

RoundTripOCR: A Data Generation Technique for Enhancing OCR Error Correction in Low-Resource Devanagari Languages

Anonymous ACL submission

Abstract

Optical Character Recognition (OCR) technology has revolutionized the digitization of printed text, enabling efficient data extraction and analysis across various domains. Just like Machine Translation systems, OCR systems are prone to errors, stemming from factors such as poor image quality, diverse fonts, and language variations. In this work, we address the challenge of data generation and post-OCR error correction, specifically for low-resource languages. We propose a novel approach for synthetic data generation for Devanagari languages, RoundTripOCR, that tackles the scarcity of the OCR Error Correction dataset. In this work, we release a post-OCR text correction dataset for Hindi, Marathi, Bodo, Nepali, Konkani and Sanskrit. We also present a novel approach for OCR error correction by leveraging techniques from machine translation. Our method involves translating the erroneous OCR output into a corrected form by treating the OCR errors as mistranslations in a parallel text corpus. We employ a state-of-the-art pre-trained transformer model, mBART, to learn the mapping from erroneous to correct text pairs, effectively correcting OCR errors. The sample dataset can be accessed using the link¹.

1 Introduction

Devanagari script is the most widely used script in India and other Asian countries. There is a rich collection of ancient Devanagari manuscripts, which is a wealth of knowledge. To make these manuscripts available to people, efforts are being made to digitize these documents. Optical Character Recognition (OCR) technology has revolutionized the digitization and processing of written or printed text by enabling machines to automatically convert scanned documents into editable and searchable text formats. With the proliferation of document

digitization efforts across various domains such as finance, healthcare, education, and government, OCR plays a crucial role in enhancing document accessibility, information retrieval, and automation of document-intensive workflows. However, despite significant advancements in OCR technology over the years, the accurate recognition of text from scanned documents remains a challenging task due to inherent complexities in document layouts, font variations, noise, and other distortions.

Traditional OCR systems typically follow a pipeline approach comprising image preprocessing, feature extraction, character segmentation, and recognition stages. While these systems have achieved remarkable success in many applications, they are susceptible to errors, especially when dealing with degraded or low-quality document images. OCR errors can manifest in various forms, including misrecognitions, substitutions, omissions, and insertions, leading to inaccuracies in the recognized text output. These errors not only impede the reliability of OCR systems but also pose significant challenges for downstream tasks such as information extraction, text mining, and content analysis. Addressing OCR errors requires robust error detection and correction mechanisms that can effectively handle a wide range of error patterns and variations.

Our contributions are as follows:

1. **RoundTripOCR**, a novel approach to generate a Post-OCR error correction dataset artificially for Devanagari languages.
2. Post-OCR text correction dataset containing around 13.1 million sentences in Hindi, 1.58 million sentences in Marathi, 2.54 million sentences in Bodo, 2.97 million sentences in Nepali, 1.95 million sentences in Konkani and 8.93 million sentences in Sanskrit.
3. Benchmarks for the Post-OCR error correction task based on the pre-trained Seq2Seq

¹<https://drive.google.com/drive/folders/1EmXGHRHo2-hRqxTN8OwXJK-zdcX-1lYE?usp=sharing>

079 language model for all six languages.

080 2 Background

081 Addressing OCR errors is crucial for improving
082 the overall quality and usability of digitized docu-
083 ments.

084 2.1 Types of OCR Errors

085 As mentioned by Volk et al. (2011) and Jatowt et al.
086 (2019), OCR systems are prone to various types
087 of errors that can occur during the process of text
088 recognition from scanned documents. Understand-
089 ing these errors is essential for developing effective
090 error correction techniques and improving the over-
091 all accuracy of OCR systems. The most common
092 types of OCR errors include:

- 093 • **Substitution Errors:** Substitution errors oc-
094 cur when the OCR system incorrectly recog-
095 nizes a character and substitutes it with a
096 different character. These errors often result
097 from similarities between characters in terms
098 of shape or visual appearance, making it chal-
099 lenging for the OCR system to distinguish
100 between them accurately. For example, "o"
101 might be substituted for "0" or "l" for "1" in
102 the alphanumeric character set.
- 103 • **Omission Errors:** Omission errors occur
104 when the OCR system fails to detect and recog-
105 nize certain characters or words in the input
106 image. This can happen due to factors such
107 as poor image quality, low resolution, or in-
108 distinct boundaries between characters. Omis-
109 sion errors can significantly affect the read-
110 ability and integrity of the recognized text,
111 especially in documents with dense text or
112 complex layouts.
- 113 • **Insertion Errors:** Insertion errors occur when
114 the OCR system mistakenly inserts additional
115 characters or words that are not present in
116 the original image. These errors often arise
117 from misinterpretations of noise or artifacts
118 in the scanned document, leading to spurious
119 insertions in the recognized text. Insertion er-
120 rors can distort the meaning of the text and
121 introduce inconsistencies in downstream pro-
122 cessing tasks.
- 123 • **Deletion Errors:** Deletion errors occur when
124 the OCR system erroneously removes or

