

Bangla Handwritten Text Tokenization and Recognition: A Hybrid Approach Using Regular Expressions and Cloud-Based Vision APIs

Mahbubur Rahman Khan

Department of Electrical and Computer Engineering

North South University

Dhaka, Bangladesh

mahbubur.khan@northsouth.edu

Abstract—This project presents a comprehensive pipeline for Bangla handwritten text tokenization and recognition. We develop a dual-level tokenization system that handles both word-level and character-level segmentation using regular expression-based finite state automata, addressing the unique challenges of Bangla script including compound consonants, vowel diacritics, and contextual character formations. For handwritten text recognition, we propose a hybrid approach that combines Google Vision API with Gemini LLM for context-aware error correction, eliminating the need for extensive model training. Our tokenizer achieves a type-token ratio of 0.746 for word-level and 0.212 for character-level tokenization on test datasets. The recognition system demonstrates variable performance across different handwriting styles, with character accuracy ranging from 3.4% to 71.2% depending on dataset complexity. The complete implementation and datasets are available at: <https://github.com/mahbubkousar/bangla-handwritten-text-tokenizer>.

Index Terms—Bangla tokenization, handwritten text recognition, regular expressions, finite state automata, OCR, natural language processing

I. INTRODUCTION

Bangla, one of the world's most spoken languages with over 230 million native speakers, presents unique challenges for digital text processing due to its complex script characteristics. The language employs a rich orthographic system with 11 vowels, 39 consonants, numerous compound characters (যুক্তাক্ষর), and contextual vowel modifiers (কার) that significantly complicate both tokenization and handwritten text recognition tasks.

This project addresses two fundamental challenges in Bangla text processing: (1) developing a robust tokenization system that correctly segments Bangla text at both word and character levels while preserving orthographic integrity, and (2) creating an efficient handwritten text recognition pipeline that leverages cloud-based APIs to avoid the computational overhead of training deep learning models from scratch.

Our contribution includes a formal definition of Bangla tokens using regular expressions, implementation of deterministic finite automata for token recognition, and a

hybrid approach combining Google Vision API with large language model correction for handwritten text recognition. We evaluate our system on custom datasets representing various handwriting styles and demonstrate its effectiveness through comprehensive metrics.

II. RELATED WORK

Recent advances in Bangla handwritten text recognition have been driven primarily by deep learning architectures. Alif et al. [1] proposed a modified ResNet-34 architecture achieving 98.70% accuracy for vowels, 97.34% for consonants, and 99.02% for numeric characters. Their approach optimized dropout layer placement, applying dropout after every two layers within ResNet blocks.

Opu et al. [2] developed a lightweight 74-layer CNN architecture specifically for Handwritten Bangla Character Recognition (HBCR), achieving 96.87% accuracy with a significantly smaller model size of 13.2 MB compared to existing models ranging from 51.6 MB to 403.2 MB. This work emphasized efficiency alongside accuracy, demonstrating faster loading (1.911 seconds) and testing times (7.956 seconds).

For tokenization, Chakraborty et al. [3] explored Unicode block-based pattern matching for Bangla text segmentation, establishing foundational principles for handling conjunct consonants and vowel diacritics. Moll et al. [4] provided the theoretical framework for formal language theory applications in text processing.

While these deep learning approaches achieve impressive accuracy, they require substantial computational resources and training data. Our work diverges by proposing a training-free hybrid approach that combines existing cloud services with context-aware correction, making it more accessible for resource-constrained environments typical in undergraduate research settings.

III. METHODOLOGY

Our methodology includes two main components: a formal tokenization system based on regular expressions and finite state automata, and a hybrid handwritten text recognition pipeline combining cloud vision APIs with

LLM-based correction. Figure 1 illustrates the complete system architecture.

A. Formal Token Definition

We formally define the Bangla tokenization problem using principles from formal language theory. Our approach establishes a clear mathematical foundation for token recognition.

1) *Alphabet Definition*: The alphabet Σ encompasses all valid Bangla characters:

$$\Sigma = \{\text{অ, আ, ..., ই, ও, এ, ..., ক, ঙ, ি, ?, ,, '}\}$$

This includes:

- Base characters (11 vowels, 39 consonants)
- Modifiers and diacritics
- Bangla numerals (০-৯)
- Punctuation marks (both native and borrowed)
- Whitespace characters

2) *Language of Tokens*: The set of all valid tokens forms a regular language \mathcal{L} over alphabet Σ . Each token is a string following specific formation rules defined through regular expressions. Our tokenizer partitions input strings into sequences of valid tokens from \mathcal{L} .

B. Token Classification

We implement two tokenization levels:

Word-level tokens: Space-separated lexical units including words, numbers, punctuation, and whitespace.

Character-level tokens: Individual graphemes including single consonants, vowels, conjunct characters, and modifiers. A grapheme represents the smallest meaningful visual unit in Bangla script.

IV. UNICODE BLOCK-BASED PATTERN RECOGNITION

Our tokenization system employs Unicode block-based pattern matching to identify distinct lexical units. Each pattern corresponds to a specific Unicode range ensuring accurate token boundary detection.

