# Complexity Metric for Code-Mixed Social Media Text

Souvick Ghosh[1], Satanu Ghosh[2], and Dipankar Das[3]

[1]Rutgers University, New Brunswick, NJ 08901, US
`{souvick.ghosh}@rutgers.edu`
[2]MAKAUT, Kolkata 700064, WB, India
`{satanu.ghosh.94}@gmail.com`
[3]Jadavpur University, Kolkata 700032, WB, India
`dipankar.dipnil2005}@gmail.com`

**Abstract.** An evaluation metric is an absolute necessity for measuring the performance of any system and complexity of any data. In this paper, we have discussed how to determine the level of complexity of code-mixed social media texts that are growing rapidly due to multilingual interference. In general, texts written in multiple languages are often hard to comprehend and analyze. At the same time, in order to meet the demands of analysis, it is also necessary to determine the complexity of a particular document or a text segment. Thus, in the present paper, we have discussed the existing metrics for determining the code-mixing complexity of a corpus, their advantages and shortcomings as well as proposed several improvements on the existing metrics. The new index better reflects the variety and complexity of a multilingual document. Also, the index can be applied to a sentence and seamlessly extended to a paragraph or an entire document. We have employed two existing code-mixed corpora to suit the requirements of our study.

## 1 Introduction

Social media text differs from other conventional text in many ways. It is noisy in nature and requires comprehensive processing from a multilingual point of view. In social media communication, a multilingual speaker often uses more than one language. Therefore, the communication, inherently informal in nature, presents the scientific community with a challenging yet interesting problem.

First of all, we need to understand the necessity of language switching. Is it motivational? Or is it circumstantial? Although mixing of languages was prevalent in verbal communication, it was the proliferation of social media which accelerated the use of multiple languages in a single written communication and this is motivated by both social and conversational needs [1]. Sometimes, the speaker is not competent in the language he is writing. The lack of vocabulary provokes him to use words from his native language as a substitute. On some other occasions, the need is purely social and is used by the writer to mark him as part of a large group.

Automatic identification of the code switching points is important as it helps to understand the frequency of code switching or code mixing and subsequent complexity of the text. Also, it would allow us to determine the language specific models which are better suited for the analysis of such text. It is also important to understand the

differences between code switching and code mixing as both the terms are used inter-changeably in the literature. In our work, the term 'Code- Switching' refers to inter-sentential language shifts while the term 'Code-Mixing' refers to intra-sentential shifts of language.

In the present work, we have used short utterances collected from Facebook pages and Twitter data for our analysis. As the dataset is purely based on Indian social media text, it is essential that we give a brief statistics about the degree of multilingualism in India. There are more than 20 officially recognized languages in India. The number of Hindi speakers range from 14.5% to 24.5% of total population (source:Wikipedia[1]). Other languages are spoken by 10% or less out of the total population. English and Hindi are mostly used for official communication in India. Similarly, it has been also observed that the diversity of Indian languages and the necessity for faster and efficient communication motivates the mixing of languages in Indian social media context. This brings us to one of the challenging problems, i.e. transliteration. Most of the time, the languages like Hindi or Bengali are not written using their native scripts - Devanagari for Hindi and Eastern Neo Brahmi script for Bengali. Instead, the users prefer using Roman script as it is more convenient with a regular keyboard. While analyzing a code-mixed transliterated text, it is often useful to determine the complexity of the corpus. For any task on code-mixed corpora such as language identification, part-of-speech tagging, information retrieval or question answering, it is important for the researchers to compare the difficulty of their work with regards to the level of language mixing in the text. Also, it is expected that along with the increasing complexity and more code-mixing in a text, the accuracy of text processing would decrease and the error rates would increase.

In Section 2, we have discussed the previous work done in related area. In Section 3, we provide a brief statistics of the code-mixed corpora and its preparation. The index is presented in Section 4. The working of the index is further elaborated using various cases and examples in Section 5. Section 6 contains the results and finally, Section 7 concludes the paper and discusses the scopes of future work.

