

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

Harshvivek Kashid and Pushpak Bhattacharyya

Indian Institute of Technology Bombay

{harshvivek,pb}@cse.iitb.ac.in

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. In this work, we address the challenge of data generation and post-OCR error correction, specifically for low-resource languages. We propose an approach for synthetic data generation for Devanagari languages, **RoundTripOCR**, that tackles the scarcity of the post-OCR Error Correction datasets for low-resource languages. We release post-OCR text correction datasets 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 erroneous OCR output into a corrected form by treating the OCR errors as mis-translations in a parallel text corpus, employing pre-trained transformer models to learn the mapping from erroneous to correct text pairs, effectively correcting OCR errors.

1 Introduction

The Devanagari script is the most extensively used writing system in the Indian subcontinent. It was the principal script for Sanskrit, the ancient literary language of Indian civilization. Sanskrit was used to write a wide range of texts covering various domains, including literature, philosophy, science, art, architecture, and mathematics. This includes the Vedas, Upanishads, and epics like Mahabharata and Ramayana. Devanagari script originated from ancient Brahmi script through various transformations (Jayadevan et al., 2011). Apart from vowels, modifiers and consonants (Figure 1), it has a rich set of conjunct consonants, known as ligatures, where multiple characters combine to form new glyphs. These are difficult to segment and recognize because they don't correspond directly to

अ (a)	आ (ā)	इ (i)	ई (ī)	उ (u)	ऊ (ū)	ऋ (ṛ)
ए (e)	ऐ (ai)	ओ (o)	औ (au)	अं (am)	अः (ah)	
Vowels						
ा (ā)	ि (i)	ी (ī)	ु (u)	ू (ū)	ृ (ṛ)	ै (e)
ै (ai)	ो (o)	ौ (au)	ं (m)	ः (h)	ँ (m)	
Modifiers						
क (ka)	ख (kha)	ग (ga)	घ (gha)	ङ (ṅa)		
च (ca)	छ (cha)	ज (ja)	झ (jha)	ञ (ña)		
ट (ṭa)	ठ (ṭha)	ड (ḍa)	ढ (ḍha)	ण (ṇa)		
त (ta)	थ (tha)	द (da)	ध (dha)	न (na)		
प (pa)	फ (pha)	ब (ba)	भ (bha)	म (ma)		
य (ya)	र (ra)	ल (la)	व (va)	श (śa)		
ष (ṣa)	स (sa)	ह (ha)	ळ (ḷa)			
Consonants						

Figure 1: Vowels, modifiers and consonants of Devanagari script.

individual letters. Devanagari characters often have vowel signs (*matras*) and other diacritical marks that appear above, below, or beside the base character. These modifiers can be challenging to detect, segment, and associate correctly with the base character. Devanagari script includes a horizontal line (called the *Shirokekha*) that connects the characters in each word. Unlike in Latin scripts, where spaces clearly divide words, in Devanagari, words often connect via the headline, making word segmentation harder for OCR systems. Proper segmentation of Devanagari words, characters, and sub-components (such as vowels and consonants) is difficult because components often overlap, connect through ligatures, or blend with the *Shirokekha* line. This is less common in simpler scripts like Latin, where individual letters are often spaced apart and stand independently. Many Devanagari characters look quite similar, especially in certain fonts or degraded images, leading to higher chances of OCR errors. OCR technology has revolutionized the digitization and processing of written or printed text by enabling machines to automatically

convert scanned documents or handwritten texts into editable and searchable text formats. However, despite significant advancements 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 machine translation (Kolak et al., 2003; Laique et al., 2021; Nguyen et al., 2021b; Ignat et al., 2022). 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:

1. **RoundTripOCR**¹, a technique to artificially generate post-OCR error correction data for low-resource Devanagari script languages in the form $\langle T, T' \rangle$, where T' is the OCR output text and T is the correct OCR output text (Section 3).
2. Post-OCR error correction dataset, containing 3.1 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 4.07 million sentences in Sanskrit (Table 1).
3. Benchmarks for the Post-OCR error correction task based on the pre-trained Seq2Seq language models for all six languages (Section 5.2).