A. Word-Level Patterns

We define four primary patterns for word-level tokenization:

Word Token Pattern:

1 WORD: [\u0980-\u09FF]+

This pattern captures sequences containing consonants (ক্ষনবর্ণ), vowels (স্বরবর্ণ), vowel diacritics (কার), and conjunct consonants (যুক্তাক্ষর).

Numeric Token Pattern:

1 NUMBER: [\u09E6-\u09EF]+

Recognizes Bangla numerals (০-৯) in the dedicated Unicode range.

Punctuation Token Pattern:

1 PUNCTUATION: [,?;!,:] | \u0964 | \u0965

Handles both native punctuation (দাঁড়ি , দিদাঁড়ি !!) and borrowed Latin punctuation.

Whitespace Token Pattern:

1 WHITESPACE: \s+

Matches standard Unicode whitespace categories.

B. Character-Level Patterns

Character tokenization requires careful handling of compound formations and diacritics:

Compound Consonant Pattern:

1 COMPOUND: [\u0995-\u09B9] (?:\u09CD[\u0995-\u09B9])+

Consonant-Vowel Pattern:

1 CONSONANT_VOWEL: [\u0995-\u09B9] [\u09BE-\u09CC]

Independent Vowel Pattern:

1 VOWEL: [\u0985-\u0994]

Pattern matching follows a strict precedence hierarchy:

- 1) Compound consonants (highest priority)
- 2) Vowel-consonant combinations
- 3) Individual consonants
- 4) Independent vowels
- 5) Modifiers and diacritics (lowest priority)

This ordering prevents incorrect decomposition of valid compound characters.

V. TOKENIZATION ALGORITHMS

We implement two algorithms for word and character-level tokenization using finite state automata principles.

A. Word-Level Tokenization

Algorithm 1 processes input text character by character, maintaining state transitions based on Unicode character categories. The algorithm implements the maximal munch principle, ensuring the longest valid token is recognized at each position.

B. Character-Level Tokenization

Character tokenization (Algorithm 2) recognizes grapheme clusters rather than individual Unicode code points. This is crucial for Bangla where single user-perceived characters may comprise multiple code points.

Key features include:

- Unicode Normalization Form C (NFC) preprocessing for consistent representation
- Precedence-based pattern matching for correct compound recognition
- Position tracking for maintaining token boundaries
- Unknown token handling for unrecognized characters

VI. HANDWRITTEN TEXT RECOGNITION PIPELINE

Our recognition pipeline combines traditional image processing with modern cloud-based OCR and large language model correction, eliminating the need for extensive model training.

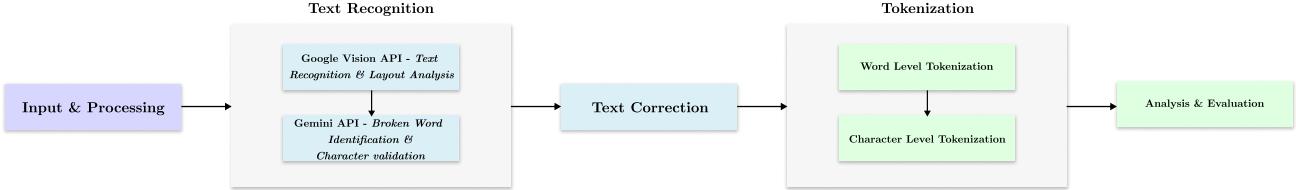


Figure 1. Complete system architecture showing the dual pipeline for tokenization and handwritten text recognition. The tokenization branch processes digital text through word and character-level segmentation, while the recognition branch handles handwritten images through preprocessing, OCR extraction, and LLM-based correction.

Algorithm 1 Bangla Word-Level Tokenization

```

Require: Text string  $T$ , word patterns  $P_{word}$ 
Ensure: List of word tokens  $W$ 
1:  $T_{norm} \leftarrow \text{NORMALIZE}(T)$  {Unicode NFC}
2:  $W \leftarrow \emptyset$ 
3:  $pos \leftarrow 0$ 
4: while  $pos < |T_{norm}|$  do
5:    $matched \leftarrow \text{false}$ 
6:   for each pattern  $(label, regex) \in P_{word}$  do
7:      $match \leftarrow \text{REGEX\_MATCH}(regex, T_{norm}, pos)$ 
8:     if  $match \neq \text{null}$  then
9:        $token \leftarrow \text{EXTRACT\_TEXT}(match)$ 
10:      if  $label \neq \text{WHITE SPACE}$  OR  $token \neq \text{empty}$  then
11:         $W \leftarrow W \cup \{(token, label, pos, |token|)\}$ 
12:      end if
13:       $pos \leftarrow \text{END\_POSITION}(match)$ 
14:       $matched \leftarrow \text{true}$ 
15:      break
16:    end if
17:   end for
18:   if  $matched = \text{false}$  then
19:      $W \leftarrow W \cup \{(T_{norm}[pos], UNKNOWN, pos, 1)\}$ 
20:      $pos \leftarrow pos + 1$ 
21:   end if
22: end while
23: return  $W$ 

```

A. Image Preprocessing

The preprocessing pipeline optimizes input images for OCR:

- 1) **Grayscale Conversion:** Reduces computational complexity from $O(3n)$ to $O(n)$
- 2) **Denoising:** Fast non-local means denoising with $h = 10$ and 21×21 search window removes Gaussian noise while preserving edges
- 3) **Adaptive Thresholding:** Local threshold computation using $T(x, y) = \mu(x, y) - C$ within an 11×11 Gaussian kernel ($\sigma = 1.5$, $C = 5$)

This approach handles illumination variations exceeding

Algorithm 2 Bangla Character-Level Tokenization

```

Require: Text string  $T$ , character patterns  $P_{char}$ 
Ensure: List of character tokens  $C$ 
1:  $T_{norm} \leftarrow \text{NORMALIZE}(T)$  {Unicode NFC}
2:  $C \leftarrow \emptyset$ 
3:  $pos \leftarrow 0$ 
4: while  $pos < |T_{norm}|$  do
5:    $matched \leftarrow \text{false}$ 
6:   for each pattern  $(label, regex) \in P_{char}$  by precedence do
7:      $match \leftarrow \text{REGEX\_MATCH}(regex, T_{norm}, pos)$ 
8:     if  $match \neq \text{null}$  then
9:        $token \leftarrow \text{EXTRACT\_TEXT}(match)$ 
10:       $C \leftarrow C \cup \{(token, label, pos, |token|)\}$ 
11:       $pos \leftarrow \text{END\_POSITION}(match)$ 
12:       $matched \leftarrow \text{true}$ 
13:      break
14:    end if
15:   end for
16:   if  $matched = \text{false}$  then
17:      $C \leftarrow C \cup \{(T_{norm}[pos], UNKNOWN, pos, 1)\}$ 
18:      $pos \leftarrow pos + 1$ 
19:   end if
20: end while
21: return  $C$ 

```

20% intensity range, critical for Bangla handwriting where ink density fluctuations cause up to 15% recognition errors with global thresholding.

B. Hybrid Recognition Algorithm

Algorithm 3 describes our hybrid approach combining Google Vision API with Gemini LLM for context-aware correction.

VII. EXPERIMENTAL SETUP AND RESULTS

A. Dataset Preparation

We created custom datasets with four categories for each text sample:

- 1) **Clean Handwriting:** Carefully hand-written text
- 2) **Fast Writing:** Natural daily handwriting style

Algorithm 3 Hybrid Bengali Handwriting Recognition

Require: Input image I , confidence threshold τ

Ensure: Corrected Bengali text $T_{corrected}$

- 1: **Image Preprocessing:**
- 2: $I_{gray} \leftarrow \text{RGB2Gray}(I)$
- 3: $I_{denoised} \leftarrow \text{FastNLMeansDenoising}(I_{gray})$
- 4: $I_{binary} \leftarrow \text{AdaptiveThreshold}(I_{denoised})$
- 5: **Primary Text Extraction:**
- 6: $\mathcal{W} \leftarrow \text{GoogleVisionAPI}(I_{binary}, \text{lang}=[\text{bn}, \text{en}])$
- 7: $T_{original} \leftarrow \text{concatenate}(\mathcal{W})$
- 8: Extract confidence scores: $C_i \leftarrow \text{confidence}(w_i)$
- 9: **Error Detection:**
- 10: $\mathcal{B} \leftarrow \emptyset$ {Broken words set}
- 11: **for** each word $w_i \in \mathcal{W}$ **do**
- 12: **if** $C_i < \tau$ OR $\text{HasUnusualPatterns}(w_i)$ **then**
- 13: $\mathcal{B} \leftarrow \mathcal{B} \cup \{w_i\}$
- 14: **end if**
- 15: **end for**
- 16: **Context-Aware Correction:**
- 17: **if** $|\mathcal{B}| > 0$ **then**
- 18: prompt $\leftarrow \text{ConstructPrompt}(T_{original}, \mathcal{B})$
- 19: $T_{corrected} \leftarrow \text{GeminiLLM}(\text{prompt})$
- 20: **if** $\text{IsReasonable}(T_{corrected}, T_{original})$ **then**
- 21: $T_{final} \leftarrow T_{corrected}$
- 22: **else**
- 23: $T_{final} \leftarrow T_{original}$
- 24: **end if**
- 25: **else**
- 26: $T_{final} \leftarrow T_{original}$
- 27: **end if**
- 28: **Post-processing:**
- 29: $T_{final} \leftarrow \text{RemoveExtraSpaces}(T_{final})$
- 30: $T_{final} \leftarrow \text{FixPunctuation}(T_{final})$
- 31: **return** T_{final}

3) **Ruled Paper Writing:** Text written on lined paper

4) **Ground Truth:** Typed version for evaluation

Five different text samples (D1-D5) were collected, each written in all three styles, providing 15 handwritten images with corresponding ground truth texts.

B. Implementation Environment

All experiments were conducted on Google Colab platform with:

- Google Vision API for OCR extraction
- Gemini API for context-aware correction
- Python with OpenCV for image preprocessing
- Regular expression library for tokenization

C. Tokenization Results

Table I presents tokenization statistics across different handwriting styles compared to ground truth.

The tokenizer successfully identified distinct token types with the following distribution for ground truth text:

Table I
TOKENIZATION STATISTICS SUMMARY

Metric	File Type			
	Clean	Fast	Ruled	Ground Truth
Word-Level				
Total Tokens	107.0	108.2	105.4	122.6
Type-Token Ratio	0.772	0.738	0.753	0.744
Character-Level				
Total Tokens	434.2	445.6	428.6	500.6
Type-Token Ratio	0.267	0.261	0.265	0.254

- Word tokens: 85.9% of total
 - Punctuation: 10.8%
 - Other tokens: 3.2%
- Character-level tokenization revealed complex grapheme distributions:
- Consonants: 25.9%
 - Consonant-vowel combinations: 24.8%
 - Whitespace: 20.2%
 - Compound characters: 8.6%
 - Dependent vowels: 6.8%
 - Independent vowels: 5.8%

D. Recognition Performance

Table II summarizes recognition accuracy across datasets and handwriting styles.

Table II
RECOGNITION ACCURACY BY DATASET

Dataset	Character Acc. (%)	Word Acc. (%)	Total Words
D1	48.5	49.2	465
D2	4.2	0.9	318
D3	3.4	2.7	333
D4	58.1	55.6	324
D5	71.2	72.6	135
By Style			
Clean	38.0	38.1	525
Fast	32.3	32.1	525
Ruled	40.9	38.5	525

Recognition accuracy varied significantly across datasets, with D5 achieving the highest performance (71.2% character accuracy) and D2-D3 showing poor results (below 5%). This variation reflects differences in handwriting complexity and text content. Ruled paper writing showed marginally better performance (40.9%) compared to clean (38.0%) and fast writing (32.3%) styles.

VIII. DISCUSSION

Our dual-pipeline approach demonstrates both strengths and limitations in handling Bangla text processing tasks.

A. Tokenization Analysis

The tokenization system successfully handles complex Bangla orthographic features including compound consonants and vowel diacritics. The type-token ratio of 0.744–0.772 for word-level tokenization indicates good vocabulary coverage. Character-level tokenization effectively identifies grapheme clusters, crucial for maintaining semantic integrity in Bangla text.

The precedence-based pattern matching prevents common errors in compound character decomposition. For example, ক্ষ is correctly recognized as a single unit rather than being split into ক and ষ.

B. Recognition Challenges

The high variability in recognition accuracy (3.4% to 71.2%) reveals several challenges:

- 1) **Handwriting Quality:** Complex cursive styles in D2 and D3 significantly degraded performance
- 2) **API Limitations:** Google Vision API, while robust for printed text, struggles with handwritten Bangla script
- 3) **Context Dependency:** The LLM correction improves results but requires sufficient context for accurate prediction

The marginal improvement with ruled paper (40.9% vs 38.0% for clean) suggests that writing guidelines provide minimal benefit for OCR accuracy.

C. Comparison with State-of-the-Art

While deep learning approaches achieve 96-98% accuracy on standardized datasets [1], [2], our training-free approach trades accuracy for accessibility and ease of deployment. This makes it suitable for rapid prototyping and resource-constrained environments.

IX. CONCLUSION

This project presents a comprehensive pipeline for Bangla text tokenization and handwritten recognition. Our contributions include:

- 1) A formal definition of Bangla tokens using regular expressions and finite state automata
- 2) Implementation of dual-level tokenization handling both words and grapheme clusters
- 3) A novel hybrid recognition approach combining cloud APIs with LLM correction
- 4) Custom datasets representing various handwriting styles with ground truth annotations

The tokenization system demonstrates robust performance with consistent type-token ratios across different text styles. While the recognition accuracy varies significantly, the approach provides a practical alternative to resource-intensive deep learning methods.

REFERENCES

- [1] M. A. R. Alif, "State-of-the-art bangla handwritten character recognition using a modified resnet-34 architecture," *Int. J. Innov. Sci. Res. Technol.*, vol. 9, pp. 438–448, 2024.
- [2] M. N. I. Opu, M. E. Hossain, and M. A. Kabir, "Handwritten bangla character recognition using convolutional neural networks: a comparative study and new lightweight model," *Neural Computing and Applications*, vol. 36, no. 1, pp. 337–348, 2024.
- [3] S. Chakraborty, P. Das, S. M. Diplo, M. A. Pramanik, and J. Noor, "An analytical review of preprocessing techniques in bengali natural language processing," *IEEE Access*, 2025.
- [4] R. N. Moll, M. A. Arbib, and A. J. Kfoury, *An introduction to formal language theory*. Springer Science & Business Media, 2012.

APPENDIX A IMPLEMENTATION RESOURCES

A. GitHub Repository

The complete project implementation is available at:

<https://github.com/mahbubkousar/bangla-handwritten-text-tokenizer>

The repository contains:

- **Datasets:** Five text samples with three handwriting styles each
- **Code Implementation:**
 - Complete Bangla Tokenizer Pipeline
 - Bangla Text Recognition System
- **Results:** Tokenization outputs and recognition results for all datasets

APPENDIX B DETAILED EXPERIMENTAL RESULTS

A. Tokenization Performance Metrics

Table III
COMPLETE TOKENIZATION STATISTICS (PART 1)

File Type	Text Length				Word Tokens		Word TTR	
	Mean	Std	Min	Max	Mean	Std	Mean	Std
clean_output	612.2	198.969	273	757	107.0	30.725	0.772	0.048
fast_output	629.0	205.515	267	760	108.2	28.297	0.738	0.029
ground_truth	701.4	294.085	279	1106	122.6	44.590	0.744	0.031
ruled_output	609.8	181.908	287	723	105.4	29.314	0.753	0.066

Table IV
COMPLETE TOKENIZATION STATISTICS (PART 2)

File Type	Character Tokens		Character TTR	
	Mean	Std	Mean	Std
clean_output	434.2	139.810	0.267	0.053
fast_output	445.6	137.442	0.261	0.067
ground_truth	500.6	209.920	0.254	0.059
ruled_output	428.6	127.997	0.265	0.056

B. Comparison with Ground Truth

Table V
WORD-LEVEL METRICS COMPARISON WITH GROUND TRUTH

File Type	Precision		Recall		F1 Score	
	Mean	Std	Mean	Std	Mean	Std
clean_output	0.482	0.417	0.425	0.374	0.450	0.391
fast_output	0.436	0.384	0.376	0.320	0.400	0.342
ruled_output	0.487	0.415	0.415	0.360	0.446	0.380

Table VI
CHARACTER-LEVEL METRICS COMPARISON WITH GROUND TRUTH

File Type	Precision		Recall		F1 Score	
	Mean	Std	Mean	Std	Mean	Std
clean_output	0.810	0.194	0.767	0.169	0.787	0.177
fast_output	0.761	0.160	0.719	0.118	0.737	0.126
ruled_output	0.792	0.178	0.741	0.160	0.764	0.163

APPENDIX C
HANDWRITTEN TEXT RECOGNITION RESULTS

A. Detailed Recognition Performance

Table VII
RECOGNITION RESULTS BY DATASET AND HANDWRITING STYLE

Dataset	Image Type	Char Acc. (%)	Word Acc. (%)	Correct Words	Total Words	Correct Chars	Total Chars
D1	clean	49.2	51.0	79	155	542	1101
D1	fast	29.4	46.5	72	155	323	1101
D1	ruled	66.8	50.3	78	155	735	1101
D2	clean	3.5	0.9	1	106	26	756
D2	fast	3.5	0.9	1	106	26	756
D2	ruled	5.6	0.9	1	106	42	756
D3	clean	3.6	2.7	3	111	24	688
D3	fast	3.3	2.7	3	111	22	688
D3	ruled	3.3	2.7	3	111	22	688
D4	clean	69.5	53.7	58	108	463	666
D4	fast	57.0	54.6	59	108	379	666
D4	ruled	47.8	58.3	63	108	318	666
D5	clean	64.1	82.2	37	45	178	279
D5	fast	68.3	55.6	25	45	190	279
D5	ruled	81.3	80.0	36	45	226	279

B. Aggregate Performance Analysis

Table VIII
SUMMARY PERFORMANCE BY DATASET

Dataset	Char Acc. (%)	Word Acc. (%)	Correct Words	Total Words	Correct Chars	Total Chars
D1	48.5	49.2	229	465	1600	3303
D2	4.2	0.9	3	318	94	2268
D3	3.4	2.7	9	333	68	2064
D4	58.1	55.6	180	324	1160	1998
D5	71.2	72.6	98	135	594	837

Table IX
SUMMARY PERFORMANCE BY HANDWRITING TYPE

Image Type	Char Acc. (%)	Word Acc. (%)	Correct Words	Total Words	Correct Chars	Total Chars
clean	38.0	38.1	178	525	1233	3490
fast	32.3	32.1	160	525	940	3490
ruled	40.9	38.5	181	525	1343	3490

APPENDIX D
SAMPLE DATASET IMAGES

This section presents sample images from each of the five datasets used in our experiments. Each dataset includes three handwriting styles: clean, fast, and ruled paper writing.