deletes characters or words from the input im- 125
age. These errors can occur due to segmenta- 126
tion errors, where the OCR system incorrectly 127
identifies boundaries between characters or 128
words, leading to the omission of valid text 129
fragments. Deletion errors can result in loss 130
of information and inaccuracies in the recog- 131
nized text output. 132

Over the years, researchers have explored var- 133
ious approaches to mitigate OCR errors, includ- 134
ing rule-based post-processing techniques (Khos- 135
robeigi et al., 2020), statistical language models 136
(Mei et al., 2018), and machine learning-based 137
methods (Virk et al., 2021). While these ap- 138
proaches have shown promise in certain scenarios, 139
they often rely on handcrafted rules or linguistic 140
resources, limiting their generalization to diverse 141
document types and languages. 142

In recent years, there has been growing interest in 143
applying advanced machine learning and natural 144
language processing techniques to address OCR er- 145
rors effectively. One promising direction is to lever- 146
age techniques from machine translation, which 147
aims to automatically translate text from one lan- 148
guage to another. By treating OCR errors as mis- 149
translations and modelling the correction process as 150
an automatic post-editing (APE) task, it is possible 151
to harness the power of neural machine translation 152
models to learn the mapping from erroneous to cor- 153
rect OCR text output. This paradigm shift not only 154
enables end-to-end error correction but also facili- 155
tates the integration of contextual information and 156
linguistic knowledge into the correction process, 157
leading to more accurate and robust OCR systems. 158
In the upcoming sections of the paper, we will look 159
at what are the different types of errors that exist 160
in the OCR output and how we will generate the 161
OCR error correction data and use it to mitigate the 162
OCR errors. 163

164 2.2 Machine Translation

Machine translation (MT) describes the automatic 165
process of translating text or speech from one lan- 166
guage to another utilizing algorithmic methods and 167
technology without the need for human translators. 168
As mentioned by Bhattacharyya (2015) MT's objec- 169
tive is to facilitate communication and understand- 170
ing amongst multilingual individuals by translating 171
written or spoken text automatically. MT systems 172
typically analyze the input text or speech using nat- 173
ural language processing (NLP) techniques, break- 174

ing it down into smaller linguistic units such as words, phrases, or sentences. These units are then processed and translated into the target language based on predefined rules, statistical models, or more advanced methods such as neural networks.

2.3 Automatic Post-editing and OCR Error Correction

Automatic Post-Editing (APE) is a methodology that utilizes various techniques to improve the quality of Machine Translation output automatically, including rule-based, statistical, and neural-based techniques (Vu and Haffari, 2018). APE systems are trained on human-edited translations, allowing them to identify and correct errors in grammar, fluency, and terminology. While MT systems have advanced significantly, they often produce translations that contain errors or lack fluency, especially with complex or domain-specific content. Output generated by a machine translation system is not always perfect and hence requires further editing (Parton et al., 2012). The task of editing a machine-translated text is referred to as post-editing (PE), which is time-consuming and costly. Hence, an efficient and automated system is required for post-editing. APE addresses these limitations by refining MT output, thereby enhancing the overall quality and usability of translations. APE can be easily integrated into existing machine translation workflows, making it a natural extension of the MT process.

While OCR systems play a crucial role in digitizing text, inherent limitations lead to errors in the extracted text. This necessitates post-processing techniques to refine the OCR output and achieve higher accuracy (Nguyen et al., 2021). Viewing this process through the lens of APE offers a valuable framework for developing effective error correction methods. Post-OCR error correction can be considered an Automatic Post Editing task. Similar to a machine translation system generating a translated sentence from a source language, the OCR system produces a "translated" text from the visual information in an image. This "translation" process is prone to errors due to limitations in both OCR systems, image quality, and stylistic variation. Just like an APE system refines a machine-translated sentence to improve fluency and accuracy, the post-OCR correction system aims to refine the text generated by the OCR system to remove errors and achieve a more accurate representation of the original document. Both MT and OCR error correction

face common challenges like handling ambiguity, dealing with rare words, and adapting to stylistic variations.