## 2  Related Work

In 2001, Kilgarriff [3] discussed and pointed out that corpus linguistics do not have proper methods for comparing corpora. Most of the corpus descriptions are textual and based on the opinions of the researchers. Such impressions are highly subjective and not a proper measure of corpus similarity or complexity. In contrast, Whenever we work on a new corpus, the questions that are inevitably raised ae about the limitations and benefits of using that corpus. The size and homogeneity of the data are some of the factors which have been used intensively. However, such approaches are mainly word based and are applicable for monolingual texts.

Measuring corpus similarity has a wide array of applications. It has theoretical as well as research applications where one can judge the complexity of the dataset before performing their technical analysis. Also, one may want to replace a dataset with an-other. It is beneficial only if there is some way to determine whether the two datasets

---

[1] https://en.wikipedia.org/

are similar and comparable in terms of their complexity and usage. This would in turn help in inter-domain portability of NLP systems too.

To the best of our knowledge, Gambäck and Das [2] proposed the first index for code-mixed social media text. Termed as Code Mixing Index (CMI), the index tries to assess the level of code-switching in an utterance. The measure aimed at comparing one code-switched corpus with another. Gambäck and Das [2] worked on Hindi/Bengali-English Facebook data collected from various chat groups. The corpora introduced by them has 28.5% of the messages written in at least two languages. CMI can be described as the fraction of total words that belong to languages other than the most dominant language in the text,

$$CMI = 100 * \left[1 - \frac{max w_i}{n-u}\right] \text{, if n > u}$$
$$CMI = 0 \text{, if n=u}$$

where n-u is the sum of N languages present in the utterance of their respective number of words and $max\{w_i\}$ is the highest number of words belonging to a particular language where n is total number of tokens and u is the number of language independent tags. Gambäck and Das [2] averaged the CMI values for all the sentences to obtain 'CMI all' and for only the code-mixed sentences to obtain 'CMI mixed', respectively.

However, CMI considers only the fraction of words in the corpus which are code-switched. We have used CMI as an initial parameter and have suggested some improvements which would take into account the number of languages and the number of code-switching points present in the corpus.

## 3   Corpora Details

The present task requires a social media corpus which has diversity in terms of languages and code-mixed content. Forum for Information Retrieval Evaluation [2] (FIRE) organized a shared task on Mixed Script Information Retrieval. The data set used for training and test suited our purpose perfectly. Another shared task was organized by the Twelfth International Conference on Natural Language Processing [3] (ICON-2015). This data set was bilingual in nature and used code-mixed social media text. Therefore, We have modified these two corpora in order to accomplish our task.

### 3.1   FIRE 2015 Shared Task Corpus

We have modified the transliterated corpus which was provided by the organizers of FIRE 2015 Shared Task on Mixed Script Information Retrieval. The dataset contains 3701 sentences and 63526 word tokens. Each word may belong to one of the nine languages present in the entire dataset. The nine languages were Bengali (Bn), English (En), Gujarati (Gu), Hindi (Hi), Kannada (Ka), Malayalam (Ml), Marathi (Mr), Tamil (Ta) and Telugu (Te). The dataset is extremely multilingual in nature. The languages

---

[2] http://fire.irsi.res.in/fire/2015/home
[3] http://ltrc.iiit.ac.in/icon2015/

present in the dataset are the most prevalent ones that we can find in Indian social media context. It must be noted that the words of a single query usually came from 1 or 2 languages and rarely from 3 different languages. This is in line with the language mixing trends that we have witnessed in social media context. The users, even if familiar with multiple languages, rarely use more than three languages while writing their posts or tweeting. As a matter of fact, most of the sentences are bilingual in nature with one of the dominant languages as either English or Hindi. Thus, there are sentences that mix Tamil and English words, or Bengali and Hindi words, but not for example, Gujrati and Kannada words. The named entities (marked as NE), language independent words (marked as X) and mixed words containing intra-word language switches (marked as MIX), were all considered undefined and assigned with UN (universal) tag. The numbers of utterances, tokens for each of the language pairs in the training set are given in Table 1.

