

Challenges in the Design of Optical Character Recognition for Medical Image Modalities.

¹Efosa Osagie, ²Shemi Ayo-Ogbor,

¹Visiting Lecturer, School of Physics, Engineering and Computer Science, University of Hertfordshire, Hatfield, United Kingdom.

²Medical Officer, Ministry of Health, Al Taif, Kingdom of Saudi Arabia.

ABSTRACT: The aim of this review study is to analyse and evaluate the challenges faced in the research and design of optical character recognition systems for computer vision tasks in processing medical image modalities. In this regard, this study provides a summarized comprehensive foundation of optical character recognition, identifies the challenges faced in the design of optical character recognition models in medical image modalities with attention to the burned-in text data in these images and in conclusion, provide benefits of this technology in health system, identify possible future research directions.

KEYWORDS: Optical character recognition, medical image processing, deep learning, artificial intelligence, computer vision, medical image modalities.

I. INTRODUCTION

With the vast growth in computing power available, there have been a rapid application of artificial intelligence in developing useful models for various medical imaging modalities. Medical imaging modalities includes ultrasounds and X-ray images taken from the human body during clinical examination. One of such modelling areas is the optical character recognition (OCR) which is becoming more and more applied for medical image processing. The traditional OCR process is a process of converting images of machine printed or handwritten characters into a format accessible by a computer for post-processing operations (Srihari et al, 2013). There has been a long history of research in this area as an important application of machine learning and more recently, deep learning algorithms, with the goal to access restricted form of text data which is more commonly burned-in on medical images. OCR has also been applied to access patient demographics and annotations for various purposes such as creating a searchable database (Silva et al, 2018), anonymization for privacy reasons (Newhauser et al, 2014), building structured patient record system (Li et al, 2015), amongst others. Major Electronic health record (HER) systems are already integrating OCR in accessing, extracting, and processing text in medical image modalities. Even with the vast development of OCR modelling techniques, there exist notable challenges been faced by researchers and developers, especially when dealing with medical image data containing burned-in text. Such medical image data containing burned-in text includes ultrasound, magnetic resonance imaging (MRI) and any other imaging acquired using a medical imaging acquisition scanner. This study aims to clarify the challenges faced in the research for the design of OCR for medical image processing. This study presented the medical image background problem in section II, described the issue of privacy and limited datasets in section III. In section IV, varied font features in medical images are discussed, section V provides benefits of an efficient OCR system in health system and VI, this study draws discussion and conclusion to possible research gaps identified in dealing with these challenges.

II. COMPLEX MEDICAL IMAGE BACKGROUND

Medical image data comprises of human body parts captured using various imaging acquisition modalities such as ultrasound and x-rays. To ensure a uniform distribution and viewing of medical images, the digital imaging and communication (DICOM) format was created by the national electrical manufacturers (NEMA), and this has been widely adopted since then (Bhagat et al, 2012). As part of clinical examination process, which involves capturing, relevant patient data including demographics are burned into the image by the acquisition machine (Monteiro et al, 2015). This results to an image consisting of a luminous background due to the lightening condition during capture. The patient health information exists as burned-in text data on the pixel content of the image and this is the text data, OCR seeks to access and extract. Traditional OCR models find it difficult to extract features from the non-uniform background of these medical images, which also have weak contrast owing to the poor background quality with a relatively small region of interest with predominant unclear boundaries in the entire image. These complex nature of the background leads to poor performance in traditional OCR models in medical image processing and has opened a direction for further research.

III. LIMITED MEDICAL IMAGE DATASET

This remains a major limitation in the application of artificial intelligence in general to the field of medical imaging, which also include the OCR for text detection and recognition for various reasons such as protected health information removal, classification, and pseudo-anonymization. The availability of clinical data between medical centres and researchers is greatly hindered by privacy protection requirements and accessibility. This has greatly resulted to only limited datasets available for the design, and evaluation of OCR models, as researchers must make do with limited data collections available, leading to a slow progress in the field. Most medical centres with databases of medical images, follow a regional regulatory framework such as the Health Insurance Portability and Accountability Act of 1996 in the United States (HIPAA), as mandate on privacy protection of patient's medical records which ensures a guarantee is given on the confidentiality of data during storage and transmission via any secured or unsecured means (Li et al, 2005). Due to this privacy problem, obtaining a sufficiently large-scale, balanced and properly annotated images for efficient training, testing and evaluation of OCR models, remains a major challenge (Qin et al, 2019). The summary from the first annual conference on machine intelligence in medical imaging (C-MIMI) was held in September 2016 and the common agreement was data starvation in research requiring medical image evaluation (Kohli et al, 2017). The ideal dataset for OCR evaluation in medical imaging should have adequate quantity, highly visible annotation, possible burned-in data which can be synthetic for privacy concerns and every other feature that makes up a typical clinical examination. During performance evaluation of an OCR model, these burned-in data are required along with any other printed patient's demographics on the medical image, to determine the model's accuracy using suitable metrics. There is at present, very limited number of fully annotated medical images containing fictitious patient's data existing in the public space, and this was confirmed from a wide systematic search using popular medical image databases.

