

# Evolution of the Modern Phase of Written Bangla: A Statistical Study

Paheli Bhattacharya  
pahelibhattacharya@gmail.com  
Govt. College of Engineering  
and Textile Technology,  
Serampore, Hooghly,  
India.

Arnab Bhattacharya  
arnabb@iitk.ac.in  
Dept. of Computer Science  
and Engineering,  
Indian Institute of Technology, Kanpur,  
India.

## Abstract

Active languages such as Bangla (or Bengali) evolve over time due to a variety of social, cultural, economic, and political issues. In this paper, we analyze the change in the written form of the modern phase of Bangla quantitatively in terms of character-level, syllable-level, morpheme-level and word-level features. We collect three different types of corpora—classical, newspapers and blogs—and test whether the differences in their features are statistically significant. Results suggest that there are significant changes in the length of a word when measured in terms of characters, but there is not much difference in usage of different characters, syllables and morphemes in a word or of different words in a sentence. To the best of our knowledge, this is the first work on Bangla of this kind.

## 1 Introduction

Bangla (or Bengali) is one of the most widely spoken languages. It belongs to the Indo-European family of languages and is believed to have been derived from Prakrit in around 650 CE. The history of Bangla is divided into three phases: Old Bangla (till 1350 CE), Medieval Bangla (1350-1800 CE) and Modern Bangla (1800 CE-).<sup>1</sup>

Since its inception, Bangla, like any other active language, has undergone a lot of changes due to a variety of social, cultural, economic and political causes. The changes happen mostly in vocabulary and pronunciation, one of the big catalysts for which is the adoption of words of foreign origin either directly or indirectly into the

language. For example, in Bangla, there is no word that depicts the concept “football” directly, and consequently, the English word has been adopted verbatim and has become part of the language now. Similarly, the English word “box” has been incorporated in Bangla as *বাক্স* (bAksa<sup>2</sup>) by suitably modifying its pronunciation.

A particularly remarkable source of variety in Bangla is the two clearly distinct forms of written prose in the modern phase – *Sadhu Bhasha* (chaste language) and *Chalit Bhasha* (colloquial language). The chaste language was used earlier (by the likes of Bankimchandra Chattopadhyay, Rabindranath Tagore, Saratchandra Chattopadhyay and others) and has been now replaced in almost all communications in Bangla by the colloquial version. The most notable change has happened in the form of verbs and pronouns which has become shorter and can be more easily pronounced. For example, the verb *করিয়াছি* (kariAChi) has become *করেছি* (karaChi) and the pronoun *তাহাদের* (tahadera) has been transformed to *তাদের* (tader).

With the advancement of digital world, the electronic media have imparted a large impact on the modern language which is clearly reflected in newspapers, blogs and social networking forums. It is extremely rare to find longer words such as *যৌবনতেজোদীপ্ত* (JaubanatejodIpta) now and *চৈতানিশীথশশী* (chaitaranishIthashashi) than in the classical literature.

However, while all these notions of change are commonly believed to be true, to the best of our knowledge, there is no work that tests whether these perceptions about the differences are *statistically*

<sup>1</sup>The history and genesis of the language can be found in [http://www.bpedia.org/B\\_0137.php](http://www.bpedia.org/B_0137.php).

<sup>2</sup>We have used the ITRANS transliteration mechanism to specify the words in Bangla font (<http://www.aczoom.com/itrans/>). The rules for Bangla are available at <http://www.aczoom.com/itrans/html/beng/node4.html>.

*significant*. In this paper, we precisely aim to fill this gap. Our main contribution, thus, is to study the changes in the modern phase of written Bangla in a statistically robust manner.

We collect three different corpora – one consisting of classical literature, and the other two that of newspapers and blogs (the details are in Section 3). We then extract different features at the word and sentence levels and test whether the changes across the corpora are significant when viewed from a statistical standpoint.