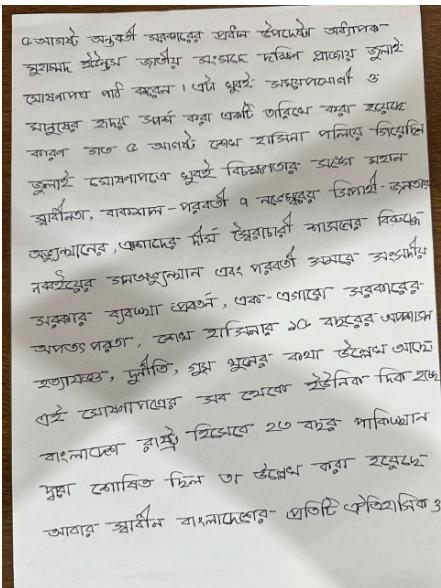


Figure 2. *
(a) Clean Handwriting

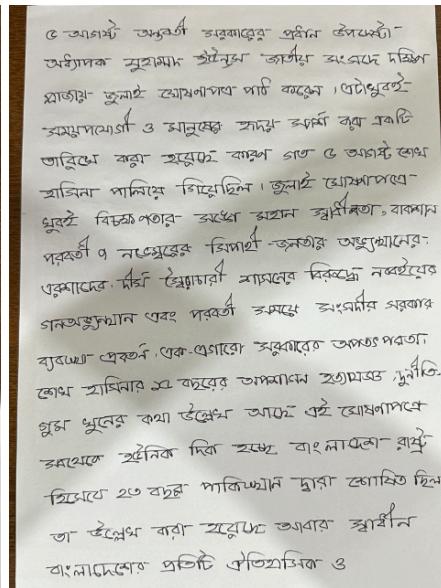


Figure 3. *
(b) Fast Writing

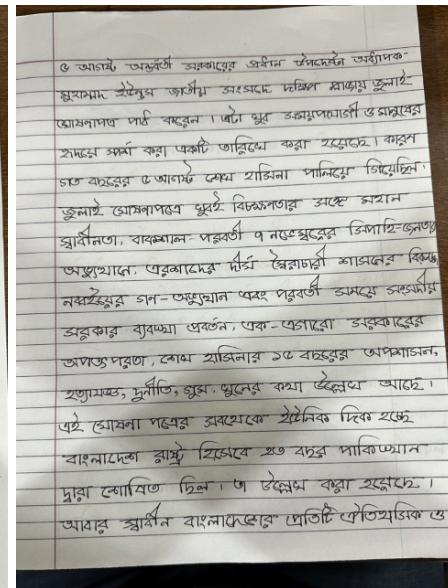


Figure 4. *
(c) Ruled Paper

Figure 5. Dataset 1 (D1): Sample handwriting styles showing variation in writing quality

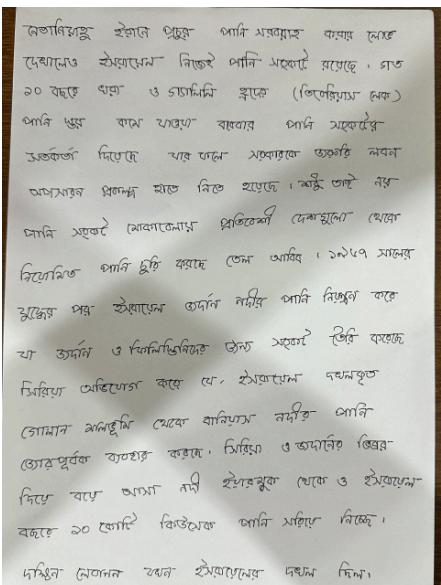


Figure 6. *
(a) Clean Handwriting

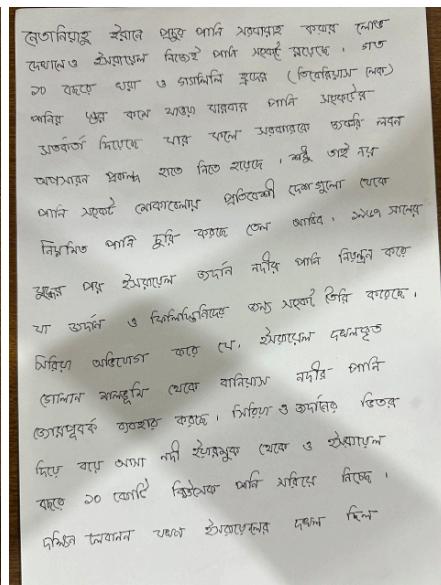


Figure 7. *
(b) Fast Writing

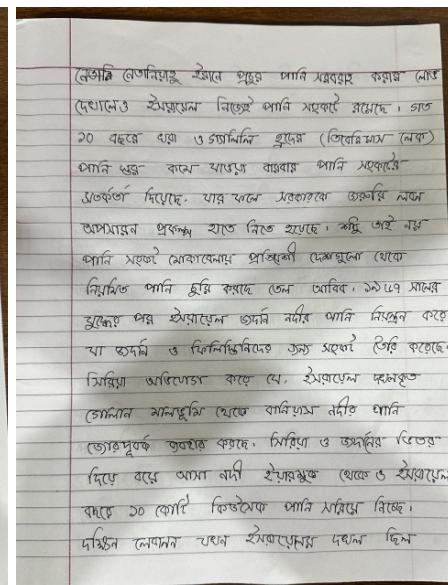


Figure 8. *
(c) Ruled Paper

Figure 9. Dataset 2 (D2): Complex cursive handwriting style with lower recognition accuracy