VARIED FONT FEATURES OF IMAGING ACQUISITION MACHINE : There has been a great difficulty for the state-of-the-art OCR to recognise all different font types, sizes and family including too small or big characters which may be tedious to detect and identify in medical images. One major reason for this difficulty is the low resolution and extreme small size in which these text data appear on the pixel data of these medical images, this has made research in the design of OCR less efficient in identifying these varied fonts. Monteiro et al (2017) applied convolutional neural network (CNN) to design an OCR for character recognition, though it was able to identify characters in the validation dataset used, but all the images in the validation dataset were sourced from one medical institution using only their unique font feature. The issue of no standard font feature adapted as default by the electrical manufacturers of these imaging acquisition machines has become a challenge in the design of a robust, generalized, and adaptable OCR system for medical images.

BENEFITS OF OCR IN HEALTH SYSTEMS : Technology has advanced over the years and one area of keynote and importance is in that of healthcare management. Consequently, it has improved the quality of life over time, and saved many lives in the process. The OCR technology has been a solution for extracting text from medical image modalities for various data entry purposes useful for medical diagnostics and a robust electronic health management system (Bhure, 2021). Visual impairment which is a decreased ability to see to a degree has caused a lot of problems not fixable by usual means, such as glasses. According to a report from the global blindness and visual impairment data in 2015, there were an estimated 253 million people with visual impairment worldwide, out of which 36 million were blind and a further 217 million had various cases moderate to severe visual impairment (MSVI) include problem with seeing in low contrast and brightness condition (Ackland et al, 2017). Vision is important to not only see objects but also for dark adaptations, contrast sensitivity, balance, and colour perceptions. An efficient and highly accurate OCR system can provide health workers who are blind or visually impaired with the capacity to scan text on medical images and modalities and then have it converted into easily accessible forms such as synthetic speech and plain text. It can help them recognise abnormally without the help of a third party verification thereby enabling improvement of diagnosis speed and mobility. With the rapid entry into electronic health records (EHRs) systems which replaced the old paper-based storage and retrieval processes, which was designed to make patients' management more accurate, safer, and more accessible, this system involves large number of documentations, medical images, investigation reports and prescription, there is always the difficulty tracking files and keeping inventory (Dash et al, 2019). This can cause issues in getting statistical report on files especially where medical image modalities are involved and inefficient records disposition. OCR technology can bridge this issue by effectively categorising these clinical and non-clinical datasets into meaningful categories. Despite technological innovation in medical imaging and information system technologies, the radiology report has remained stagnant for many years. Imaging plays a pivotal role in the diagnostic process for many patients. With an average estimates of average diagnostic error rates ranging from 3% to 5%, there are approximately 40 million diagnostic errors involving

imaging annually worldwide (Jn Itret al, 2018). Physicians face these problems because of the ever-increasing burden of the number of patients and the qualified personnel available to attend to their cases. The potential to improve diagnostic performance and reduce patient harm by identifying and learning from these errors is very important. A highly accurate OCR can improve diagnostic performance and reduce patient harm by identifying and learning from these errors where it involves interpreting radiological findings relating to the burned-in text by the acquisition machine and discerning diseases using the technology thereby improving in terms of speed, reduced cost and man-power, easy access to information, easy patient's data management, and improved health security.

IV. DISCUSSION AND CONCLUSION

It is very important to apply the innovating and rapidly growing advantage of machine learning in the design of an optimal and highly accurate OCR for purpose of extracting text data from medical image which is usually needed for various post-processing needs such as recognition of sensitive character and pseudo-anonymization of the identified patients' data, amongst others. It is our view that this review article gives researchers and developers, a detailed summary of notable challenges been faced, and to seek various technological answers to these problems. Also, the cooperation of medical centres and research institutes around the world is required, especially in the issue of providing datasets for the design and evaluation stage, this can be achieved by replacing sensitive patient's data on the medical images, with synthetic data to prevent any breach of privacy regulation in their region. Hence, it is highly recommended that these points raised are considered in the design and modelling of OCR systems for medical image modalities.