2.4 Round-trip translation

Information that is intentionally manufactured rather than derived from actual events is referred to as synthetic data. Synthetic data generation techniques are generally employed to generate artificial data for training machine learning models and neural networks.

Due to insufficient post-editing data available for the WMT APE 2016 shared task (Bojar et al., 2016) to train neural models, Junczys-Dowmunt and Grundkiewicz (2016), created two phrase-based translation models: English-German and German-English, using provided parallel training data to conduct round-trip translation. Using them in the Round-trip Translation approach resulted in the generation of artificial post-editing triplets $\langle src, mt, pe \rangle$. This artificial data creation method assisted in resolving the problem of insufficient training data, which frequently arises in NMT-based systems. Inspired by the Round-trip Translation approach, we propose a novel synthetic data generation technique, RoundTripOCR, which we discuss in detail in the following section.

3 RoundTripOCR: Data Generation Technique

The creation of artificial OCR data involves a systematic process aimed at simulating real-world scenarios keeping in consideration the common OCR error types and generating diverse datasets for training and evaluation purposes.

We use monolingual corpora taken from Technology Development for Indian Languages (TDIL)² and Maheshwari et al. (2022) to generate sentences of different lengths for generating artificial OCR data. To introduce variability and diversity into the dataset, 50 different Devanagari font combinations were selected from Google Fonts³. Each font style offered unique characteristics, such as varying stroke thickness, serif styles, and overall aesthetics as shown in Figure 1. Utilizing the selected Devanagari font combinations, 50 images could potentially be generated from a single sentence. PIL provides a comprehensive set of image processing functionalities, enabling the programmatic creation

²<https://www.tdil-dc.in>

³<https://fonts.google.com/?subset=devanagari>

Hindi Sentence: "चुनावो के बाद सरकार ने मुंबई में करों के माध्यम से अपने राजसव को बढ़ाया"
 Transliteration: chunaawo ke baad sarkar ne Mumbai me karoM ke maadhyam se apne raajaswa ko badhaayaa
 Gloss: Elections after government Mumbai in taxes through its revenue increased.
 Translation: After the elections, the government increased its revenue through taxes in Mumbai.

चुनावो के बाद सरकार ने मुंबई में करों के माध्यम से अपने राजसव को बढ़ाया

चुनावो के बाद सरकार ने मुंबई में करों के माध्यम से अपने राजसव को बढ़ाया

चुनावो के बाद सरकार ने मुंबई में करों के माध्यम से अपने राजसव को बढ़ाया

चुनावो के बाद सरकार ने मुंबई में करों के माध्यम से अपने राजसव को बढ़ाया

चुनावो के बाद सरकार ने मुंबई में करों के माध्यम से अपने राजसव को बढ़ाया

चुनावो के बाद सरकार ने मुंबई में करों के माध्यम से अपने राजसव को बढ़ाया

चुनावो के बाद सरकार ने मुंबई में करों के माध्यम से अपने राजसव को बढ़ाया

चुनावो के बाद सरकार ने मुंबई में करों के माध्यम से अपने राजसव को बढ़ाया

Figure 1: Examples of images generated with different fonts during RoundTripOCR data generation process.

of images with text rendered in specific font styles as described in the Appendix B. The generated images were subjected to optical character recognition (OCR) using the Pytesseract library as described in the Appendix C. Pytesseract is not supported for Bodo, Nepali, and Konkani languages. Thus, we used Pytesseract-Hindi for Bodo and Nepali, and Pytesseract-Marathi for Konkani due to similarities in these languages. We used Pytesseract-Sanskrit for the Sanskrit language. Pytesseract leverages machine-learning algorithms to extract text from images and convert them into machine-readable formats, including the Devanagari texts. The OCR process is aimed at simulating real-world OCR scenarios and generating text outputs from the rendered images. Since we can get 50 $\langle \text{Text } T, \text{OCR output } T' \rangle$ datapoints from a single sentence $\langle \text{Text } T \rangle$, this approach can be extended to any low-resource language.