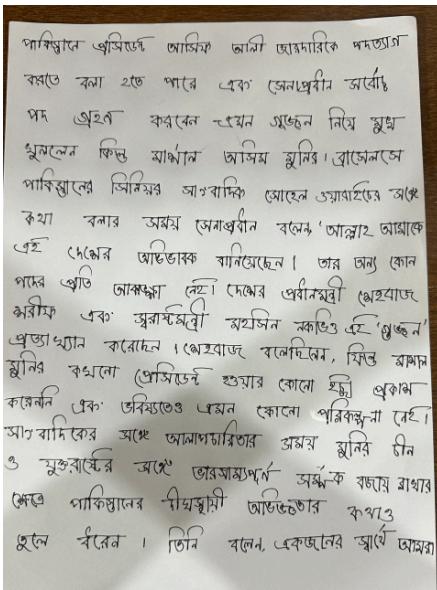


Figure 10. *
(a) Clean Handwriting

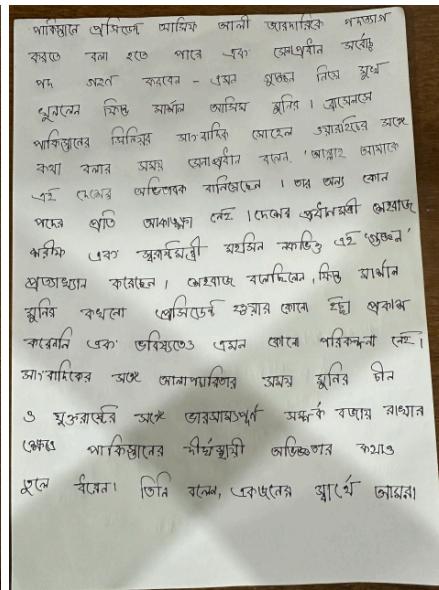


Figure 10. *
(b) Fast Writing

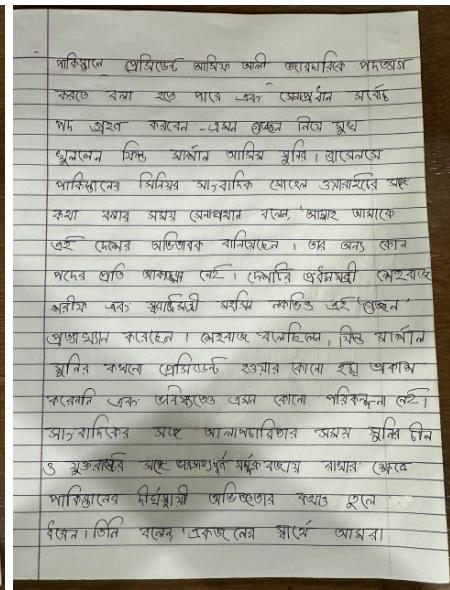


Figure 10. *
(c) Ruled Paper

Figure 13. Dataset 3 (D3): Highly stylized handwriting presenting significant recognition challenges

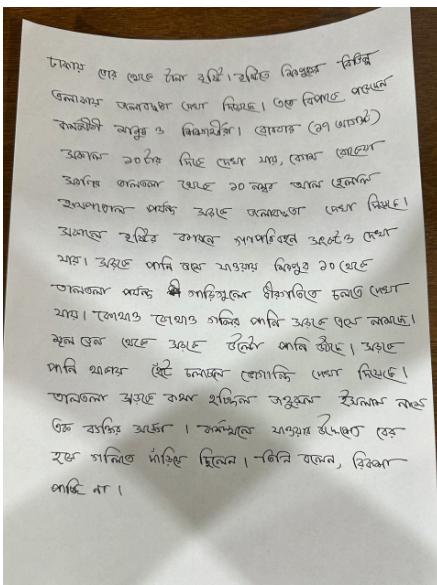


Figure 14. *
(a) Clean Handwriting

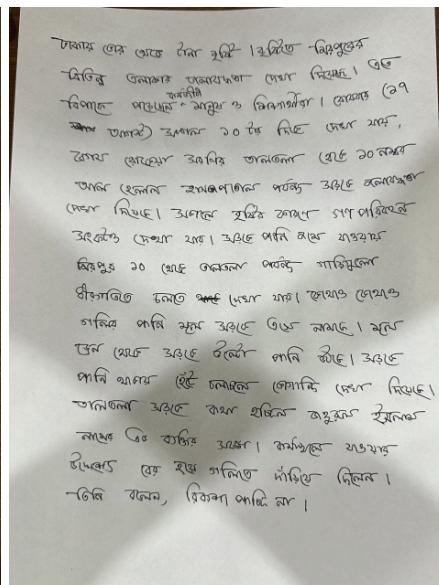


Figure 14. *
(b) Fast Writing

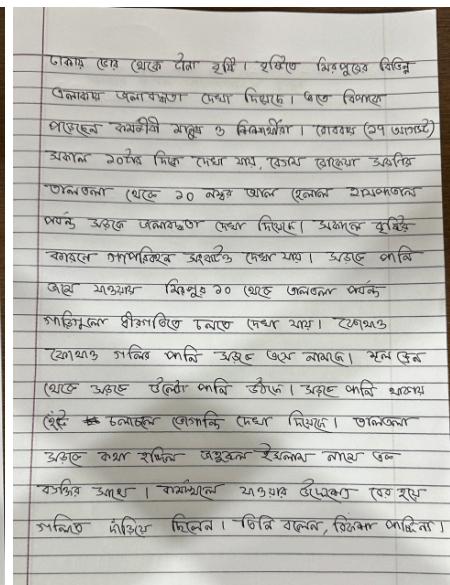


Figure 14. *
(c) Ruled Paper

Figure 17. Dataset 4 (D4): Moderate complexity handwriting with reasonable recognition performance