REFERENCES

1. Anton Patyuchenko (2019), Medical Image Processing: From Formation to Interpretation, pp 3-4. Available from: <https://www.analog.com/media/en/technical-documentation/tech-articles/Medical-Image-Processing-From-Formation-to-Interpretation.pdf>
2. Ackland, P., Resnikoff, S., & Bourne, R. (2017). World blindness and visual impairment: despite many successes, the problem is growing. *Community eye health*, 30(100), 71–73.
3. Bhagat, A. P. and Atique, M. (2012) Medical images: Formats, compression techniques and DICOM image retrieval a survey, 2012 International Conference on Devices, Circuits and Systems (ICDCS). IEEE. DOI:10.1109/icdcsyst.2012.6188698.
4. Bhure, A. (2021), A Review of Optical Character Recognition (OCR) in Healthcare, *International Journal for Research in Applied Science and Engineering Technology*. International Journal for Research in Applied Science and Engineering Technology (IJRASET). DOI:10.22214/ijraset.2021.34142.
5. Dash, S., Shakyawar, S. K., Sharma, M. and Kaushik, S. (2019, June 19) Big data in healthcare: management, analysis and future prospects, *Journal of Big Data*. Springer Science and Business Media LLC. DOI:10.1186/s40537-019-0217-0.
6. Huang, L.-C., Chu, H.-C., Lien, C.-Y., Hsiao, C.-H. and Kao, T. (2009) Privacy preservation and information security protection for patients' portable electronic health records, *Computers in Biology and Medicine*, 39 (9), pp. 743–750. DOI:10.1016/j.combiomed.2009.06.004.
7. Kohli, M. D., Summers, R. M. and Geis, J. R. (2017, May 17) Medical Image Data and Datasets in the Era of Machine Learning—Whitepaper from the 2016 C-MIMI Meeting Dataset Session, *Journal of Digital Imaging*. Springer Science and Business Media LLC. DOI:10.1007/s10278-017-9976-3.
8. Li, M., Poovendran, R. and Narayanan, S. (2005) Protecting patient privacy against unauthorized release of medical images in a group communication environment, *Computerized Medical Imaging and Graphics*. Elsevier BV. DOI:10.1016/j.compmedimag.2005.02.003.
9. Li, X., Hu, G., Teng, X., & Xie, G. (2015). Building Structured Personal Health Records from Photographs of Printed Medical Records. *AMIA ... Annual Symposium proceedings*. AMIA Symposium, 2015, 833–842.
10. Monteiro, E., Costa, C. and Oliveira, J. L. (2017). A De-Identification Pipeline for Ultrasound Medical Images in DICOM Format. *Journal of Medical Systems*, 41(5). Available from: <http://dx.doi.org/10.1007/s10916-017-0736-1>.
11. Monteiro, E., Costa, C. and Oliveira, J. L. (2015) A machine learning methodology for medical imaging anonymization, in: 2015 37th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC). IEEE,.
12. Newhauser, W., Jones, T., Swerdloff, S., Newhauser, W., Cilia, M., Carver, R., Halloran, A. and Zhang, R. (2014) Anonymization of DICOM electronic medical records for radiation therapy, *Computers in Biology and Medicine*, 53, pp. 134–140. DOI:10.1016/j.combiomed.2014.07.010.

13. Qin, X., Bui, F. M. and Nguyen, H. H. (2019) Learning from an Imbalanced and Limited Dataset and an Application to Medical Imaging, 2019 IEEE Pacific Rim Conference on Communications, Computers and Signal Processing (PACRIM). IEEE. DOI:10.1109/pacrim47961.2019.8985057.
14. Silva, J. M., Pinho, E., Monteiro, E., Silva, J. F. and Costa, C. (2018) Controlled searching in reversibly de-identified medical imaging archives, *Journal of Biomedical Informatics*, 77, pp. 81–90. DOI:10.1016/j.jbi.2017.12.002.
15. Srihari, S. N., Shekhawat, A. en Lam, S. W. (2003) “Optical Character Recognition (OCR)”, in *Encyclopedia of Computer Science*. GBR: John Wiley and Sons Ltd., bll 1326–1333.
16. Zhou, S. K., Greenspan, H., Davatzikos, C., Duncan, J. S., Van Ginneken, B., Madabhushi, A., Prince, J. L., Rueckert, D. and Summers, R. M. (2021, May) A Review of Deep Learning in Medical Imaging: Imaging Traits, Technology Trends, Case Studies With Progress Highlights, and Future Promises, *Proceedings of the IEEE*. Institute of Electrical and Electronics Engineers (IEEE). DOI:10.1109/jproc.2021.3054390.