The rest of the paper is organized as follows. Section 2 discusses the related work. Section 3, Section 4 and Section 5 describe the corpora, the features and the statistical testing method respectively. Section 6 discusses the results before Section 7 concludes.

## 2 Related Work

*Evolutionary linguistics* or the study of evolution of languages has long fascinated human beings. In addition to numerous studies that have been developed, there are whole conferences—EvoLang, the Evolution of Language International Conferences (evolang.org)—that are devoted for this. Among the computational studies for language evolution and change, an overwhelming majority of the work is focused on European languages [7, 2, 5].

Indian languages, due to the relative paucity of digital resources, had not been studied deeply. The recent ease of using Unicode and the surge of excellent work in the field of natural language processing (NLP) have, however, changed the situation dramatically. Sikder analyzed the change in type of words in Bangla and showed how the frequency of foreign words are increasing, especially for young people living in urban areas [6]. Choudhury et al. studied the change of Bangla verb inflections for the single verb কর (kara), which means “to do” in English [1]. They gave a functional explanation for the rise in several dialects of Bangla because of phonological differences while uttering the verb inflections.

To the best of our knowledge, computational studies describing the changes in the Bangla language, however, have not been undertaken so far. The evolution of a language is most visible when changes occur in the discourse. Since we are still far off from that for Bangla, in this paper, we study some micro-level aspects of the written language in terms of characters, syllables, morphemes and words.

Corpus	Number of words
Classical	1,58,807
Newspaper	7,71,989
Blog	5,18,485

Table 1: Number of words in the corpora.

## 3 Corpora

For our work, we collected three different corpora:

1. *Classical Corpus*: It includes the literary works of 4 eminent authors.
2. *Newspaper Corpus*: It includes the news articles from 7 leading newspapers of both India and Bangladesh.
3. *Blog Corpus*: It includes blog articles (but not the comments) from 11 blogs.

Appendix A list the details of the three corpora. The total number of words in each of the corpus are listed in Table 1.

## 4 Features

### 4.1 Character-level Features

In Bangla, there are two types of characters—vowels and consonants. The consonants cannot be pronounced on their own and must always end with the sound of a vowel. Vowels, on the other hand, can be pronounced on their own and are written either as independent letters or as diacritical marks on the consonant they attach to. For example, ক্ (k) is a consonant. When it is joined with আ (A), it is written as কা (kA). Thus, the diacritical mark for the vowel আ (A) is ৃ (A<sup>3</sup>). The vowel অ (a) has an invisible diacritical mark. Its only effect is to remove the ্ (the consonant-ending marker) from the consonant it attaches to. Thus, ক্+অ (k + a) is written as ক (ka).

We distinguish the diacritical mark of a vowel from the vowel itself as the latter can stand on its own. For example, the correct parsing of খুশীতে (khushIte) is খ্+ু+শ্+ী+তে (kh+u+sh+I+t+e) and that of আলোক (Aloka) is আ+ল্+ো+ক্+া (A+l+o+k+a) where ্ is used to represent the invisible diacritical mark of the vowel অ (a). The four consonants, ত্, ন্, হ্, ্ (t,h, .n, H, .N respectively), are treated differently in that they do

<sup>3</sup>The ITRANS coding for the diacritical marks remain the same.

not have the consonant-ending marker  $\cdot$ . Thus, বাংলা (bA.nIA) is parsed as  $\bar{b}+I+\bar{n}+\bar{l}+I$  (b+A+.n+l+A).

Conjunct characters where two (or three) consonants are joined together are parsed differently. There is no vowel at the end of the first (respectively, the first two) and only the last consonant has a vowel ending written as a diacritical mark. Hence, the correct parsing of সন্ত্রস্ত (santrasta) is  $\bar{s}+\bar{t}+\bar{n}+\bar{t}+\bar{r}+\bar{a}+\bar{s}+\bar{t}+\bar{a}$ . (s+a+n+t+r+a+s+t+a).