2 Related work

As mentioned by Volk et al. (2011) and Jatowt et al. (2019), OCR systems are prone to various types

of errors that can occur during the process of text recognition from scanned documents. The most common types of OCR errors include: substitution errors, insertion errors, deletion errors and segmentation errors.

Even state-of-the-art OCR models are susceptible to making recognition errors (Dong and Smith, 2018). Errors are particularly frequent in the case of low-resource languages because most off-the-shelf OCR tools do not directly support these languages, and training a high-performance OCR system is challenging given the small amount of data that is typically available (Rijhwani et al., 2020). We use post-OCR error correction tools and techniques to correct these errors and improve the quality of the transcription. Over the years, researchers have explored various approaches to mitigate OCR errors, including rule-based post-processing techniques (Khosrobeigi et al., 2020), statistical language models (Mei et al., 2018), and machine learning-based methods (Virk et al., 2021). While these approaches have shown promise in certain scenarios, they often rely on handcrafted rules or linguistic resources, limiting their generalization to diverse document types and languages.

In recent years, there has been growing interest in applying advanced machine learning and natural language processing techniques to address OCR errors effectively. One promising direction is to leverage techniques from machine translation, which aims to automatically translate text from one language to another (Lyu et al., 2021). By treating OCR errors as mistranslations and modelling the correction process as an automatic post-editing (APE) task, it is possible to harness the power of neural machine translation models to learn the mapping from erroneous to correct OCR text output. This paradigm shift not only enables end-to-end error correction but also facilitates the integration of contextual information and linguistic knowledge into the correction process, leading to more accurate and robust OCR systems.

The emergence of Transformer architecture and attention mechanisms (Vaswani et al., 2023) has led to the adoption of deep learning models for post-OCR tasks. Post-OCR tasks have been reframed as Sequence-to-Sequence tasks in recent studies, whereby researchers have applied Machine Translation models (Amrhein and Clematide, 2018). The BERT (Devlin et al., 2019) and BART (Lewis et al., 2020) models were used by Nguyen et al. (2020)

¹RoundTripOCR code and dataset details are on GitHub: <https://github.com/harshvivek14/RoundTripOCR>

and Soper et al. (2021), respectively. Maheshwari et al. (2022) compared standard Sequence-to-Sequence models with pre-trained models. A lot of data is needed to train these models, and the predominant method for obtaining post-OCR training data has been crowdsourcing (Clematide et al., 2016). Although this can yield extremely accurate training data, the procedure often proves costly and time-consuming. Thus, synthetic data generation has been widely employed in this task (D’hondt et al., 2017; Jasonarson et al., 2023; Guan and Greene, 2024). The sentence or line-level OCR error correction by using the sentence or line-level dataset has also proven to be effective in addressing segmentation and word errors of the OCR output (Dwivedi et al., 2020; Lyu et al., 2021; Rijhwani et al., 2021; Ignat et al., 2022; Thomas et al., 2024).

2.1 Automatic Post-Editing and OCR Error Correction

Automatic Post-Editing (APE) uses techniques to improve the quality of Machine Translation (MT) output automatically, including rule-based, statistical, and neural-based techniques (Chollampatt et al., 2020; Deoghare et al., 2023). 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; Läubli et al., 2013; Pal et al., 2016).

OCR systems play a crucial role in digitizing text, but 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., 2021a). 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 text present in an image. This process is prone to errors due to limitations in 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.

2.2 Round-trip translation

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 $\langle src, mt, pe \rangle$, where *src* is source sentence, *mt* is machine translated sentence and *pe* is post-edited sentence. 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 and image-based synthetic data generation technique for the OCR system by Etter et al. (2019), which promises unlimited training data at zero annotation cost, we propose a synthetic data generation technique for post-OCR error correction, **RoundTripOCR**, which we discuss in detail in the following section.

3 RoundTripOCR

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

To introduce variability into the dataset, 50 different Devanagari font combinations were selected from Google Fonts². Each font style offered unique characteristics, such as varying stroke thickness, serif styles, and overall aesthetics, as shown in Figure 3. 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 of images with

²<https://fonts.google.com/?subset=devanagari>

# of Sent.	Hindi	Marathi	Bodo	Nepali	Konkani	Sanskrit
Train dataset	3,129,200	1,581,405	2,541,649	2,970,148	1,950,874	4,070,000
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: Distribution of dataset generated using RoundTripOCR technique for all six languages.

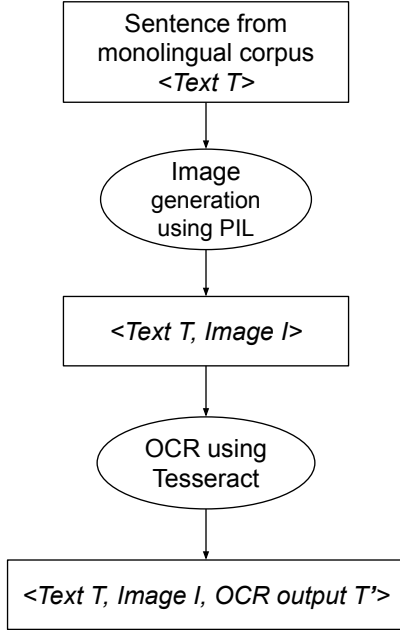


Figure 2: **RoundTripOCR**: Artificial post-OCR error correction data generation process. We get $\langle \text{Text } T, \text{Image } I, \text{OCR output } T' \rangle$ as output, where $\langle \text{Text } T \rangle$ will be used as corrected OCR output text and $\langle \text{OCR output } T' \rangle$ as OCR output.

text rendered in specific font styles. The generated images were subjected to optical character recognition (OCR) using the Pytesseract library. Pytesseract is not supported for Bodo, Nepali, and Konkani languages. Thus, we use Pytesseract-Hindi for Bodo and Nepali along with Hindi and Pytesseract-Marathi for Konkani and Marathi 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 $\langle \text{Text } T, \text{OCR output } T' \rangle$ data points 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 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

We generate post-OCR error correction datasets for Bodo, Nepali, Konkani, Hindi, Sanskrit and Marathi texts. The corpora for Hindi was sourced from the CC-100 corpus (Conneau et al., 2020), and Konkani, Nepali, Bodo and Marathi texts were sourced from Technology Development for Indian Languages (TDIL)³ and Sanskrit texts were sourced from Maheshwari et al. (2022). Leveraging the RoundTripOCR technique, we generate datasets containing around 3.1 million sentence pairs in Hindi, 1.58 million sentence pairs in Marathi, 2.54 million sentence pairs in Bodo, 2.97 million sentence pairs in Nepali, 1.95 million sentence pairs in Konkani and 4.07 million sentence pairs in Sanskrit as mentioned in the Table 1. Each pair have $\langle \text{Text } T \rangle$, which is the corrected OCR output text, and $\langle \text{OCR output } T' \rangle$, which is the OCR output sentence⁴.

4 Sequence to Sequence models

We conducted a series of experiments employing sequence-to-sequence models: **mBART**, **mT5** and **IndicBART**. These are powerful models designed for multilingual tasks, particularly in low-resource languages.

mBART (Multilingual BART) is a sequence-to-sequence denoising autoencoder that pre-trains on a variety of languages by corrupting and reconstructing text, making it highly effective for tasks like machine translation and text generation across different languages (Liu et al., 2020). It has been extensively used for tasks involving noisy inputs,

³<https://www.tdil-dc.in>

⁴Datasets are available at: <https://github.com/harshvivek14/RoundTripOCR>

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

such as post-OCR error correction, due to its ability to learn contextual representations and perform cross-lingual transfer (Soper et al., 2021; Maheshwari et al., 2022). We used *mbart-large-50* version of mBART.

mT5 (Multilingual T5) extends the T5 model’s text-to-text framework to a massively multilingual setting (Xue et al., 2021). With the capacity to handle over 100 languages, mT5 is effective for multilingual NLP tasks, including translation, summarization, and post-OCR error correction (Madarász et al., 2024). This model leverages the original T5 framework, where every NLP task is reframed as a text generation problem, allowing for consistent and flexible handling of a wide range of tasks across languages. The version we use in our experiments is *mT5-base*.

IndicBART is a variant of mBART that is specifically tailored for Indic languages like Hindi, Bengali, Marathi, and others (Dabre et al., 2022). It adapts the pre-training and fine-tuning processes to better handle the linguistic and scriptural characteristics of these languages, which are often underrepresented in large-scale language models. IndicBART has proven to be highly effective for tasks such as machine translation in Indic scripts.

5 Experiments and Results

The pre-trained models were sourced from Hugging Face⁵ and finetuned using NVIDIA A100 GPU for 2 to 3 epochs. A learning rate of $5e-4$ was applied, managed by a polynomial learning rate scheduler. The training was conducted with 32-bit floating-point precision, and the best-performing model from each run was saved for evaluation. To facilitate effective model training and evaluation, we partitioned the dataset into training, testing, and validation sets. The testing set and validation set contained 10,000 pairs each.

We further curated the second dataset exclusively featuring a single font style; in particular, 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 as shown in Figure 4. This 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.

⁵<https://huggingface.co/models>

Model	Hindi		Marathi		Konkani		Nepali		Bodo		Sanskrit	
	CER	WER	CER	WER	CER	WER	CER	WER	CER	WER	CER	WER
OCR (<i>Tesseract</i>)	2.25%	5.83%	4.10%	15.37%	4.22%	16.80%	5.78%	24.29%	5.89%	24.03%	8.77%	44.73%
IndicBART (<i>single</i>)	2.30%	5.65%	4.23%	15.04%	3.78%	14.64%	5.56%	22.94%	4.69%	16.69%	6.84%	32.31%
IndicBART (<i>all fonts</i>)	2.19%	5.33%	4.08%	12.95%	3.51%	12.70%	4.04%	15.04%	3.84%	12.65%	6.38%	31.89%
mT5 (<i>single</i>)	2.08%	5.50%	3.65%	15.01%	3.81%	15.38%	3.88%	16.51%	4.45%	16.99%	6.57%	30.23%
mT5 (<i>all fonts</i>)	1.95%	4.88%	2.91%	10.51%	3.13%	12.95%	3.37%	14.29%	4.20%	15.53%	6.41%	29.32%
mBART (<i>single</i>)	2.11%	5.82%	3.59%	14.47%	3.28%	13.06%	3.19%	14.27%	3.68%	13.02%	6.43%	29.29%
mBART (<i>all fonts</i>)	1.56%	3.47%	2.46%	9.89%	2.27%	8.52%	2.39%	10.65%	2.36%	6.82%	5.67%	25.50%

Table 2: Comparison of mBART, mT5, and IndicBART for Hindi, Marathi, Konkani, Nepali, Bodo, and Sanskrit test datasets based on CER and WER metrics. Tesseract OCR is the baseline. Models for which training is done using a single font style data are indicated as: *single*. Models trained on data with all fonts are indicated as *all fonts*. The best results are highlighted in bold.

5.1 Evaluation metric

In OCR error correction, performance is commonly measured using *Character Error Rate* (CER) and *Word Error Rate* (WER). Both metrics evaluate the edit distance between predicted and ground truth text.

CER is defined as:

$$CER = \frac{Sc + Dc + Ic}{N}$$

where Sc , Dc , and Ic are the number of character-level substitutions, deletions, and insertions, respectively, and N is the total number of characters in the reference text.

