

AnciDev: A Dataset for High-Accuracy Handwritten Text Recognition of Ancient Devanagari Manuscripts

Authors: Vriti Sharma, Rajat Verma and Rohit Saluja

Introduction

- Indian manuscripts contain **millennia of accumulated knowledge** in mathematics, astronomy, medicine, linguistics, philosophy, and logic.
- India's manuscripts represent **hundreds of languages and scripts** making them essential for understanding linguistic evolution and preserving endangered languages.



मिकांतनिबंधको॑ समयसारञ्चनिनेदि॥३॥ गवेदिसमाधिसुन्त्रको॑ नमिसमन्नाव॥
स्वरूप्यव्यावउञ्चनयेप्रेवंदिको॑ चउशारण्लेष्ट्रष्ट्रवउत्तममंगलष्ट्रणमिकह्वै॥
कियाञ्चिवस्मृ॥४॥ देवधर्मगुरुष्ट्रणतिकरि॥ यादवादञ्चव्लोकि॥ कियाकोञ्चनाथा
कह्वै॥ कुंदकुंदपददोकि॥५॥ अरचौंअरचाजैंनकी॑ चस्तौचरचाजैंनको॑ ओधलेन्तर
लमोहमद॥ त्यगिगाङ्गेगुननैनाधण्कर्त्तमओरञ्चकृत्तमा॥ जिनष्ट्रिमाजिनगेह॥ तिन
सबको॑ परलांगमकरि॥ धारुंधर्मसेह॥६॥ धर्मांकुचउविधिदानसुनगांकंदसधाधर्म॥
नांकेषोडसजावनो॥ नमिरतनञ्चयपर्मी॥७॥ मतकंसर्वयतीश्वराविनञ्च-आर्यासर्व॥

Figure1: Ancient Indian Pandulipi Manuscripts

Unavailability of open-source dataset on Indian manuscripts

Goal: To propose a 500-manuscript-page dataset with approximately 3000 lines from diverse manuscripts in Devanagari script

Advantages of this dataset:

- **Sufficient training data:** 3000 lines provides a meaningful corpus for training deep learning models, particularly for tasks like manuscript analysis or character segmentation.
- **Diversity enables generalization:** Drawing from diverse manuscripts creates a robust dataset that helps models generalize specific patterns or document types.
- **Benchmark potential:** A well-curated dataset of this size could serve as a standard benchmark for the research community.

Dataset Samples

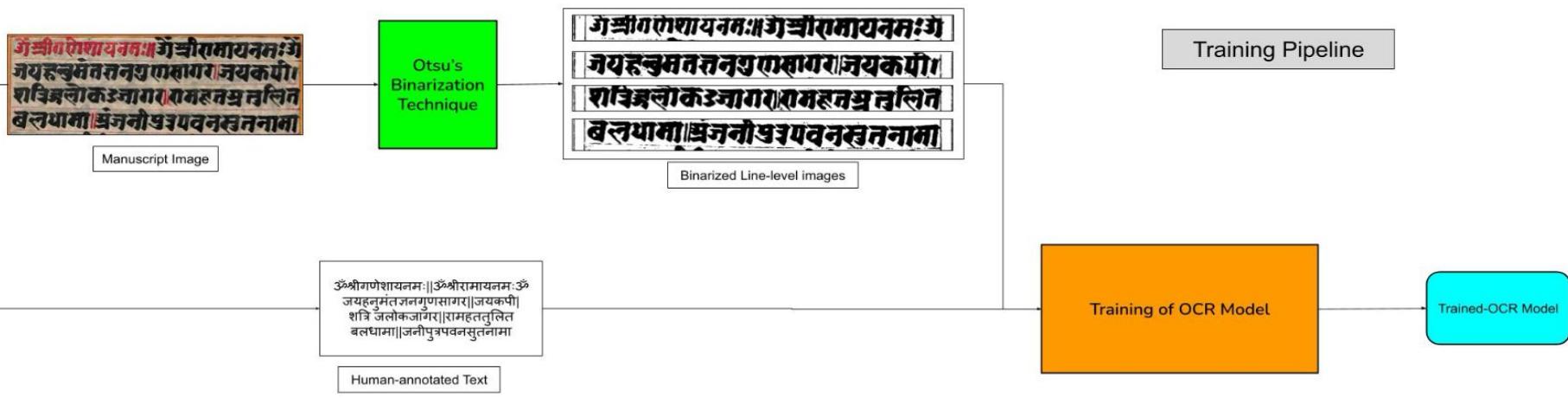
लउसारेकहांगयेपाडोसीहमारे॥२७॥
 ससास्वामनिर्जीर्जीतरै॥कुजीमूर्गीआ
 ईनहीमेरै॥झंबेदपुरांणपड्योत्तहाबांणी
 ॥मुकुंकहावुलावैरांणी॥३०॥षषा-
 षावधणीहैषगीसियांणी॥वांहपकडि

लनु सारे कहां गये पाडोसी हमारे ॥२९॥
 ससाखामनिर्वभीर्वभीतरे दुजामूर आ
 ईन ही मेरे॥ झंबेदपुराण पद्यो न ही बाणी
 मुकुंकहावुलावैरांणी॥३०॥षषा
 षावधणी हैषरी सियांणी॥वांहपकडि