#### 4.1.1 Character frequencies

We count the frequencies of all the characters—consonants, vowels and diacritical marks—using the parsing system discussed above for the three corpora. The number of distinct characters is 61 that includes the 39 consonants (the consonant  $\bar{b}$  (b) is counted only once), the 11 vowels (the vowel  $\bar{a}$  (no ITRANS code) is not used any more) and the corresponding 11 diacritical marks.

We also count the frequencies of bi-gram and tri-gram characters. For example, the bi-grams in the word বাংলা (bA.nIA) are  $\bar{b}A$  (bA),  $\bar{A}n$  (A.n),  $\bar{n}I$  (.nl) and  $\bar{l}A$  (IA). The tri-grams are extracted similarly.

We arrange the uni-gram characters (and bi-grams and tri-grams) in descending order of their frequencies. When comparing corpus  $C_1$  with  $C_2$ , we consider the top-50 entries from the sorted list of  $C_1$  and find their frequencies in  $C_2$ . Thus, the comparison of  $C_1$  with  $C_2$  differs from that of  $C_2$  with  $C_1$  as, in the later case, the frequencies of the top-50 entries of  $C_2$  are considered. The frequencies from the two corpora form the two non-parametric distributions between which the changes are statistically tested. Instead of using the raw counts as frequencies, we compute the relative ratios by dividing by the total number of characters in the corpus; this makes two corpora of differing sizes comparable.

#### 4.1.2 Character-based word length

For each corpus, we produce a count of words that have a particular length in terms of characters. Thus, if there are 300 words of length 4, the frequency corresponding to 4 in the non-parametric distribution is 300. The distribution consists of all the word lengths and their frequencies. The comparison of corpus  $C_1$  with  $C_2$  is symmetric for this feature.

## 4.2 Morpheme-level Features

A *morpheme* is the smallest meaning-bearing unit in a language. A morpheme may not be able to stand on its own, although a word necessarily does. Every word is

composed of a root word (sometimes called a *lexeme*) and possibly one or more morphemes.

To extract morphemes, we used the unsupervised program Undivide++<sup>4</sup>, which is based on the work by Dasgupta et al. [4]. Unfortunately, the program have many parameters, and even after repeated tuning and discussion with the authors, we could not replicate the accuracies as reported in [4] on our corpora for the 4110-word test-set provided by them. Although the program is only about 50% accurate on average, we still use it to extract all the morphemes from the words in the corpora. (Appendix B reports the performance on the metrics as proposed in [3].)

#### 4.2.1 Morpheme frequencies

For every morpheme, we get a count of words that have it. Similar to the character frequencies, we then extract the top-50 (normalized) frequencies from each corpus.

#### 4.2.2 Morpheme-based word length

Using the program Undivide++, every word is segmented into a list of prefix(es), root word and suffix(es). The “length” of the word is then counted as the number of such segments. For example, if প্রদেশটিকে (pradeshaTike) is segmented into the prefix  $\bar{p}r$  (pra), the root  $\bar{d}e$  (desha) and the two suffixes  $\bar{t}i$  (Ti) and  $\bar{k}e$  (ke), its length is counted as 4.

## 4.3 Syllable-level Features

*Syllables* are the smallest subdivisions uttered while pronouncing a word. Since syllables are phonetic units, they cannot be extracted completely correctly without speech analysis. To bypass the problem, we employ a very simple and intuitive heuristic which is almost always correct.

We assume that any combination of characters till the next vowel is a syllable. Thus, each vowel, each consonant with its vowel ending (encoded as a diacritical mark), and each conjunct character is a separate syllable.

The consonants,  $\bar{t}h$ ,  $\bar{n}$ ,  $\bar{H}$  (t.h, .n, H respectively), are treated as single syllables since they do not have the consonant-ending marker  $\cdot$ . However,  $\bar{N}$  (.N) is considered part of the preceding syllable. Thus, the word অকস্মাৎ (akasmAt.h) has three syllables  $\bar{a}$  (a),  $\bar{k}a$  (ka),  $\bar{s}mA$  (smA) and  $\bar{t}h$  (t.h) while বাঁধা (bA.NdhA) has two syllables  $\bar{b}A$  (ba.N) and  $\bar{d}hA$  (dhA).

<sup>4</sup>Available from <http://www.hlt.utdallas.edu/~sajib/Morphology-Software-Distribution.htm>

### 4.3.1 Syllable frequencies

For every uni-gram syllable (and also bi-grams of syllables), we get a count of words that have it. We again consider only the top-50 (normalized) frequencies from each corpus.

### 4.3.2 Syllable-based word length

Similar to characters, the word length is also counted in terms of syllables. The “inverted list”, i.e., the number of words having a particular syllable-length is then used as the feature for that syllable length.

## 4.4 Word-level Features

The words are parsed from the sentences using orthographic word boundaries (i.e., the white-space characters including `?`, `!`, `.` and the Bangla character `।`).

### 4.4.1 Word frequencies

The words are for sentences what the characters are for words. Thus, this feature is computed in exactly the same way as characters.

### 4.4.2 Word-based sentence length

Similar to word length, the sentence length is counted in terms of number of words.

## 5 Statistical Testing

All the features that are used in this paper are summarized in Table 2.

To test whether the distributions of the various features for the different corpora are statistically different from each other, we employ the non-parametric two-sample *Kolmogorov-Smirnov (K-S) test*<sup>5</sup>. For each pair of corpora, we perform three tests. Suppose the corpora are  $C_1$  and  $C_2$ . The *null hypothesis*  $H_0$  for all the three tests state that the samples observed empirically for  $C_1$  and  $C_2$  come from the *same* distribution.

There can be three ways by which the *alternate hypothesis* can vary. For the non-equal ( $\neq$ ) test, the alternate hypothesis  $H_A^\neq$  states that the empirical values  $x_i^{(1)}$  and  $x_i^{(2)}$  for the distributions from  $C_1$  and  $C_2$  are different, i.e., for every  $i$ ,  $x_i^{(1)} \neq x_i^{(2)}$ . For the greater than ( $>$ ) test, the alternate hypothesis  $H_A^>$  states that the empirical values for  $C_1$  are greater than the corresponding

values for  $C_2$ , i.e., for every  $i$ ,  $x_i^{(1)} > x_i^{(2)}$ . The less than ( $<$ ) test is similar where the alternate hypothesis  $H_A^<$  tests whether for every  $i$ ,  $x_i^{(1)} < x_i^{(2)}$ .

The K-S test returns a *p-value* that signifies the confidence with which the null hypothesis can be rejected. The lower the p-value, the more statistically significant the result is. Thus, for the  $H_A^\neq$  case, it means the two distributions are more different. If the result of a  $\neq$  test is statistically significant at a particular level of significance, then the result of either the  $>$  test or the  $<$  test (but not both) must be significant as well at the same level of significance.

## 6 Results

The differences between the word lengths in terms of number of characters between the three corpora are found to be statistically significant<sup>6</sup> for the alternate hypothesis  $H_A^\neq$ . (The tables in Appendix C list all the p-values.) More interestingly, the alternate hypothesis  $H_A^<$  is found to be very significant for classical versus blog, classical versus newspaper and blog versus newspaper comparisons. This shows that the frequency of words having a shorter length is less in classical than in blogs which, in turn, is less than newspapers. Thus, this shows that longer words were more common in the classical literature than in newspapers which are more than that in blogs.

Although the classical corpus exhibits longer words in terms of syllables (due to the  $H_A^<$  test), the non-equality test ( $H_A^\neq$ ) is not significant. This, thus, indicates that the use of conjunct characters were more in classical literature which led to longer words in terms of characters but not in terms of syllables. The blogs and newspapers differ in terms of number of syllables though.

The differences in number of morphemes is again not significant. Thus, contrary to popular perception, words with many suffixes and prefixes are not more abundant in the classical literature as compared to the current scenario. Similarly, the number of words per sentence for classical is not statistically different either.

Frequencies of uni-gram characters, bi-gram characters and uni-gram syllables are not significantly different across the corpora. Frequencies of tri-gram characters, bi-gram syllables, uni-gram words and bi-gram words of classical are significantly different from both blogs and newspapers for the alternate hypotheses  $H_A^\neq$  and  $H_A^<$ .

<sup>5</sup>Unless otherwise mentioned, we consider the level of significance to be 5%, i.e., a result is statistically significant when the p-value of the test is less than or equal to 0.05.

<sup>5</sup>We use the Octave software to perform the tests.

Feature type	Level			
	Character	Syllable	Morpheme	Word
Uni-gram frequency	yes	yes	yes	yes
Bi-gram frequency	yes	yes	no	yes
Tri-gram frequency	yes	no	no	no
Length of word or sentence	yes	yes	yes	yes

Table 2: Features used.

The newspaper and blog corpora show little statistical difference in frequencies indicating that the current styles of formal and informal writing are quite alike.

## 7 Conclusions

In this paper, we provided a model of statistically testing the differences of writing styles across various phases of a language. To the best of our knowledge, this is the first work of its kind in Bangla. This work has aimed at building a basic foundation on which more analysis in terms of higher-level features can be carried out in the future. Also, bigger corpora will allow robust and more detailed analyses of the results.

## References

- [1] M. Choudhury, V. Jalan, S. Sarkar, and A. Basu. Evolution, optimization, and language change: The case of Bengali verb inflections. In *ACL SIG Computational Morphology and Phonology*, pages 65–74, 2007.
- [2] M. Christiansen. *Language evolution*. Oxford University Press, 2003.
- [3] S. Dasgupta and V. Ng. Unsupervised morphological parsing of Bengali. *Language Resources and Evaluation*, 40(3-4):311–330, 2006.
- [4] S. Dasgupta and V. Ng. High-performance, language-independent morphological segmentation. In *HLT-NAACL*, pages 155–163, 2007.

- [5] P. Niyogi. *The Computational nature of language learning and evolution*. MIT Press, 2006.

- [6] S. Sikder. Contemporary bengali language. Amor Ekushey: <http://archive.thedailystar.net/suppliments/2013/Amor%20Ekushey%20Thursday,February21,2013>.

- [7] L. Steels. The synthetic modeling of language origins. *Evolution of Communication*, 1(1):1–34, 1997.

## Appendices

### A Corpora

Table 3, Table 4 and Table 5 list the details of the three corpora.

### B Morphological Parsing

Table 6 shows the performance of Undivide++ on our corpora.

### C Results

Table 7 to Table 10 show the p-values of all the symmetric tests. Table 11 to Table 18 show the p-values of all the non-symmetric tests.

Author	URL	On
Rabindranath Tagore	www.rabindra-rachanabali.nltr.org	14 <sup>th</sup> July, 2013
Bankimchandra Chattopadhyay	www.bankim.rachanabali.nltr.org	14 <sup>th</sup> July, 2013
Saratchandra Chattopadhyay	www.sarat-rachanabali.nltr.org	14 <sup>th</sup> July, 2013
Swami Vivekananda	www.dduttamajumder.org/baniorachana	14 <sup>th</sup> July, 2013

Table 3: Classical corpus.

Name	Website	From	To
Anandabazar Patrika	www.anandabazar.com	1 <sup>st</sup> June, 2011	12 <sup>th</sup> July, 2013
Akhon Samay	www.akhonsamoy.com	14 <sup>th</sup> July, 2013	14 <sup>th</sup> July, 2013
Jana Kantha	www.dailyjanakantha.com	1 <sup>st</sup> January, 2010	13 <sup>th</sup> July, 2013
Inqilab	www.dailyinqilab.com	1 <sup>st</sup> June, 2013	11 <sup>th</sup> July, 2013
Jugantor	www.jugantor.com	1 <sup>st</sup> July, 2013	12 <sup>th</sup> July, 2013
Naya Diganta	www.dailynayadiganta.com	30 <sup>th</sup> June, 2013	14 <sup>th</sup> July, 2013
Pratham Alo	www.prothom-alo.com	1 <sup>st</sup> January, 2007	10 <sup>th</sup> July, 2013

Table 4: Newspaper corpus.

Name	Blog	On
AmarBlog	www.amarblog.com	8 <sup>th</sup> July, 2013
Bokolom	www.bokolom.com	10 <sup>th</sup> July, 2013
CoffeeHouserAdda	www.coffeehouseradda.in	7 <sup>th</sup> July, 2013
CadetCollege	www.cadetcollegeblog.com	16 <sup>th</sup> July, 2013
ChoturMatrik	www.choturmatrik.net	7 <sup>th</sup> July, 2013
MuktoBlog	www.muktoblog.net	9 <sup>th</sup> July, 2013
MuktoMona	www.mukto-mona.com	10 <sup>th</sup> July, 2013
NagarikBlog	www.nagorikblog.com	9 <sup>th</sup> July, 2013
Nirman	www.nirmaaan.com	8 <sup>th</sup> July, 2013
Sachalayatan	www.sachalayatan.com	8 <sup>th</sup> July, 2013
SomeWhereIn	www.somewhereinblog.net	15 <sup>th</sup> July, 2013

Table 5: Blog corpus.

Corpus	Accuracy	Recall	Precision	F-Score
Classical	48.80%	40.00%	49.78%	44.38%
Blog	55.60%	48.80%	53.53%	51.05%
Newspaper	54.30%	47.70%	52.68%	50.06%
Merged	56.40%	50.31%	54.00%	52.08%

Table 6: Performance of Undivide++ [4] on our corpora.

	Blog			Newspaper		
	≠	>	<	≠	>	<
Classical	1.66E-2	8.25E-1	8.33E-3	6.56E-4	9.14E-1	3.28E-4
Blog	-	-	-	2.09E-2	9.11E-1	3.28E-4

Table 7: K-S test results for frequency of characters per word.

	Blog			Newspaper		
	≠	>	<	≠	>	<
Classical	8.96E-2	4.43E-2	2.84E-2	3.03E-1	9.71E-1	1.52E-1
Blog	-	-	-	2.82E-3	9.74E-1	1.41E-3

Table 8: K-S test results for frequency of syllables per word.

	Blog			Newspaper		
	≠	>	<	≠	>	<
Classical	9.79E-1	8.94E-1	6.41E-1	9.79E-1	8.94E-1	6.41E-1
Blog	-	-	-	9.99E-1	8.94E-1	8.94E-1

Table 9: K-S test results for frequency of segments (morphemes plus root word) per word.

	Blog			Newspaper		
	≠	>	<	≠	>	<
Classical	9.97E-1	7.26E-1	7.26E-1	9.99E-1	8.35E-1	8.35E-1
Blog	-	-	-	8.64E-1	4.86E-1	6.06E-1

Table 10: K-S test results for frequency of words per sentence.

	Classical			Blog			Newspaper		
	≠	>	<	≠	>	<	≠	>	<
Classical	-	-	-	4.63E-1	2.37E-1	7.93E-1	0.60E-1	3.09E-1	6.96E-1
Blog	9.97E-1	7.30E-1	7.30E-1	-	-	-	7.30E-1	8.38E-1	7.30E-1
Newspaper	9.67E-1	6.12E-1	7.30E-1	9.97E-1	7.30E-1	8.38E-1	-	-	-

Table 11: K-S test results for frequency of uni-grams of characters.