WER is similarly defined at the word level:

$$WER = \frac{Sw + Dw + Iw}{W}$$

Sw , Dw , and Iw are the number of word-level substitutions, deletions, and insertions, respectively, and W is the total number of words in the reference text. Lower CER and WER indicate better OCR error correction performance.

5.2 Results

We evaluated the performance of several models, IndicBART, mT5, and mBART, on six languages: Hindi, Marathi, Konkani, Nepali, Bodo, and Sanskrit. The models were assessed using two metrics: CER and WER. Tesseract output was considered as the baseline. Across all languages, mBART (*all fonts*) consistently outperformed other models, showing the lowest CER and WER, followed by mT5 (*all fonts*). We present detailed results of all conducted experiments in Table 2 comparing the finetuned models with the baseline in the test dataset.

For instance, in Hindi, Tesseract recorded a CER of 2.25% and a WER of 5.83%, whereas the neural models significantly reduced the errors. Among

them, *mBART (all fonts)* consistently demonstrated the best performance with a CER of 1.56% and a WER of 3.47%. Similar trends were observed in Marathi, where Tesseract had a CER of 4.10% and a WER of 15.37%, while *mBART (all fonts)* outperformed with a CER of 2.46% and WER of 9.89%.

In Konkani, Tesseract’s error rates were even higher, with a CER of 4.22% and a WER of 16.80%. However, *mBART (all fonts)* again achieved the best results with a CER of 2.27% and a WER of 8.52%, illustrating its robust performance across different scripts. Nepali, being another challenging language for OCR, saw a high error rate from Tesseract (CER of 5.78% and WER of 24.29%), but *mBART (all fonts)* reduced these errors to a CER of 2.39% and WER of 10.65%. For Bodo, Tesseract recorded a CER of 5.89% and WER of 24.03%, while *mBART (all fonts)* again provided substantial improvements, bringing the CER down to 2.36% and the WER to 6.82%.

Sanskrit presented the greatest challenge, with Tesseract yielding high error rates of 8.77% CER and 44.73% WER. Even here, *mBART (all fonts)* outperformed the other models with a CER of 5.67% and a WER of 25.50%, marking a significant improvement. We also tested our best-performing model on 1,000 randomly selected unseen sentences from Maheshwari et al. (2022), which were obtained by OCRing Sanskrit books using Tesseract, resulting in a 6.34% CER and 41.8% WER. After OCR error correction, we achieved a 3.42% CER and 25.7% WER. This improvement in error rates confirms the efficacy of our proposed RoundTripOCR technique in real-world use cases as well. In summary, *mBART (all fonts)* consistently delivered the best results across all languages, reducing both CER and WER considerably compared to raw OCR output from Tesseract, followed