मिक्षांतनिबंधकौं समयसारअभिनंदि॥२८॥ वंदिसमाधिसुतंत्रकौं नमिसमनाव।
 सरूपवप्यवउ-अनयोगैवंदिकैं चउशरण्लेषुष्वचउत्तममंगलश्लभिकहौं।
 क्रियाअविरुद्ध॥२९॥ देवधर्मगुरप्रणतिकरि॥सादवाद-अवलोकि॥ क्रियाकोशभाषा
 करुं। कुंदकुंदयदोकि॥३०॥ अरचौंअरचाजैनकी॥चरचौंचरचाजैन॥ कोधलेनच
 लमोहमद॥त्यागिगङ्गंगुननैन॥३१॥ कर्त्तमओरअकृत्तमा॥जिनशतिमाजिनगेहतिन॥

सिद्धांतनिबंधकौं। समयसारअभिनंदि ॥३१॥ वंदिसमाधिसुतंत्रकौंनमिसमभाव।
 स्वरूप ॥३२॥ चउअनयोगैवंदिकैं। चनशरण्लेशुछ। चउत्तममंगलश्लभिकहौं।
 क्रियाअविरुद्ध ॥३३॥ देवधर्मगुरप्रणतिकरि। स्यादवादअवलोकि ॥ क्रियाकोशभाषा
 करुं। कुंदकुंदयदोकि ॥३०॥ अरचौंअरचाजैनकी। चरचौंचरचाजैन। क्रोधलोभछ
 लमोहमद ॥ त्यागिगफूंगुननैन ॥३१॥ कर्त्तमओरकृत्तमा ॥ जिनप्रतिमाजिनगेह ॥ तिन
 ।

Training Pipeline



Attention-OCR and CNN-RNN Finetuning

Finetuning of Attention-OCR and CNN-RNN:

No. of Epochs	200
Batch Size	32
Total pre-training dataset size (synthetic)	820 (fonts) x 5000 per Font = 4100000 line-level dataset
Total Dataset Size	3085 lines
Training Data Size (80:20 split)	2458 lines
Test Dataset Size (80:20 split)	627 lines

Tesseract-5 Finetuning

Fine-tuning of Tesseract-5:

No. of Epochs	200
Batch Size	32
Total pre-training dataset size (synthetic)	7000 synthetic lines created using real verses text + 3000 real lines
Total Dataset Size	3085 lines
Training Data Size (80:20 split)	2458 lines
Test Dataset Size (80:20 split)	627 lines

Quantitative Analysis

1. **Character Error Rate:** (Number of Incorrect characters in predicted text / total number of characters in reference text) * 100
2. **Word Error Rate:** (Number of Incorrect words in predicted text / total number of word in reference text) * 100

Model Name	Test Set	
	WER	CER
CNN-RNN	98.20	48.59
Attention-LSTM	96.73	46.33
Tesseract-5	87.42	30.06

m, s denotes manuscript and synthetic data respectively

Qualitative Analysis

Input image	Ground Truth	Tesseract-5's Predictions
राधन तणा॥ चेष्टैदवा न्वधिकारा विजञ्चाणीने आदरो लि॥	राधन तणा॥ अछैदश अधिकार॥ चित्तआणीने आदरो जि	राधन तणे॥ शैदश प्रधिकार॥ चित्तआणीने प्रादरो लि
धत्तेस्वतंत्रःपरिपूर्णआत्मा ५० नमोस्तुतेरामतवाप्रिपंकजं	धत्तेस्वतंत्रःपरिपूर्णआत्मा ५० नमोस्तुतेरामतवाप्रिपंकजं	धक्तेस्वतंत्रःपरिपूर्ण आत्मा ५० नमोस्तुतेरामतवाप्रिपंकजं
थायच्यपुनर्द्वारामंराजीवलोचनम् पुलकां कितसर्वांगागिरा	त्थाय चपुनर्द हवा रामं राजीवलोचनम् पुल को कि त सर्वांगा गिरा	त्थायचपुनर्द्वारामंराजीवलोचनम् पुलकां कितसर्वांगागिरा
५८ नमस्तेपुरुषाध्यक्षनभस्तेभक्तवत्सल नमस्तुहृषी	५८ नमस्तेपुरुषाध्यक्षनभस्तेभक्तवत्सल नमस्तुहृषी,	५६ नमस्तेपुरुषाध्यक्षनभस्तुभक्तवत्सल नमस्तुहृषी,
बोपबाढीसखबलकभूत्यै बद्धते कांतिकाढी। परमबदनरभारे	चोप बाढी सुख कवलक भूत्यै चंद्र ते कांति काढी परम बदन रभारे	बोप बाढी सुख कवलक भूत्यै चंद्र ते कांति काढी परम बदन एभारे
शामि नमभुम वदने। सिंहिणजिचगम नही	शशि निभशुभ वदने। सिंहिणजिचगम नही	शमि नमभुम वदने। सिंहिणजिचगम नही

Green and red color for correct and incorrect prediction

Conclusion and Future Work

Contributions:

- First public dataset for ancient Devanagari HTR (3K lines, 500 pages)
- Reproducible baselines with 3 HTR models
- New state-of-the-art: 30.06% CER (Tesseract-5)

Future directions:

- **Dataset expansion:** More manuscript types and time periods with diverse scribal hands for better generalization
- **Model improvements:** Advanced transformer architectures, specialized data augmentation for historical writing, language model-based post-processing
- **Evaluation:** Multi-annotator studies for reliability analysis

Thank you

For more info, scan the given QR code:

