



# DNN-HMM BASED LARGE VOCABULARY ONLINE HANDWRITTEN ASSAMESE WORD RECOGNITION SYSTEM



{S. MANDAL, H. CHOUDHURY, SRM PRASANNA & S. SUNDARAM} IIT GUWAHATI, INDIA

## OBJECTIVES

1. Development of large vocabulary online handwriting Assamese word recognition system.
2. The use of DNN-HMM framework for on-line handwriting recognition.
3. Selection of 173 basic units that can characterizes 20K most frequent Assamese words.
4. Creation of Assamese word database.

## INTRODUCTION

- The majority of the studies in Indic script considered to recognize isolated characters.
- Few works on word recognition task has reported on some of the Indic scripts such as Devanagari, Bangla and Tamil [1].
- The number of characters ranges from few hundreds to several thousand.
- The challenges include lack of common definition in basic unit (BU), availability database etc.

## ASSAMESE SCRIPT AND DATABASE CREATION

অ (1)	আ (2)	ই (3)	ঈ (4)	উ (5)	ঊ (6)	ঋ (7)	ঌ (8)	ঐ (9)
ঔ (10)	ঋ (11)	ঌ (12)	য (13)	গ (14)	ঘ (15)	ঙ (16)	চ (17)	ছ (18)
জ (19)	ঝ (20)	ট (21)	ঠ (22)	ড (23)	ঢ (24)	ণ (25)	ত (26)	থ (27)
দ (28)	ধ (29)	ন (30)	প (31)	ফ (32)	ব (33)	ভ (34)	ম (35)	য (36)
ৰ (37)	ল (38)	ৱ (39)	শ (40)	ষ (41)	স (42)	হ (43)	ক্ষ (44)	য় (45)
ড় (46)	ঢ় (47)	ৎ (48)	ং (49)	ঙ (50)	ঞ (51)	ঠ (52)	ড (53)	ণ (54)
঵ (55)	শ (56)	ষ (57)	স (58)	হ (59)	঺ (60)	঻ (61)	় (62)	ঽ (63)
ি (64)	ি (65)	ি (66)	ি (67)	ি (68)	ি (69)	ি (70)	ি (71)	ি (72)
ি (73)	ি (74)	ি (75)	ি (76)	ি (77)	ি (78)	ি (79)	ি (80)	ি (81)
ি (82)	ি (83)	ি (84)	ি (85)	ি (86)	ি (87)	ি (88)	ি (89)	ি (90)
ি (91)	ি (92)	ি (93)	ি (94)	ি (95)	ি (96)	ি (97)	ি (98)	ি (99)
ি (100)	ি (101)	ি (102)	ি (103)	ি (104)	ি (105)	ি (106)	ি (107)	ি (108)
ি (109)	ি (110)	ি (111)	ি (112)	ি (113)	ি (114)	ি (115)	ি (116)	ি (117)
ি (118)	ি (119)	ি (120)	ি (121)	ি (122)	ি (123)	ি (124)	ি (125)	ি (126)
ি (127)	ি (128)	ি (129)	ি (130)	ি (131)	ি (132)	ি (133)	ি (134)	ি (135)
ি (136)	ি (137)	ি (138)	ি (139)	ি (140)	ি (141)	ি (142)	ি (143)	ি (144)
ি (145)	ি (146)	ি (147)	ি (148)	ি (149)	ি (150)	ি (151)	ি (152)	ি (153)
ি (154)	ি (155)	ি (156)	ি (157)	ি (158)	ি (159)	ি (160)	ি (161)	ি (162)
ি (163)	ি (164)	ি (165)	ি (166)	ি (167)	ি (168)	ি (169)	ি (170)	ি (171)
ি (172)	ি (173)							

- (a) ি + ক = কি (s53) (s12) (CV unit)
- (b) গ + ঞ = গ্ৰ (s14) (s165) (CC unit)
- (c) য় + ি = য়ি (s71) (s166) (New conjunct)
- (d) ষ + ঙ = ষ্ণ (s41) (s169) (New conjunct)

Figure 2: Generation of new basic unit (BU) by combining the BUs defined in Fig. 1

Set	Frequency (%)
Vowels	5.96
Vowels+Consonants	61.49
Vowels+Consonants+Conjuncts	64.16
Vowels+Consonants+Conjuncts+Modifiers (All BUs + their combination of Fig. 1)	99.76

Table 1: Number of characters (in %) covered by the selected BU set of Fig. 1 in Assamese OCR [2] corpus.

The set of 182 words is considered for data collection and it covers all the BUs in Fig. 1 and the characters formed by the BUs.

Set	# Writer	# Samples
Training	96	16208
Testing	48	8414
Validation	19	2941

Table 2: Samples and writers in Assam. word database.

## ASSAMESE HANDWRITTEN WORD RECOGNITION SYSTEM

- The preprocessing steps are smoothing, resampling, and size normalization.
- The following 16 features are considered: pre-processed  $x$  and  $y$  coordinates (2); first derivative of  $x$  and  $y$  coordinates (2); second derivative of  $x$  and  $y$  coordinates (2); writing direction (2); curvature (2); aspect ratio (1); linearity (1); slope (1); context map (1); ascender (1); and descender (1).
- In the proposed system, a deep neural network is trained to output HMM state probabilities.
- For training, GMM-HMM system is used to obtain alignment of the training data to HMM states.

- Next, we supply the label of each frame to the neural network for DNN training.

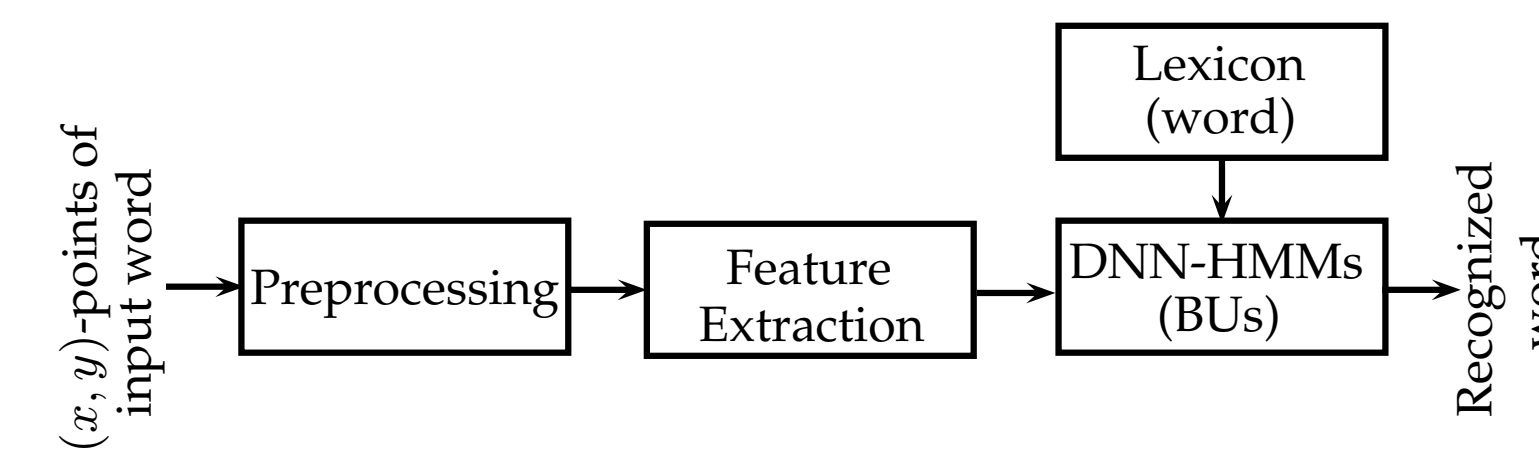


Figure 4: Block schematic of the developed online handwritten Assamese word recognition system.

- The word with highest score among all the lexicon words is declared as recognized word.

Figure 1: The basic unit (BU) set defined for recognizing 20K most frequent Assamese words.

Fig. 1 presents the identified BU that are used to write Assamese. The BUs 1—11 are vowels, BUs 12—51 are consonants, BUs 52—60 are vowel modifiers, BUs 61 and 62 are consonant modifiers, BUs 63—163 and 171—173 are conjuncts, and BUs 164—170 are few strokes.

## RESULT AND DISCUSSION

### Evaluation Methodology:

- Lexicon: 1K, 5K, 10K and 20K (close set).
- The most frequent of Assamese OCR is considered to create the lexicon.
- Accuracy =  $100 \times (t/n)$ ,  $t$  is the number of correct classification and  $n$  is the number of samples tested.
- The Kaldi toolkit is used for DNN-HMM.
- A left-to-right topology with 20 GMMs in each state is used.
- Highest accuracy is obtained for 17 states.

HMM State	GMM-HMM	DNN-HMM (# Hidden layer conf.)		
		L1:300	L2:300-300	L3:300-300-300
13	89.10	93.70	94.15	94.46
15	90.04	94.33	94.42	93.93
17	90.67	94.69	94.69	94.63
19	89.24	94.60	94.62	94.24

Table 3: Tuning of GMM-HMM and DNN-HMM models on validation set, considering 1K lexicon.

- DNN input consists of  $(2L + 1)$  feature vectors and chosen  $L=4$  as shown in Fig. 3.
- DNN input node = 144; output node = 2941 (173 HMMs with 17 states in each model).

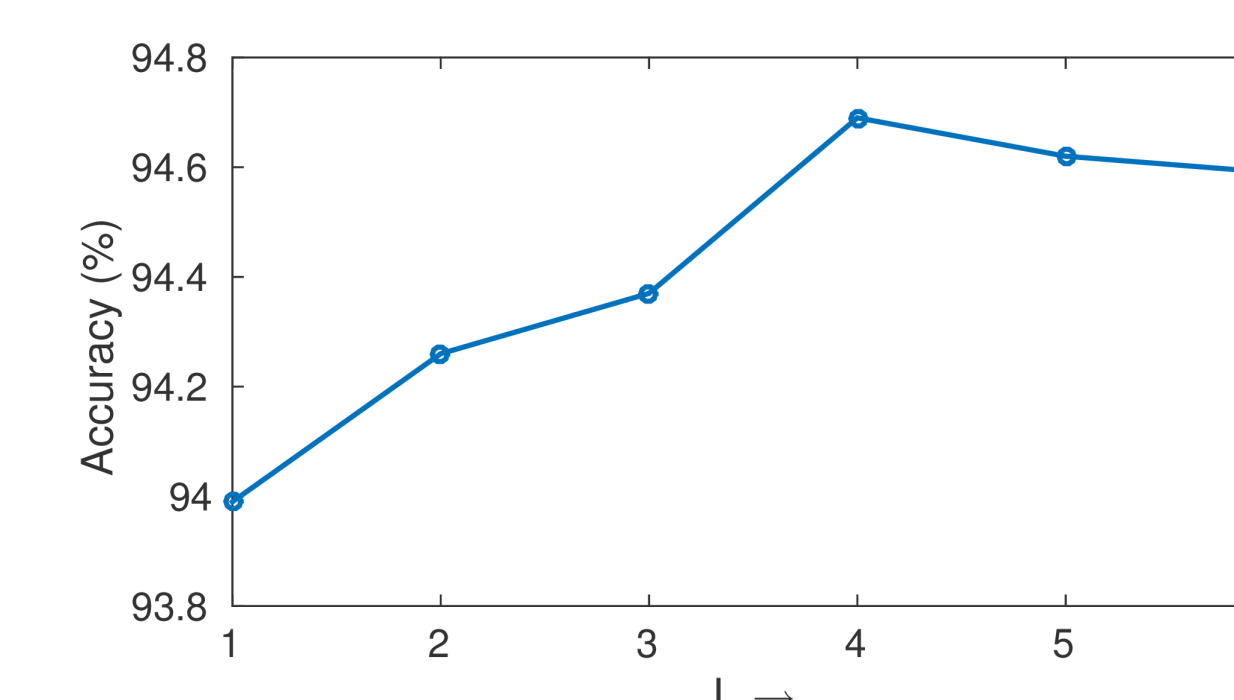


Figure 3: Tuning of the length of context window ( $L$ ) for DNN training on validation set.

- Best accuracy is obtained for two hidden layers as given in Table 3.

Lexicon size (# word)	GMM-HMM	DNN-HMM
1K	81.15	89.75
5K	77.08	88.72
10K	73.87	88.03
20K	71.13	87.36

Table 4: Word recognition accuracy (in %) of the DNN-HMM system and GMM-HMM system with various lexicon word on test set.

- Proposed DNN-HMM system significantly outperform GMM-HMM system for all lexicon size.

## FUTURE RESEARCH

- The other preprocessing steps such as shirorekha detection can be incorporated.
- Offline features can also be explored to make the model less sensitive to the temporal changes of the data.

## REFERENCES

[1] S. Bhattacharya, D. S. Maitra, U. Bhattacharya, and Swapan K Parui. An end-to-end system for Bangla online handwriting recognition. In *Proc. of Int. Conf. on Frontiers in Handwriting Recognition*, pages 373–378. IEEE, 2016.

[2] Subhankar Ghosh, P. K. Bora, Sanjib Das, and B. B. Chaudhuri. Development of an Assamese OCR using Bangla OCR. In *Proc. of the Workshop on Document Analysis and Recognition*, pages 68–73. ACM, 2012.

## CONTACT INFORMATION

**Subhasis Mandal**, Dept. of EEE, Indian Institute of Technology Guwahati, India - 781039  
**Email** subhasis.mandal@iitg.ac.in  
**Phone** +91-9085285110