large-scale and well annotated bill datasets are provided, that is important to the achievement of deep getting to know primarily based totally OCR structures. While many datasets were gathered for OCR studies and competitions, to the nice of our knowledge, there may be no publicly to be had receipt dataset. Compared to the prevailing ICDAR and different OCR datasets, the brand new dataset has a few unique capabilities and challenges.

## LITERATURE REVIEW

In Manual Process of PDF extraction we want to do all project manually.PDF automation in Ui Path allows for doing all project concerning pdf documents with Robotic Process. That can lessen time and human force, RPA project is controlled through a specific

**GLYMPSE – SHARE GPS LOCATION**: - This is the latest software advanced on January 28, 2015. This app is a fast, loose and an easy manner to percentage our place the usage of GPS monitoring in actual time with buddies and family. This app does now no longer want any join up and do now no longer want any contacts to manage.

**Block Diagram: Document Understanding Framework**



## CONCLUSION AND FUTURE WORK

We prepared the one of the first competitions at the OCR and statistics extraction for scanned receipts. For the opposition SROIE we organized new datasets and assessment protocols for 3 opposition tasks. A true variety of submissions had been obtained for all 3 tasks, which confirmed extensive hobbies on the subject from the instructional and industry. And it's far exciting to examine many new thoughts and strategies are proposed for the brand new opposition mission of key statistics extraction. According to the overall performance of the submissions, we agree with there may be nonetheless a big hole at the predicted statistics extraction overall performance. The mission of key statistics extraction remains very difficult and may be set for lots different critical file evaluation packages. It may be exciting to increase in this opposition with extra difficult and large datasets and packages within side the future. The new datasets used on this opposition may be made to be had after the event.

## REFERENCE

[1]     He J., Chen H., et al. (2010) 'Adaptive congestion control for DSRC vehicle networks', IEEE Comm. Lett., 14, (2), p.127-129.

[2]     D. Karatzas, F. Shafait, S. Uchida, M. Iwamura, L. Gomez, S. Robles, J. Mas, D. Fernandez, J. Almazan, L.P. de las Heras (2013). ICDAR 2013 Robust Reading Competition. ICDAR.

[3]     D. Karatzas, L. Gomez-Bigorda, A. Nicolaou, D. Ghosh , A. Bagdanov, M. Iwamura, J. Matas, L. Neumann, VR. Chandrasekhar, S. Lu, F. Shafait, S. Uchida, E. Valveny (2015): ICDAR 2015 robust reading competition. ICDAR.

[4]     Everingham, M. and Eslami, S. M. A. and Van Gool, L. and Williams, C.K. I. and Winn, J. and Zisserman, A. (2015). The Pascal Visual Object Classes Challenge: A Retrospective. IJCV, 2015

[5]     D. Karatzas, L. Rushinol (2015). The Robust Reading Competition Annotation and Evaluation Platform

**Corresponding Author**

**Raste Suchitra***

TY. Students, Department of Computer Engineering, Sahakar Maharshi Shankarrao Mohite Patil Institute of Technology and Research, Akluj, Solapur, Maharashtra, India

**Raste Suchitra[1]* Nalwande Nikita[2] Mr. Phule U. G.[3]**