| Language Tags | # Sentences | # Words | Percentage (%) Of Corpus |
|---|---|---|---|
| English (En) | 2665 | 21996 | 34.63 |
| Bengali (Bn) | 355 | 4919 | 7.74 |
| Gujarati (Gu) | 165 | 1075 | 1.69 |
| Hindi (Hi) | 614 | 5897 | 9.28 |
| Kannada (Ka) | 373 | 2212 | 3.48 |
| Malayalam (Ml) | 151 | 1390 | 2.19 |
| Marathi (Mr) | 229 | 2414 | 3.8 |
| Tamil (Ta) | 342 | 3694 | 5.82 |
| Telugu (Te) | 603 | 7002 | 1.1 |
| Language Independent | 2582 | 12927 | 20.35 |

Table 1: Statistics of FIRE 2015 Corpus.

### 3.2 ICON 2015 Shared Task Corpus

Another recent shared task was conducted by Twelfth International Conference on Natural Language Processing (ICON-2015), for part-of-speech tagging of transliterated social media text. In the shared task, the code-mixed data was collected from Bengali-English Facebook chat groups. The sentences are in mixed English-Bengali and English-Hindi - and have been obtained from the "JU Confession" Facebook group, which contains posts in English-Bengali with a few Hindi words in some cases.

We have modified the ICON Shared Task Corpora for developing our index metric. The dataset contains three languages - Bengali, Hindi and English. It contains 2341 sentences and 38199 word tokens in total. The statistics for the dataset have been presented in Table 2.

| Language Tags | Number of Sentences | Number of Words Present | Percentage (%) Of Corpus |
|---|---|---|---|
| English (En) | 1563 | 15435 | 40.41 |
| Bengali (Bn) | 1059 | 13002 | 34.04 |
| Hindi (Hi) | 153 | 1006 | 2.63 |
| Language Independent | 2268 | 8756 | 22.92 |

Table 2: Statistics of ICON 2015 Corpus.

## 4 Complexity Factor (CF)

We introduce an index, termed hereafter as Complexity Factor (CF), to measure the complexity of a multilingual corpus. This index can be applied to any sentence, paragraph or document which contains multiple languages. The index uses the concept of CMI as proposed by Gambäck and Das [2] and makes some practical additions on it.

Complexity Factor(CF) considers three different aspects while analyzing any text - Language Factor (LF), Switching Factor (SF) and Mix Factor (MF). CF can be calculated for sentences and easily extended to paragraphs and entire documents. In the next section, we have proposed three variants of Complexity Factor. Complexity Factor 1 (henceforth mentioned as CF1) is a simple baseline which considers LF, SF and MF. Complexity Factor 2 (CF2) and Complexity Factor 3 (CF3) are the two indexes which have been carefully fine-tuned to efficiently represent the complexity of any transliterated text.

### 4.1 Language Factor (LF)

This factor represents the number of different languages present in a sentence as a fraction of the total number of words in the sentences. It is evident that if a sentence becomes more multilingual, the complexity increases manifold. For example, For any given sentence, Language Factor can be defined as,

$$LF = \frac{W}{N}$$

where W is the number of words and N is the number of distinct languages in the sentence.

Sentence 1: "*Boss, **<Bn> ajkal ki korchis </Bn>***? We have been getting no news about you!*" (English Translation: Boss, what are you doing in these days? We have been getting no news about you! )

Sentence 2: "***<Bn> Kal khela dekhli? </Bn>** What a game! <u><Hi> Virat ne toh kamaal kar diya! </Hi></u>*" (English translation: Did you watch the match yesterday? What a game! Virat was simply superb!)

Sentence 1 contains two languages, Bengali and English, while Sentence 2 contains three languages - English, Bengali and Hindi. In both the sentences, Bengali words are boldfaced and Hindi words are underlined. Language Factor is 6 for Sentence 1 (W=12, N=2) and 4 for Sentence 2 (W=12, N=3). It must be noted that longer the text block we

are considering, it has more probability of finding multiple languages in it. This factor is inversely proportional to Complexity Factor (CF) and rewards shorter sentences with more distinct languages in it. The LF can range from W (in case of a monolingual text) to 1 (when each word belongs to a different language).

### 4.2   Switching Factor (SF)

It is essential to consider the number of times the writer switches from one language to the other. As the number of switches increases, it becomes more complex to analyze the text for various tasks like language identification, part-of-speech tagging, question-answering, summarization, etc.

For any given sentence, Switching Factor is defined as the ratio of number of switching points present in the sentence to the maximum number of switching points possible for that sentence. For a block of W words, the maximum number of code-switches occurs when each alternate words belong to different languages. So the maximum number of switching points for a W-word sentence is W-1. Switching Factor, denoted by SF, can be written as:

$$SF = \frac{S}{W-1} \text{ , if W > 1}$$

$$SF = 0 \text{ , if W = 1}$$

where S is the number of code-switches and W is the number of words in the sentences or block of text.

Consider the following examples,

Sentence 1: ***Ki** post **korcho**? Public forum **eta*** (English translation: What are you posting? This is a public forum)

Sentence 2: *It is painful **je khelata harlam*** (English translation: It is painful that we lost the match)

Both the sentences contain a mix of Bengali and English words (Bengali words are boldfaced). It should be noted that while both sentences contain 3 words each in Bengali and English, the relative arrangement of the words make Sentence 1 more complex than Sentence 2. For sentence 1, SF is 0.8 (S=4, W=6) while for sentence 2, it is only 0.2 (S=1, W=6). Thus, we can observe that Switching Factor captures this complexity factor and it is directly proportional to Complexity Factor (CF). For a single word sentence, SF = 0. SF can reach the maximum value of 1 when no two consecutive words belong to the same language.

### 4.3   Mix Factor (MF)

Mix Factor, referred to as MF for the rest of the paper, is based on Code Mixing Index (CMI). It is the ratio of number of words which are not written in the dominant language of the sentence to the total number of language-dependent words present in the sentence. It can be written as:

$$MF = \frac{W' - max\{w\}}{W'} \text{ , if W' > 0}$$

$$MF = 0 \text{ , if W' = 0}$$

where W' is the number of words in distinct languages, i.e., the number of words except the undefined ones, max{w} is the maximum number of words belonging to the most frequent language in the sentence.

Sentence 1: *"Boss, **ajkal ki korchis**? We have been getting no news about you!"* (English Translation: Boss, what are you doing these days? We have been getting no news about you! )

Sentence 2: *"**Kal khela dekhli**? What a game! Virat ne toh kamaal kar diya!"* (English translation: Did you watch the match yesterday? What a game! Virat was simply superb!)

For sentence 1, MF is 0.25 (BN: 3, EN: 9) while for sentence 2, MF is 0.5 (BN: 3, EN: 3, HI: 6)

MF can range from $1 - \frac{1}{W}$ (when every word in the sentence belongs to a different language) to 1 (for monolingual texts).

