

# PhonSenticNet: A Cognitive Approach to Microtext Normalization for Concept-Level Sentiment Analysis

Ranjan Satapathy<sup>1</sup>, Aalind Singh<sup>2</sup>, Erik Cambria<sup>1</sup>

<sup>1</sup>SCSE, Nanyang Technological University

<sup>2</sup>Vellore Institute of Technology, India

satapathy.ranjan@ntu.edu.sg, singh.aalind@gmail.com,  
cambria@ntu.edu.sg

## Abstract

With the current upsurge in the usage of social media platforms, the trend of using short text (microtext) in place of standard words has seen a significant rise. The usage of microtext poses a considerable performance issue in concept-level sentiment analysis, since models are trained on standard words. This paper discusses the impact of coupling sub-symbolic (phonetics) with symbolic (machine learning) Artificial Intelligence to transform the out-of-vocabulary concepts into their standard in-vocabulary form. The phonetic distance is calculated using the Sorensen similarity algorithm. The phonetically similar in-vocabulary concepts thus obtained are then used to compute the correct polarity value, which was previously being miscalculated because of the presence of microtext. Our proposed framework increases the accuracy of polarity detection by 6% as compared to the earlier model. This also validates the fact that microtext normalization is a necessary pre-requisite for the sentiment analysis task.

**Index Terms:** microtext normalization, phonetics, concept level sentiment analysis

## 1. Introduction

With the popularization of mobile phones and Internet social networks, the use of electronic text messaging, or texting, has reached astonishing figures such as more than 8,000 tweets produced per second<sup>1</sup>. These type of communications are usually performed in real time and over platforms which impose limitations on the length of the messages, as in the case of Twitter or the traditional SMS system. Because of this, the writing style of these messages differs from normal standards and phenomena such as word shortenings, contractions and abbreviations are commonly used both to gain writing speed and circumvent length limitations. Moreover, even in the case of messaging platforms where length restrictions do not apply (e.g. WhatsApp), it is also common to see a writing style which tries to better reflect the feelings of the writer. Given that most data today is mined from the web, microtext analysis is vital for many natural language processing (NLP) tasks. In the context of sentiment analysis, microtext normalization is a necessary step for pre-processing text before polarity detection is performed [1].

The two main features of microtext are relaxed spelling and reliance on emoticons and out-of-vocabulary (OOV) words involving phonetic substitutions (e.g., 'b4' for 'before'), emotional emphasis (e.g., 'goooooood' for 'good') and popular acronyms (e.g., 'otw' for 'on the way') [2, 3, 4]. The challenge arises when trying to automatically rectify and replace them with the correct in-vocabulary (IV) words [5]. It could be

thought that microtext normalization is as simple as performing find-and-replace pre-processing [6]. However, the wide-ranging diversity of spellings makes this solution impractical (e.g., the spelling of the word "tomorrow" is generally written as "tomorow, 2moro, tmr among others). Furthermore, given the productivity of users, novel forms which are not bound to orthographic norms in spelling can emerge. For instance, a sampling of Twitter studied in [5] found over 4 million OOV words where new spellings were created constantly, both voluntarily and accidentally. Concept-based approaches to sentiment analysis focus on a semantic analysis of text through the use of web ontologies or semantic networks, which allow the aggregation of conceptual and affective information associated with natural language opinions. The analysis at concept-level is intended to infer the semantic and affective information associated with natural language opinions and hence, to enable a comparative fine-grained feature based sentiment analysis. In this work, we propose PhonSenticNet, a concept based lexicon which advantages from phonetic features to normalize the OOV concepts to IV concepts. The International Phonetic Alphabet (IPA)<sup>2</sup> is used as the phonetic feature in the proposed framework.

The rest of the paper is organised as follows: Section 2 discusses the literature survey in microtext, Section 3 discusses the datasets used, Section 4 discusses the experiments performed and Section 5 concludes the work done with future directions for this work.

## 2. Related Work

This section discusses the work done in the domain of microtext analysis.

### 2.1. Microtext Analysis

Microtext has become ubiquitous in today's communication. This is partly a consequence of Zipf's law, or principle of least effort (for which people tend to minimize energy cost at both individual and collective levels when communicating with one another), and it poses new challenges for NLP tools which are usually designed for well-written text [7]. Normalization is the task of transforming unconventional words or concepts to their respective standard counterpart. [8] uses Soundex algorithm to transform out-of-vocabulary to in-vocabulary and shows its effect on sentiment analysis task.

In [9], authors present a novel unsupervised method to translate Chinese abbreviations. It automatically extracts the relation between a full-form phrase and its abbreviation from monolingual corpora, and induces translation entries for the ab-

<sup>1</sup><http://www.internetlivestats.com/one-second/>

<sup>2</sup><https://www.internationalphoneticassociation.org/content/full-ipa-chart>