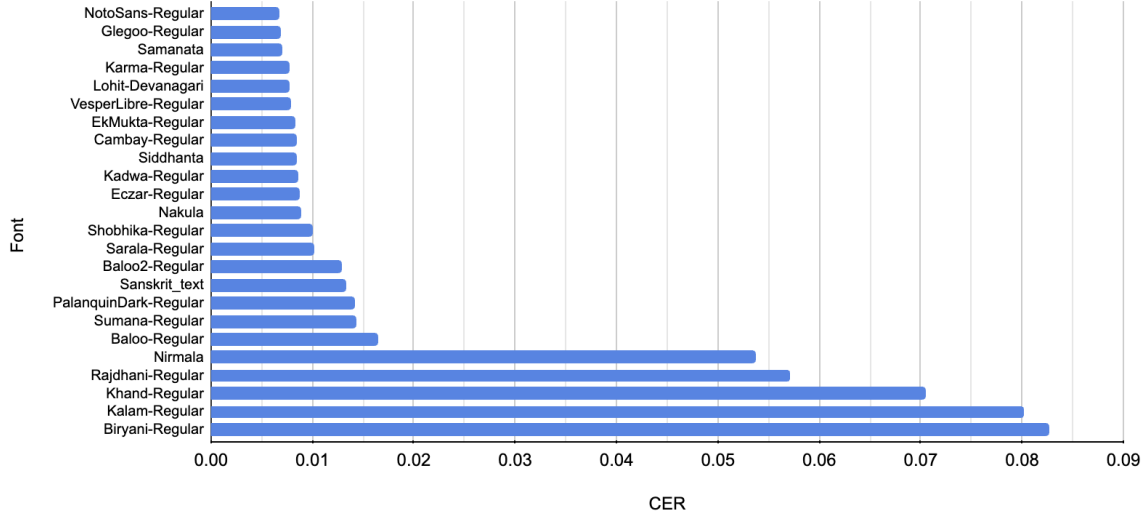


Figure 4: Comparison of different fonts and their CER in the Hindi test dataset.

closely by *mT5 (all fonts)*. These findings highlight the advantage of transformer-based models for OCR error correction.

6 Conclusion and Future Work

We introduced a novel approach for OCR error correction data generation and created a vast dataset comprising 3.1 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 4.07 million sentences in Sanskrit. 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 error output correction, specifically state-of-the-art models like *mBART*, trained on diverse datasets to substantially enhance OCR accuracy. Our research contributes to making textual content more accessible and usable, thereby facilitating broader access to information and knowledge in multilingual societies. Our findings also confirm 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 images as future work. By incorporating these images with controlled variations in font styles, noise levels, and image degradations using a synthetic data generator tool⁶ for text recognition, we can investigate the impact on model generalization and robustness towards real-world document image complexities. 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 error correction and detection, including improvements to the architecture.

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.

Ethical Statement

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

⁶<https://pypi.org/project/trdg>

generating artificial data in this study was sourced from publicly accessible repositories, ensuring no privacy or ethical concerns associated with their use. Specifically, the datasets used do not contain any personally identifiable information or sensitive data that could infringe on individual privacy.

The datasets were chosen based on their availability and openness for research purposes, aligning with ethical guidelines and best practices in data usage. By leveraging publicly available data, this study adheres to the principles of transparency and reproducibility in research while maintaining high ethical standards.

Acknowledgements

The authors thank the anonymous reviewers for their constructive feedback and discussion during the rebuttal, which helped improve this submission. We extend our sincere gratitude to the Computation for Indian Language Technology (CFILT) Lab at the Indian Institute of Technology Bombay for providing the computational resources that were indispensable for the successful completion of this research. The first author would like to thank Himanshu Dutta, Sourabh Deoghare, and P S V N Bhavani Shankar for their invaluable support and assistance in conducting the experiments, which contributed to the progress and quality of this work.

References

- Chantal Amrhein and Simon Clematide. 2018. [Supervised ocr error detection and correction using statistical and neural machine translation methods](#). *Journal for Language Technology and Computational Linguistics (JLCL)*, 33.
- Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurélie Névél, 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.
- Shamil Chollampatt, Raymond Hendy Susanto, Liling Tan, and Ewa Szymanska. 2020. [Can automatic post-editing improve NMT?](#) In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2736–2746, Online. Association for Computational Linguistics.
- Simon Clematide, Lenz Furrer, and Martin Volk. 2016. [Crowdsourcing an OCR gold standard for a German and French heritage corpus](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 975–982, Portorož, Slovenia. European Language Resources Association (ELRA).
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- Raj Dabre, Himani Shrotriya, Anoop Kunchukuttan, Ratish Puduppully, Mitesh Khapra, and Pratyush Kumar. 2022. [IndicBART: A pre-trained model for indic natural language generation](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1849–1863, Dublin, Ireland. Association for Computational Linguistics.
- Sourabh Deoghare, Diptesh Kanojia, Fred Blain, Tharindu Ranasinghe, and Pushpak Bhattacharyya. 2023. [Quality estimation-assisted automatic post-editing](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1686–1698, Singapore. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Eva D’hondt, Cyril Grouin, and Brigitte Grau. 2017. [Generating a training corpus for OCR post-correction using encoder-decoder model](#). In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1006–1014, Taipei, Taiwan. Asian Federation of Natural Language Processing.
- Rui Dong and David Smith. 2018. [Multi-input attention for unsupervised OCR correction](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2363–2372, Melbourne, Australia. Association for Computational Linguistics.
- Agam Dwivedi, Rohit Saluja, and Ravi Kiran Sarvadev-abhatla. 2020. [An ocr for classical indic documents containing arbitrarily long words](#). In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 2386–2393.

- David Etter, Stephen Rawls, Cameron Carpenter, and Gregory Sell. 2019. [A synthetic recipe for ocr](#). pages 864–869.
- Shuhao Guan and Derek Greene. 2024. [Advancing post-OCR correction: A comparative study of synthetic data](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 6036–6047, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Oana Ignat, Jean Maillard, Vishrav Chaudhary, and Francisco Guzmán. 2022. [OCR improves machine translation for low-resource languages](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1164–1174, Dublin, Ireland. Association for Computational Linguistics.
- Atli Jasonarson, Steinnþór Steingrímsson, Einar Sigurðsson, Árni Magnússon, and Finnur Ingimundarson. 2023. [Generating errors: OCR post-processing for Icelandic](#). In *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 286–291, Tórshavn, Faroe Islands. University of Tartu Library.
- 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.
- R. Jayadevan, Satish R. Kolhe, Pradeep M. Patil, and Umapada Pal. 2011. [Offline recognition of devanagari script: A survey](#). *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 41(6):782–796.
- 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.
- Okan Kolak, William Byrne, and Philip Resnik. 2003. [A generative probabilistic OCR model for NLP applications](#). In *Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*, pages 134–141.
- Sobia Nasir Laique, Umar Hayat, Shashank Sarvepalli, Byron Vaughn, Mounir Ibrahim, John McMichael, Kanza Noor Qaiser, Carol Burke, Amit Bhatt, Colin Rhodes, and Maged K. Rizk. 2021. [Application of optical character recognition with natural language processing for large-scale quality metric data extraction in colonoscopy reports](#). *Gastrointestinal Endoscopy*, 93(3):750–757.
- Samuel Lübbli, Mark Fishel, Gary Massey, Maureen Ehrensberger-Dow, and Martin Volk. 2013. [Assessing post-editing efficiency in a realistic translation environment](#). In *Proceedings of the 2nd Workshop on Post-editing Technology and Practice*, Nice, France.
- 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 Luke Zettlemoyer. 2020. [Multilingual denoising pre-training for neural machine translation](#). *Transactions of the Association for Computational Linguistics*, 8:726–742.
- Lijun Lyu, Maria Koutraki, Martin Krickl, and Besnik Fetahu. 2021. [Neural OCR post-hoc correction of historical corpora](#). *Transactions of the Association for Computational Linguistics*, 9:479–493.
- Gábor Madarász, Noémi Ligeti-Nagy, András Holl, and Tamás Váradi. 2024. [Ocr cleaning of scientific texts with llms](#). In *Natural Scientific Language Processing and Research Knowledge Graphs: First International Workshop, NSLP 2024, Hersonissos, Crete, Greece, May 27, 2024, Proceedings*, page 49–58, Berlin, Heidelberg. Springer-Verlag.
- Ayush Maheshwari, Nikhil Singh, Amrith Krishna, and Ganesh Ramakrishnan. 2022. [A benchmark and dataset for post-OCR text correction in Sanskrit](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 6258–6265, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Jie Mei, Aminul Islam, Abidrahman Moh’d, Yajing Wu, and Evangelos Milios. 2018. [Statistical learning for ocr error correction](#). *Information Processing and Management*, 54:Pages 874–887.
- Thi Nguyen, Adam Jatowt, Mickaël Coustaty, and Antoine Doucet. 2021a. [Survey of post-ocr processing approaches](#). *ACM Computing Surveys*, 54:1–37.
- Thi Nguyen, Adam Jatowt, Nhu-Van Nguyen, Mickaël Coustaty, and Antoine Doucet. 2020. [Neural machine translation with bert for post-ocr error detection and correction](#). pages 333–336.
- Thi Tuyet Hai Nguyen, Adam Jatowt, Mickael Coustaty, and Antoine Doucet. 2021b. [Survey of post-ocr processing approaches](#). *ACM Comput. Surv.*, 54(6).
- Santanu Pal, Sudip Kumar Naskar, and Josef van Genabith. 2016. [Multi-engine and multi-alignment based automatic post-editing and its impact on translation productivity](#). In *Proceedings of COLING 2016, the*