### 4.4 Complexity Factor

Finally, we have combined all the three factors to formulate the Complexity Factor as,

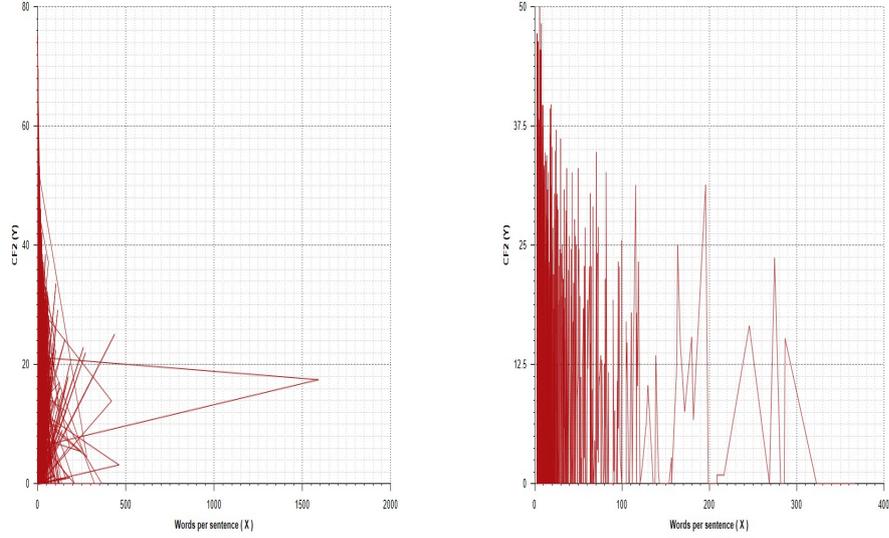$$CF = \frac{a * MF + b * SF}{f(LF)}$$

where a and b are the weights for Mix Factor (MF) and Switching Factor (SF), respectively. On the other hand, f is a function of Language Factor(LF) that we use as a dampening factor. After experimentation with the weights, we finally settled a = 50 and b = 50. We calculated CF by having f(LF) = LF, by using a linear function, $f(LF) = \left(\frac{0.25}{W-1}\right)(LF - 1) + 1$ and by using a geometric function, $f(LF) = \frac{arctan(LF)}{\pi} + 0.75$. We have calculated CF in three different ways and presented our results in Table 3 and Table 4.

$$CF1 = \frac{50 * MF + 50 * SF}{LF}$$

$$CF2 = \frac{50 * MF + 50 * SF}{\left(\frac{0.25}{W-1}\right)(LF-1)+1}$$

$$CF3 = \frac{50 * MF + 50 * SF}{\frac{arctan(LF)}{\pi} + 0.75}$$

We have considered MF and SF to be equally important while determining the code-mixing complexity of the social media text. However, the number of languages in social media texts is often limited to two or three. So, the impact of the LF on complexity has been dampened by the use of a linear function (in CF2) and a geometric function (in CF3). This ensures that the complexity of any given text is not heavily reduced by the language factor.

(a) FIRE Corpus: Graph of Words per Sentence vs. CF2.

(b) ICON Corpus: Graph of Words per Sentence vs. CF2.

Fig. 1: Performance of Complexity Factor for both corpora

## 5   Working of the Index

We have presented a few test cases to compare the performance of our index in comparison to the existing index (CMI). The first four cases presented are from a purely mathematical perspective and serves the purpose of illustrating the mathematical precision of the index which we proposed here. For the remaining cases, we have presented examples from the dataset to elucidate our objectives. In all the examples, $w_i$ and $l_i$ represents the word and the language at position i, respectively.

*Case 1: w1/l1 w2/l2 w3/l3 w4/l4 w5/l5 w6/l6 w7/l7 w8/l8 w9/l9 w10/l10*
The sentence contains 10 words each belonging to a different language. Ideally, any index should denote the complexity of such a code-mixed sentence as 100 (in a scale of 0-100). There is no better example of a more complex sentence than this from a multilingual perspective. CMI gives a value of 90 for such a sentence. CF2 and CF3 both give the complexity as 95.