	Classical			Blog			Newspaper		
	≠	>	<	≠	>	<	≠	>	<
Classical	-	-	-	6.03E-2	6.39E-1	3.01E-2	1.52E-1	7.15E-1	7.64E-2
Blog	1.32E-2	2.04E-2	6.60E-3	-	-	-	1.32E-2	2.04E-1	6.60E-3
Newspaper	1.11E-11	1.0	3.06E-17	1.64E-1	6.70E-1	8.20E-2	-	-	-

Table 12: K-S test results for frequency of bi-grams of characters.

	Classical			Blog			Newspaper		
	≠	>	<	≠	>	<	≠	>	<
Classical	-	-	-	2.16E-5	5.25E-1	1.08E-9	7.05E-8	7.64E-2	3.52E-8
Blog	2.48E-5	3.63E-2	1.24E-5	-	-	-	6.25E-5	9.32E-2	3.12E-5
Newspaper	2.48E-5	3.63E-2	1.24E-5	6.25E-5	9.32E-2	3.12E-5	-	-	-

Table 13: K-S test results for frequency of tri-grams of characters.

	Classical			Blog			Newspaper		
	≠	>	<	≠	>	<	≠	>	<
Classical	-	-	-	1.12E-1	7.26E-1	5.61E-2	1.77E-1	7.26E-1	8.89E-2
Blog	7.27E-2	8.38E-1	3.00E-2	-	-	-	7.20E-1	8.38E-1	3.82E-1
Newspaper	2.80E-1	7.30E-1	1.40E-1	2.80E-1	6.12E-1	1.40E-1	-	-	-

Table 14: K-S test results for frequency of uni-grams of syllables.

	Classical			Blog			Newspaper		
	≠	>	<	≠	>	<	≠	>	<
Classical	-	-	-	4.44E-16	1	2.09E-16	4.44E-16	1	2.09E-16
Blog	2.22E-2	5.61E-2	1.11E-2	-	-	-	2.22E-2	1.11E-2	5.61E-2
Newspaper	9.19E-8	6.06E-1	4.95E-8	1.98E-5	1.35E-1	9.92E-6	-	-	-

Table 15: K-S test results for frequency of bi-grams of syllables.

	Classical			Blog			Newspaper		
	≠	>	<	≠	>	<	≠	>	<
Classical	-	-	-	1.4E-2	6.06E-1	7.31E-4	2.95E-4	7.26E-1	1.47E-4
Blog	6.17E-3	3.08E-3	6.06E-1	-	-	-	8.64E-1	8.35E-1	4.86E-1
Newspaper	4.44E-16	2.09E-16	1	4.44E-16	2.09E-16	9.75E-1	-	-	-

Table 16: K-S test results for frequency of uni-grams of morphemes.

	Classical			Blog			Newspaper		
	≠	>	<	≠	>	<	≠	>	<
Classical	-	-	-	3.33E-16	2.65E-1	1.87E-16	3.33E-16	9.40E-2	1.87E-16
Blog	1.01E-12	8.38E-1	5.05E-13	-	-	-	1.37E-7	4.93E-1	6.89E-8
Newspaper	1.23E-13	6.06E-1	6.14E-14	0	1	1.92E-22	-	-	-

Table 17: K-S test results for frequency of uni-grams of words.

	Classical			Blog			Newspaper		
	≠	>	<	≠	>	<	≠	>	<
Classical	-	-	-	1.37E-7	3.82E-1	6.92E-8	7.40E-11	3.82E-1	3.70E-11
Blog	4.33E-8	2.04E-1	2.16E-8	-	-	-	7.40E-11	7.30E-1	3.70E-11
Newspaper	1.01E-12	7.30E-1	5.05E-13	2.48E-5	3.82E-1	1.24E-5	-	-	-

Table 18: K-S test results for frequency of bi-grams of words.