By following this methodology, as shown in Figure 2, a comprehensive artificial dataset for OCR error detection and correction was generated, encompassing a diverse range of text passages, font styles, and linguistic variations. This dataset serves as a valuable resource for training and evaluating OCR systems, enabling researchers and practitioners to develop robust OCR algorithms and assess

their performance under various conditions.

3.1 Dataset

Leveraging the RoundTripOCR technique, we generate datasets containing around 13.1 million sentences in Hindi, 1.58 million sentences in Marathi, 2.54 million sentences in Bodo, 2.97 million sentences in Nepali, 1.95 million sentences in Konkani and 8.93 million sentences in Sanskrit. To assess the quality of the artificial dataset, we provided 100 data samples from the Hindi and Marathi datasets to language and OCR experts for evaluation. Their feedback indicated that the dataset contains errors very similar to those found in real OCR outputs. We have observed that the quality of training data obtained either from real-world data or by our RoundTripOCR technique is comparable. Thus, we believe that using our technique is a convenient and appropriate way to generate a synthetic dataset to train OCR error correction systems.

Font analysis revealed significant variations in error rates. Specifically, fonts such as Khand-Regular, Rajdhani-Regular, Nirmala, and Biryani exhibited the highest CERs, exceeding 7%. Conversely, fonts like Gargi, Karma-Regular, NotoSans-Regular, and VesperLibre-Regular demonstrated remarkably low CERs, each falling below 1% as shown in Figure

# of Sent.	Hindi	Marathi	Bodo	Nepali	Konkani	Sanskrit
Train dataset	13,129,200	1,581,405	2,541,649	2,970,148	1,950,874	8,935,790
Validation set	10,000	10,000	10,000	10,000	10,000	10,000
Test dataset	10,000	10,000	10,000	10,000	10,000	10,000

Table 1: Dataset distribution for different languages

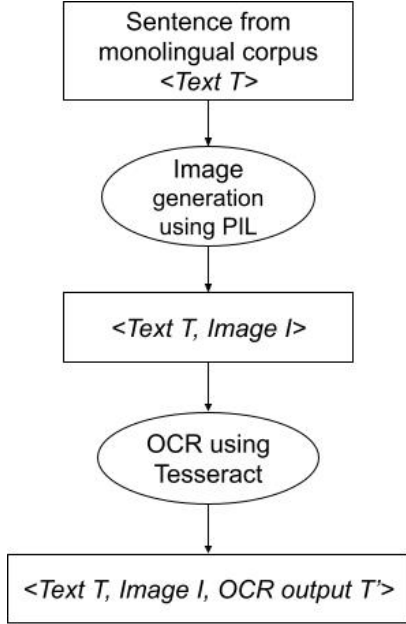


Figure 2: RoundTripOCR: The process of Artificial OCR Data Generation. At the end of the RoundTripOCR process, $\langle OCR\ output\ T' \rangle$ will act as OCR output, and $\langle Text\ T \rangle$ will act as corrected OCR output text.

3. Our findings suggest that models trained on a diverse range of fonts perform more robustly than those trained solely on a single font. This observation underscores the importance of font diversity in enhancing OCR error correction models’ performance and resilience. The sample dataset can be accessed using the link⁴.

4 Experiments and Results

4.1 Fine-tuning sequence to sequence model using our dataset

In this experiment, we harnessed a vast dataset consisting of 13 million data points, curated through our RoundTripOCR data generation methodology. Leveraging this extensive corpus, we conducted a series of experiments employing a sophisticated sequence-to-sequence model: *mBART*.

⁴<https://drive.google.com/drive/folders/1EmXGHRHo2-hRqxTN8OwXJK-zdcX-i1YE?usp=sharing>

To facilitate effective model training and evaluation, we partitioned the dataset into three distinct sets: training, testing, and validation. Specifically, the testing set contained 10,000 pairs, while the validation set comprised 10,000 pairs. This meticulous partitioning strategy enabled us to assess the performance of our models accurately and reliably across various metrics. We fine-tuned the model for 4 epochs and 3000 max_steps.

mBART (Multilingual BART) represents an extension of the BART architecture tailored specifically for multilingual text processing as explained in the Appendix D.

To facilitate an in-depth investigation, we further curated two distinct datasets. The first dataset encompassed text samples spanning various fonts, providing a rich diversity in font styles. Conversely, the second dataset exclusively features a single font style; particularly we chose the Sumana font as it shows a close to average CER when compared with all the fonts used in the creation of the dataset. This deliberate bifurcation allowed us to explore the potential advantages conferred by employing data with varying font styles, thereby enriching our understanding of the model’s performance under different font conditions.