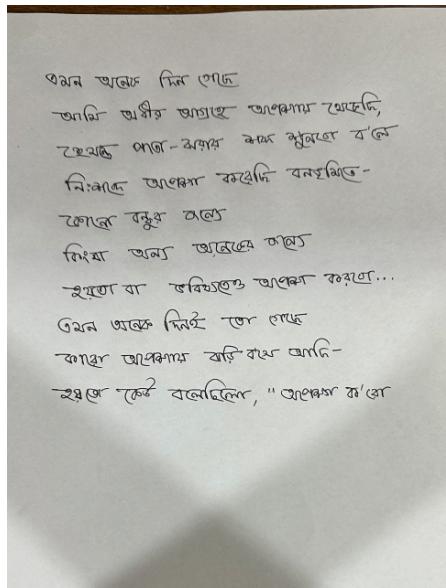


Figure 18. *
(a) Clean Handwriting

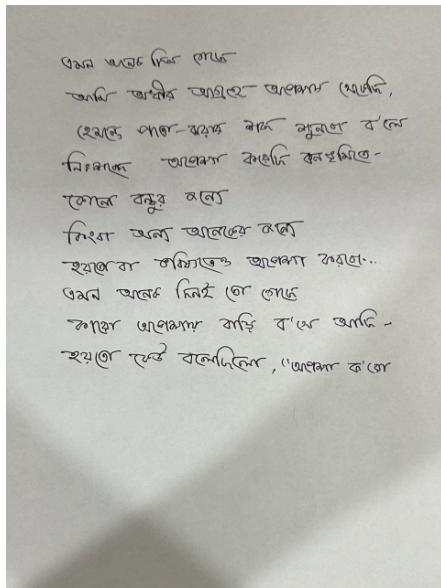


Figure 19. *
(b) Fast Writing

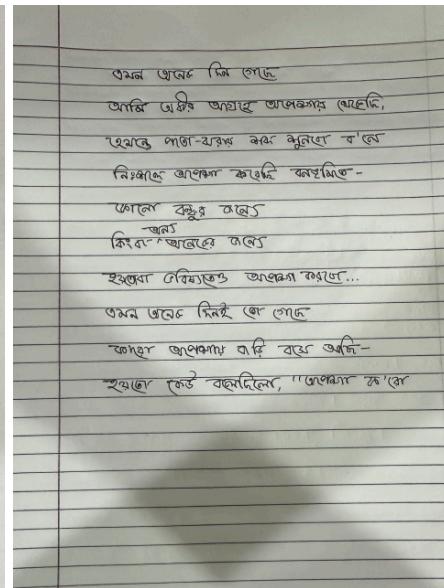


Figure 20. *
(c) Ruled Paper

Figure 21. Dataset 5 (D5): Clear handwriting achieving the highest recognition accuracy (71.2%)

APPENDIX E SAMPLE TOKENIZATION OUTPUT

Below is a sample output from our tokenization pipeline demonstrating both word and character-level tokenization:

```

1 Input text length: 1106 characters
2 Input text preview: ৫ আগস্ট অন্তর্বর্তী সরকারের প্রধান
উপদেষ্টা অধ্যাপক মুহাম্মদ ইউনুস জাতীয় সংসদের দক্ষিণ প্লাজায় ...
3
4 SUMMARY STATISTICS
5
6 Word-level Tokenization:
7   • Total tokens: 185
8   • Unique tokens: 138
9   • Type-token ratio: 0.746
10
11 Character-level Tokenization:
12   • Total tokens: 789
13   • Unique tokens: 167
14   • Type-token ratio: 0.212
15
16 TOKEN TYPE DISTRIBUTION
17
18 Word-level:
19   • WORD : 159 ( 85.9%)
20   • PUNCTUATION : 20 ( 10.8%)
21   • OTHER : 6 ( 3.2%)
22
23 Character-level:
24   • BANGLA_DIGIT : 9 ( 1.1%)
25   • WHITESPACE : 159 ( 20.2%)
26   • INDEPENDENT_VOWEL : 46 ( 5.8%)
27   • CONSONANT : 204 ( 25.9%)
28   • COMPOUND : 68 ( 8.6%)
29   • DEPENDENT_VOWEL : 54 ( 6.8%)
30   • CONSONANT_VOWEL : 196 ( 24.8%)
31   • MODIFIER : 17 ( 2.2%)
32   • OTHER : 16 ( 2.0%)
33   • PUNCTUATION : 20 ( 2.5%)

```

34
35 MOST FREQUENT TOKENS
36

37 Word-level (Top 10):

- 38 1. ' , ' → 14 times
- 39 2. 'করা' → 7 times
- 40 3. '।' → 6 times
- 41 4. 'ও' → 3 times
- 42 5. 'হয়েছে' → 3 times
- 43 6. '-' → 3 times
- 44 7. 'পরবর্তী' → 3 times
- 45 8. 'উল্লেখ' → 3 times
- 46 9. 'ড' → 2 times
- 47 10. 'আগস্ট' → 2 times

48
49 Character-level (Top 15):

- 50 1. ' ' → 151 times
- 51 2. 'র' → 41 times
- 52 3. 'ক' → 21 times
- 53 4. 'ে' → 20 times
- 54 5. 'স' → 19 times
- 55 6. 'ন' → 18 times
- 56 7. 'ঁ' → 17 times
- 57 8. 'ঁ' → 17 times
- 58 9. 'ঁ' → 17 times
- 59 10. 'ঁ' → 16 times