*Case 2: w1/l1 w2/l1 w3/l1 w4/l1 w5/l1 w6/l1 w7/l1 w8/l1 w9/l11 w10/l1*
The sentence contains 10 words each belonging to the same language. Here, we would expect the index to show complexity as zero. Each of the three indexes, CMI, CF2 and

CF3, gives the complexity as 0. In case of CF2 and CF3, two of the three components - the language factor, switching factor and the mix factor all are zero.

*Case 3: w1/l1 w2/l2 w3/l1 w4/l2 w5/l1 w6/l2 w7/l1 w8/l2 w9/l11 w10/l2*
The sentence contains 10 words belonging to two languages. The words are arranged such that no two consecutive words belong to the same language. CMI calculates the complexity of the sentence to be 50. The complexities, as given by CF2 and CF3 are 67.5 and 63.2, respectively (LF=5, MF=0.5, SF=1).

*Case 4: w1/l1 w2/l1 w3/l1 w4/l1 w5/l1 w6/l2 w7/l2 w8/l2 w9/l12 w10/l2*
The sentence contains 10 words belonging to two languages. The words are arranged such that first five words belong to the one language and the next five words belong to a second language. Once again, CMI calculates the complexity of the index as 50. CF2 and CF3 calculates the complexity to be 27.52 and 25.73, respectively (LF=5, MF=0.5, SF=0.11). It has to be mentioned that Complexity Factor correctly estimates the sentence to be less complex than the previous example (which contains more switching).

The previous examples were theoretical to mainly highlight the mathematical background of the model. We have collected a few examples from our corpus to further illustrate the robustness of our index. Once again, we compared it with respect to CMI.

*Case 5: Koi ni bhai , apne dbc wale hosla ni haarte ... "think to score goals instead of thinking abt goalkeepers"*
The above sentence contains 9 English and 8 Hindi words. The value of CMI is 47 while CF2 and CF3 are 23.21 and 21.31, respectively. Complexity Factor Indexes are less than the CMI because there is only one language switch in the sentence.

*Case 6: mari bike ma puncture padayu*
In the above sentence, there are 2 English and 3 Gujarati words with 4 language switches (each alternate word belongs to a different language). Complexity Factor for the sentence is 64.22 (as in CF2) and 61.75 (as in CF3) with the highest possible Code-Switch Factor (which is 1). CMI gives a value of 40 because it considers only the fraction of non-dominant languages present.

*Case 7: Mi maza maharastra prem dhakvla .. tu swapnil joshi la hate karun jar saharukla support karanar asel tar saalam malakun ........... I like swapnil because he's maharastrian ... also I have never unbend opinion about you ........*
In the above sentence, there are 15 English and 15 Marathi words with 5 language switches. Although the length of the sentence is quite long, it has few language switches. CF2 and CF3 recognize that aspect and assign a complexity of 28.57 and 26.02, respectively while CMI assigns it a complexity of 50.

*Case 8: Steve : 10 th anniversary celebarate pannama poiduvomo - nu .*
In this case, we have selected a smaller sentence with 3 English and 3 Tamil words. There is 1 language switch present. Once again, CMI assigns it a complexity of 50.

CF2 and CF3 are 22.97 and 24.79 respectively. These values are less than that of Case 7 because of fewer number of language switches.

*Case 9: BIG B sings the eternal journey of life well .......... " tu shola ban jo khud jalke janha rashan karde ... ekla jalo re "*
This sentence contains words from 3 different languages - English (9 words), Bengali (3 words) and Hindi (9 words). The high proportion of language mixing makes CMI value 57. However, the words of all the languages occur in clusters with only 2 language switchings. Therefore, the CF2 and CF3 values are 30.09 and 26.86, respectively (as it considers the relative ordering of language words along with the presence of non-dominant language).

*Case 10: Happy Rakshabandhan(Rakhi ) ...... Piyali Kar Lipika Bisht Lopamudra Sarkar Mandira Agrawal Payel Ghosh Trishona Vanhi*
This is another example which contains only two words which are language specific (1 English and 1 Bengali word). The remaining words are named entities. CMI assigns it a complexity of 50. CF2 and CF3 values are 25.45 and 23.55 respectively.

*Case 11: r february te amar breakup hoy .*
This sentence contains 2 English and 4 Bengali words with 4 language switches. CMI value is 33 while CF2 and CF3 values are 45.45 and 43.1 respectively. The frequent switching of languages makes this sentence more complex than usual and Complexity Factor captures that aspect.
In the following section, we discuss the range of all the indexes in both the corpora.

## 6    Results on Different Corpora

We calculated the complexity of the FIRE and ICON corpora. The results are presented in Table 3 and Table 4. The minimum and maximum values show the range of the indexes for both the corpora. It may be noted that CF2 and CF3 shows a broader range than CMI in case of FIRE corpus. The primary reason for this is the presence of more languages in the FIRE corpus. Words per sentence has been used to highlight the length of sentences present in both the corpora. The average value of the index reflects the complexity of the entire corpus. In Figure 1, graphs have been plotted where X axis represents the length of the sentence and Y axis represents the index (CF2).

### 6.1   The FIRE 2015 Shared Task Corpus

### 6.2   The ICON 2015 Shared Task Corpus

The results suggest that the FIRE Corpus was more complex than the ICON Corpus with average value of CF2 and CF3 over the entire corpus being 10.54 and 9.88 respectively. For ICON corpus, CF2 and CF3 are 4.83 and 4.71 respectively which is considerably less than that of FIRE Corpus.

| Index | Minimum Value | Maximum Value | Average |
|---|---|---|---|
| CMI | 0 | 50 | 11.65 |
| CF1 | 0 | 75 | 2.51 |
| CF2 | 0 | 75 | 10.54 |
| CF3 | 0 | 75 | 9.88 |
| Words/sentence | 1 | 1592 | 17.16 |

Table 3: Range and Mean of Each Index and Words per Sentence (In FIRE Corpus).

| Index | Minimum Value | Maximum Value | Average |
|---|---|---|---|
| CMI | 0 | 57 | 5.73 |
| CF1 | 0 | 33.5 | 1.02 |
| CF2 | 0 | 50 | 4.83 |
| CF3 | 0 | 47.38 | 4.51 |
| Words/sentence | 1 | 367 | 16.32 |

Table 4: Range and Mean of Each Index and Words per Sentence (In ICON Corpus).

## 7    Conclusion and Future Work

In this paper, we have discussed the need and application of various indexes to represent the complexity of code-mixing in transliterated social media text. We have used various examples - both mathematical and from real-life text - to demonstrate the working of Code Mixing Index. We have highlighted few of the challenges that this index face and have proposed a new index - Complexity Factor (with two variations CF1 and CF2) - which takes into account the relative ordering of words (or the number of language switches) and the number of languages present in addition to the presence of words from non-dominant languages (as done in CMI). We have proposed three variations of Complexity Factor - CF1 serves as a baseline or raw index. CF2 and CF3 are more versatile and usable. CF2 uses linear interpolation while CF3 uses geometric functions. Both of these indexes provide a more balanced view of the complexity of the text. In future, the working of the index can be further explored using more data from multilingual texts (containing more than two languages). We can also find and compare the complexities of various corpora prior performing the tasks like part-of-speech tagging or sentiment analysis. The index can also be checked for mixed texts from other regions (like Spanish-English, Mandarin-English, etc.). In future, another challenging work would be to modify the index to account for complexity caused due to intra-word mixing.

## References

1. Auer, P.: Bilingual conversation. John Benjamins Publishing (1984)
2. Gambäck, B., Das, A.: On measuring the complexity of code-mixing. In: Proceedings of the 11th International Conference on Natural Language Processing, Goa, India. pp. 1–7 (2014)
3. Kilgarriff, A.: Comparing corpora. International journal of corpus linguistics 6(1), 97–133 (2001)