4.2 Results

Results on Hindi test dataset		
Model	CER	WER
Tesseract (baseline)	2.247%	5.833%
mBART (single font)	2.105%	5.817%
mBART (all fonts)	1.556%	3.474%

Table 2: Results on Hindi test dataset

Results on Marathi test dataset		
Model	CER	WER
Tesseract (baseline)	4.100%	15.372%
mBART (single font)	3.592%	14.472%
mBART (all fonts)	2.464%	9.886%

Table 3: Results on Marathi test dataset

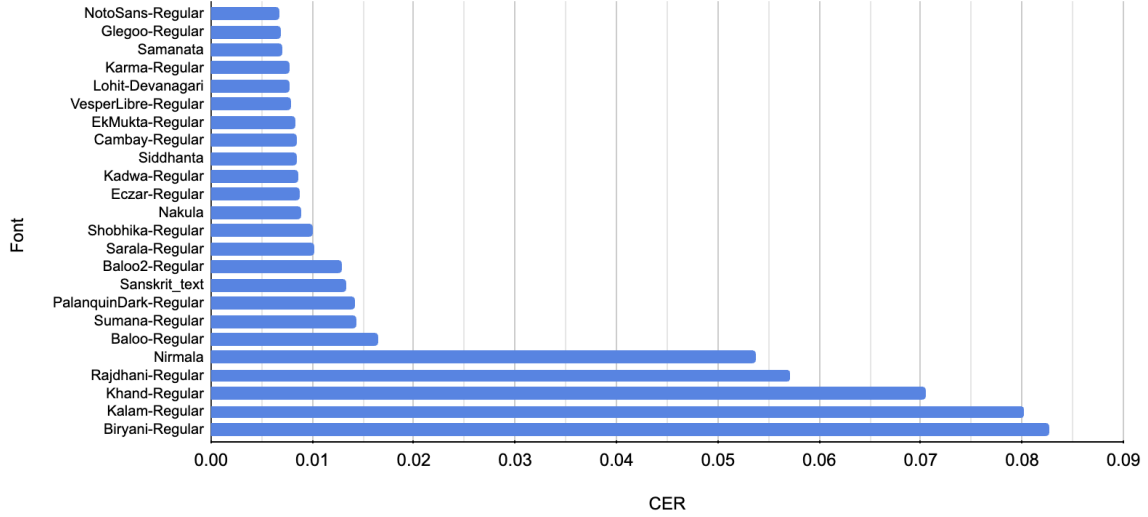


Figure 3: Comparison of different fonts and their CER in the Hindi test dataset

Results on Konkani test dataset		
Model	CER	WER
Tesseract (baseline)	4.222%	16.803%
mBART (single font)	3.279%	13.058%
mBART (all fonts)	2.265%	8.524%

Table 4: Results on Konkani test dataset

Results on Nepali test dataset		
Model	CER	WER
Tesseract (baseline)	5.783%	24.287%
mBART (single font)	3.189%	14.266%
mBART (all fonts)	2.387%	10.651%

Table 5: Results on Nepali test dataset

Results on Bodo test dataset		
Model	CER	WER
Tesseract (baseline)	5.893%	24.027%
mBART (single font)	3.677%	13.016%
mBART (all fonts)	2.359%	6.822%

Table 6: Results on Bodo test dataset

Results on Sanskrit test dataset		
Model	CER	WER
Tesseract (baseline)	8.766%	44.730%
mBART (single font)	6.428%	29.294%
mBART (all fonts)	5.670%	25.501%

Table 7: Results on Sanskrit test dataset

The evaluation of OCR model performance across multiple Indian languages reveals notable

insights. In the Hindi and Marathi test datasets, Tesseract (baseline) exhibits CERs of 2.247% and 4.100%, respectively, with corresponding WERs of 5.833% and 15.372%. Contrastingly, mBART (all fonts) achieves significantly lower error rates, recording CERs of 1.556% and 2.464% and WERs of 3.474% and 9.886% for Hindi and Marathi, respectively. Similar trends are observed in the Konkani and Nepali datasets, where mBART (all fonts) maintains a competitive edge with CERs of 2.265% and 2.387% and WERs of 8.524% and 10.651%. In more challenging languages like Bodo and Sanskrit, mBART (all fonts) continues to demonstrate superiority over the Tesseract baseline, achieving CERs of 2.359% and 5.670% and WERs of 6.822% and 25.501%, respectively. These results underscore the effectiveness of mBART, particularly when trained on diverse font variations, in enhancing OCR accuracy across a range of complex languages.