26th International Conference on Computational Linguistics: Technical Papers, pages 2559–2570, Osaka, Japan. The COLING 2016 Organizing Committee.

Human Language Technologies, pages 483–498, Online. Association for Computational Linguistics.

Kristen Parton, Nizar Habash, Kathleen McKeown, Gonzalo Iglesias, and Adrià de Gispert. 2012. [Can automatic post-editing make MT more meaningful](#). In *Proceedings of the 16th Annual conference of the European Association for Machine Translation*, pages 111–118, Trento, Italy. European Association for Machine Translation.

Shruti Rijhwani, Antonios Anastasopoulos, and Graham Neubig. 2020. [OCR Post Correction for Endangered Language Texts](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5931–5942, Online. Association for Computational Linguistics.

Shruti Rijhwani, Daisy Rosenblum, Antonios Anastasopoulos, and Graham Neubig. 2021. [Lexically aware semi-supervised learning for OCR post-correction](#). *Transactions of the Association for Computational Linguistics*, 9:1285–1302.

Elizabeth Soper, Stanley Fujimoto, and Yen-Yun Yu. 2021. [BART for post-correction of OCR newspaper text](#). In *Proceedings of the Seventh Workshop on Noisy User-generated Text (W-NUT 2021)*, pages 284–290, Online. Association for Computational Linguistics.

Alan Thomas, Robert Gaizauskas, and Haiping Lu. 2024. [Leveraging LLMs for post-OCR correction of historical newspapers](#). In *Proceedings of the Third Workshop on Language Technologies for Historical and Ancient Languages (LT4HALA) @ LREC-COLING-2024*, pages 116–121, Torino, Italia. ELRA and ICCL.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. [Attention is all you need](#).

Shafqat Mumtaz Virk, Dana Dannélls, and Azam Sheikh Muhammad. 2021. [A novel machine learning based approach for post-OCR error detection](#). In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*, pages 1463–1470, Held Online. INCOMA Ltd.

Martin Volk, Lenz Furrer, and Rico Sennrich. 2011. Strategies for reducing and correcting ocr errors. In *Language Technology for Cultural Heritage: Selected Papers from the LaTeCH Workshop Series*, pages 3–22. Springer.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. [mT5: A massively multilingual pre-trained text-to-text transformer](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics:*

7 Appendix

7.1 PIL (Python Image Library)

The Python Imaging Library, commonly known as PIL, is particularly well-suited for image archival and batch-processing applications. Pillow⁷, an extension of PIL (Python Image Library), stands out as a crucial module for image processing in Python. We generated the images using PIL with dimensions 300x300 and text with a font size of 16.

7.2 Pytesseract

Pytesseract⁸ acts as a wrapper around Tesseract OCR engine. Tesseract⁹ is an open-source OCR engine designed to extract printed or handwritten text from images. Tesseract boasts support for language recognition in over 100 languages straight out of the box. Since it's open-source, it allows flexibility for customization, integration, and experimentation, which is beneficial in research contexts like error correction. Tesseract is lightweight and can be run on various platforms without requiring extensive computational resources. In contrast, commercial models like Google Vision or OCR engines like Ocular may involve higher resource consumption or come with usage restrictions or costs.

⁷<https://pypi.org/project/pillow>

⁸<https://pypi.org/project/pytesseract>

⁹<https://github.com/tesseract-ocr/tesseract>