5 Conclusion and Future Work

In our work, we introduced a novel approach for OCR error correction data generation and created a vast dataset comprising 13.1 million sentences in Hindi, 1.58 million sentences in Marathi, 2.54 million sentences in Bodo, 2.97 million sentences in Nepali, 1.95 million sentences in Konkani and 8.93 million sentences in Sanskrit. Notably, our proposed methodology is versatile and can be extended to other low-resource languages that follow the Devanagari script. By leveraging monolingual corpora, our approach enables the generation

of OCR correction datasets, thus addressing the scarcity of data in such languages.

The findings from our experimentation underscore the efficacy of approaches from Machine Translation for the task of OCR output correction, specifically state-of-the-art models like mBART, trained on diverse datasets to substantially enhance OCR accuracy. By improving the accuracy of OCR systems, our research contributes to making textual content more accessible and usable, thereby facilitating broader access to information and knowledge in multilingual societies. Our findings suggest that models trained on a diverse range of fonts perform more robustly than those trained solely on a single font. This observation underscores the importance of font diversity in enhancing OCR error correction models' performance and resilience.

Our findings motivate the exploration of data augmentation techniques utilizing synthetically generated Devanagari script images. By incorporating these images with controlled variations in font styles, noise levels, and image degradations, we can investigate the impact on model generalization and robustness towards real-world document image complexities. This research could delve into generative adversarial networks (GANs) or other image synthesis techniques to create a diverse and realistic training dataset for enhanced OCR performance.

We propose the experimental findings in this work as a baseline, based on which future work can focus on novel and sophisticated techniques for the task of OCR correction, including improvements to the architecture.

6 Limitations

Our work focuses on improving OCR error correction for Devanagari script languages only. Extending this approach to achieve true multilingual OCR is a complex endeavour. Different languages possess unique linguistic characteristics, script variations, and language-specific nuances. Developing a single model capable of handling this vast diversity effectively remains a challenge. Future work should explore techniques for creating language-agnostic or language-adaptive models to address these limitations and achieve broader multilingual OCR applicability.

7 Ethical Statement

This research utilizes datasets that are openly available in the public domain. The data employed for

generating artificial data in this study was sourced from publicly accessible repositories, ensuring that there are no privacy or ethical concerns associated with their use. No user information was present in any of the datasets used in the work, protecting the privacy and identity of users. Also, the synthetic data generated as a part of this work will be released under the CC-BY-SA 4.0 license publicly for further research. We understand that every dataset is subject to intrinsic bias and that computational models will inevitably learn biased information from any dataset.

References

- Pushpak Bhattacharyya. 2015. *Machine translation*. CRC Press.
- Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurélie Névéol, Mariana Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor, and Marcos Zampieri. 2016. [Findings of the 2016 conference on machine translation](#). In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 131–198, Berlin, Germany. Association for Computational Linguistics.
- Alex Clark et al. 2015. Pillow (pil fork) documentation. *readthedocs*.
- Adam Jatowt, Mickael Coustaty, Nhu-Van Nguyen, Antoine Doucet, et al. 2019. Deep statistical analysis of ocr errors for effective post-ocr processing. In *2019 ACM/IEEE Joint Conference on Digital Libraries (JCDL)*, pages 29–38. IEEE.
- Marcin Junczys-Dowmunt and Roman Grundkiewicz. 2016. Log-linear combinations of monolingual and bilingual neural machine translation models for automatic post-editing. *arXiv preprint arXiv:1605.04800*.
- Z. Khosrobeigi, H. Veisi, H.R. Ahmadi, and H. Shabani. 2020. [A rule-based post-processing approach to improve persian ocr performance](#). *Scientia Iranica*, 27(6):3019–3033.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and

505	Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation . <i>Transactions of the Association for Computational Linguistics</i> , 8:726–742.	559
506		560
507		561
508		562
509	Ayush Maheshwari, Nikhil Singh, Amrith Krishna, and Ganesh Ramakrishnan. 2022. A benchmark and dataset for post-OCR text correction in Sanskrit . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 6258–6265, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	563
510		564
511		565
512		566
513		567
514		568
515		569
516	Jie Mei, Aminul Islam, Abidalrahman Moh'd, Yajing Wu, and Evangelos Milios. 2018. Statistical learning for ocr error correction . <i>Information Processing and Management</i> , 54:Pages 874–887.	570
517		571
518		572
519		573
520	Thi Nguyen, Adam Jatowt, Mickaël Coustaty, and Antoine Doucet. 2021. Survey of post-ocr processing approaches . <i>ACM Computing Surveys</i> , 54:1–37.	574
521		575
522		576
523	Kristen Parton, Nizar Habash, Kathleen McKeown, Gonzalo Iglesias, and Adrià de Gispert. 2012. Can automatic post-editing make MT more meaningful . In <i>Proceedings of the 16th Annual conference of the European Association for Machine Translation</i> , pages 111–118, Trento, Italy. European Association for Machine Translation.	577
524		578
525		579
526		580
527		581
528		582
529		583
530	Ray Smith. 2007. An overview of the tesseract ocr engine. In <i>Ninth international conference on document analysis and recognition (ICDAR 2007)</i> , volume 2, pages 629–633. IEEE.	584
531		585
532		586
533		587
534	Shafqat Mumtaz Virk, Dana Dannélls, and Azam Sheikh Muhammad. 2021. A novel machine learning based approach for post-OCR error detection . In <i>Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)</i> , pages 1463–1470, Held Online. INCOMA Ltd.	588
535		589
536		590
537		591
538		592
539		593
540		594
541	Martin Volk, Lenz Furrer, and Rico Sennrich. 2011. Strategies for reducing and correcting ocr errors. In <i>Language Technology for Cultural Heritage: Selected Papers from the LaTeCH Workshop Series</i> , pages 3–22. Springer.	595
542		596
543		597
544		598
545		599
546	Thuy-Trang Vu and Gholamreza Haffari. 2018. Automatic post-editing of machine translation: A neural programmer-interpreter approach . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3048–3053, Brussels, Belgium. Association for Computational Linguistics.	600
547		601
548		602
549		603
550		604
551		605
552		606
553		607
554	A Appendix	
555	The fonts used to generate images are: EkMukta-Regular, Arya-Regular, YatraOne-Regular, Siddhanta, Sura-Regular, Samanata, Karma-Regular, Nirmala, Asar-Regular, VesperLibre-Regular, Kurale-Regular, MartelSans-Regular,	
556		
557		
558		
	SakalBharati Normal, Biryani-Regular, Sumana-Regular, Sarai, Laila-Regular, Rajdhani-Regular, Nakula, Shobhika-Regular, Baloo-Regular, Lohit-Devanagari, Amiko-Regular, Akshar Unicode, Palanquin-Regular, Eczar-Regular, Glegoo-Regular, Mukta-Regular, Sanskrit2003, PalanquinDark-Regular, Baloo2-Regular, Kalam-Regular, Sanskrit_text, Halant-Regular, Hind-Regular, Cambay-Regular, PragatiNarrow-Regular, Kadwa-Regular, Kokila, Sahadeva, Utsaah, Sahitya-Regular, Khand-Regular, Sarala-Regular, NotoSans-Regular, Jaldi-Regular, RhodiumLibre-Regular, Yantramanav-Regular and Gargi.	
	Languages which follow devnagari script are: Apabhramsha, Angika, Awadhi, Bajjika, Bhili, Bhojpuri, Boro, Braj, Chhattisgarhi, Dogri, Garhwali, Haryanvi, Hindi, Kashmiri, Khandeshi, Konkani, Kumaoni, Magahi, Maithili, Marathi, Marwari, Mundari, Nagpuri, Newari, Nepali, Pāli, Pahari, Prakrit, Rajasthani, Sanskrit, Santali, Saraiki, Sherpa, Sindhi, Surjapuri, and many more. We can extend our approach to all these languages.	
	B PIL (Python Image Library)	
	In today’s digital era, the prevalence of digital images is ubiquitous. When working with Python, developers have access to a plethora of image-processing libraries to augment the capabilities of digital images. Among the most widely used libraries are OpenCV, Python Imaging Library (PIL), Scikit-image, and Pillow.	
	Pillow ⁵ , an extension of PIL (Python Image Library), stands out as a crucial module for image processing in Python. While PIL was once pivotal, it ceased support in 2011 and does not cater to Python 3. In contrast, Pillow offers expanded functionalities, runs seamlessly across major operating systems, and fully supports Python 3. It boasts compatibility with a diverse range of image formats, including "jpeg", "png", "bmp", "gif", "ppm", and "tiff". With Pillow, developers can perform a multitude of operations on digital images, ranging from basic tasks like point operations to advanced functionalities such as image filtering using built-in convolution kernels and colour space conversions.	
	The Python Imaging Library (Clark et al., 2015), commonly known as PIL, is particularly well-suited for image archival and batch-processing applications. Leveraging the Python Pillow package,	
	⁵ https://pypi.org/project/pillow	

608 developers can seamlessly perform tasks such as
609 creating thumbnails, converting between different
610 image formats, and printing images. The Pillow
611 library encompasses a comprehensive suite of basic
612 image-processing functionalities, including image
613 resizing, rotation, and transformation. Additionally,
614 the histogram method available in the Pillow mod-
615 ule facilitates the extraction of statistical data from
616 images, which can then be utilized for statistical
617 analysis and automatic contrast enhancement.

618 **C Tesseract**

619 Tesseract⁶ is an open-source OCR Engine designed
620 to extract printed or handwritten text from images.
621 Originally developed by Hewlett-Packard, its de-
622 velopment was later taken over by Google (Smith,
623 2007). Tesseract boasts support for language recog-
624 nition in over 100 languages straight out of the
625 box. The latest iteration, Tesseract 4.0, features
626 AI integration through LSTM Neural Networks,
627 enhancing its capability to detect and recognize
628 inputs of varying sizes with greater accuracy and
629 efficiency.

630 Tesseract’s versatility lies in its compatibility
631 with various programming languages and frame-
632 works through wrappers like Pytesseract, com-
633 monly known as Python-Tesseract. This tool not
634 only serves as an open-source OCR library for
635 Python but also acts as a wrapper for Google’s
636 Tesseract OCR Engine. Pytesseract offers the con-
637 venience of being a standalone script, enabling di-
638 rect printing of recognized text without the need
639 to convert it into a separate file. It supports a wide
640 range of image formats, including JPEG, PNG, GIF,
641 BMP, TIFF, and more, making it a popular choice
642 for image-to-text OCR tasks in Python.

643 Pytesseract played an important role in our work
644 of data generation, where it was used for perform-
645 ing the OCR of the images that we generated using
646 PIL.

647 **D mBART (Multilingual BART)**

648 Bidirectional and Auto-Regressive Transformers
649 (BART) is a sequence-to-sequence model pro-
650 posed by Lewis et al. (2020) that combines the
651 strengths of bidirectional and auto-regressive trans-
652 formers. Its architecture consists of an encoder-
653 decoder transformer model with masked self-
654 attention mechanism in the decoder.

655 Encoder: The encoder processes the input se-
656 quence bidirectionally using self-attention mecha-
657 nisms to capture contextual information efficiently.
658 It produces contextualized representations of the
659 input tokens.

660 Decoder: The decoder generates the output se-
661 quence autoregressively, conditioning on the en-
662 coder’s contextualized representations and previ-
663 ously generated tokens. It employs causal self-
664 attention mechanisms to ensure that each token is
665 generated based only on previously generated to-
666 kens, as opposed to looking ahead of the current
667 decoding step.

668 BART is pre-trained on large-scale text corpora
669 using denoising autoencoding objectives, where
670 corrupted input sequences are reconstructed to their
671 original forms. This pre-training objective encour-
672 ages the model to learn robust representations of
673 text and enables it to perform well on a wide range
674 of natural language processing tasks.

675 mBART (Multilingual BART) by Liu et al.
676 (2020) extends the BART architecture to support
677 multilingual text processing. It is designed to han-
678 dle input sequences in multiple languages and gen-
679 erate output sequences in the corresponding target
680 languages.

681 Language Embeddings: mBART incorporates
682 language embeddings into its architecture to enable
683 language-specific processing. These embeddings
684 encode information about the source and target lan-
685 guages, allowing the model to adapt its behaviour
686 based on the language of the input and output se-
687 quences.

688 Cross-lingual Pre-training: mBART is pre-
689 trained on multilingual text corpora using cross-
690 lingual objectives, where input sequences from dif-
691 ferent languages are reconstructed to their original
692 forms. This pre-training objective encourages the
693 model to learn language-agnostic representations
694 of text and enables it to perform effectively on mul-
695 tilingual tasks.

⁶<https://github.com/tesseract-ocr/tesseract>