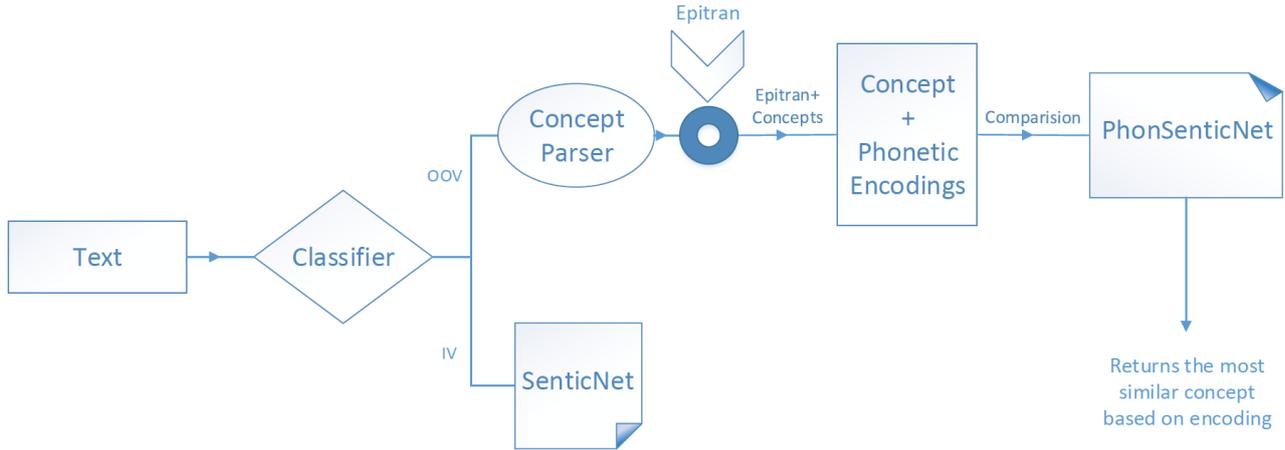


Figure 1: Architecture of the framework

abbreviation by using its full-form as a bridge. [10] uses a classifier to detect OOV words, and generates correction candidates based on morpho-phonemic similarity. The types and features of microtext are reliant on the nature of the technological support that makes them possible. This means that microtext will vary as new communication technologies emerge. In our related work, we categorized normalization into three well-known NLP tasks, namely: spelling correction, statistical machine translation (SMT), and automatic speech recognition (ASR).

### 2.1.1. Spelling Correction

Correction is executed on a word-per-word basis which is also seen as a spell checking task. This model gained extensive attention in the past and a diversity of correction practices have been endorsed by [11, 12, 13, 14, 15]. Instead, [16] and [17] proposed a categorization of abbreviation, stylistic variation, prefix-clipping, which was then used to estimate the probability of occurrence of the characters. Thus far, the spell corrector became widely popular in the context of SMS, where [18] advanced the hidden Markov model whose topology takes into account both “graphemic” variants (e.g., typos, omissions of repeated letters, etc.) and “phonemic” variants. All of the above, however, only focused on the normalization of words without considering their context.

### 2.1.2. Statistical Machine Translation

When compared to the previous task, this method appears to be rather straightforward and better since it has the possibility to model (context-dependent) one-to-many relationships which were out-of-reach previously [19]. Some examples of works include [20, 21, 22]. However, the SMT still overlooks some features of the task, particularly the fact that lexical creativity verified in social media messages is barely captured in a stationary sentence board.

### 2.1.3. Automatic Speech Recognition

ASR considers that microtext tends to be a closer approximation of the word’s phonemic representation rather than its standard spelling. As follows, the key to microtext normalization becomes very similar to speech recognition which consists of decoding a word sequence in a (weighted) phonetic framework.

For example, [19] proposed to handle normalization based on the observation that text messages present a lot of phonetic spellings, while more recently [6] proposed an algorithm to determine the probable pronunciation of English words based on their spelling. Although the computation of a phonemic representation of the message is extremely valuable, it does not solve entirely all the microtext normalization challenges (e.g., acronyms and misspellings do not resemble their respective IV words’ phonemic representation). Authors in [23] have merged the advantages of SMT and the spelling corrector model.

## 3. Datasets

This section introduces to datasets used. The twitter dataset is available on request. The concept-level lexicon SenticNet is publically available<sup>3</sup>.

### 3.1. NUS SMS Corpus

This corpus has been created from the NUS English SMS corpus<sup>4</sup>, wherein [24] randomly selected 2,000 messages. The messages were first normalized into standard English and then translated into standard Chinese. For our training and testing purposes, we only used the actual messages and their normalized English version. It also contains non-English terms, which LSTM had no problem in learning. Singlish is an English-based creole that is lexically and syntactically influenced by Hokkien, Cantonese, Mandarin, Malay and Tamil [25]. It is primarily a spoken variety in Singapore, to emerge as a means of online communication for Singaporeans [26].

### 3.2. SenticNet

SenticNet [27] is a knowledge base of 100,000 commonsense concepts. Sentic API provides the semantics and sentics (i.e., the denotative and connotative information) associated with the concepts of SenticNet 5, a semantic network of commonsense knowledge that contains 100,000 nodes (words and multiword expressions) and thousands of connections (relationships between nodes). We used concept parser [28] in order to break sentences to concepts and analyze them. The concepts in the SenticNet contains their corresponding polarities.

<sup>3</sup><http://sentic.net/senticnet-5.0.zip>

<sup>4</sup><http://github.com/kite1988/nus-sms-corpus>

Concepts	Polarity	Soundex Encoding	IPA Encoding
a_little	Negative	A000.L340	æ..litəl
abandon	Negative	A153	æbɒndæn
absolutely_fantastic	Positive	A124.F532	əbsəlu:tli.fəntəstɪk

Table 1: Sample Soundex and IPA Encodings with polarities for SenticNet5

### 3.3. Normalized Tweets

The authors in [8], built a dataset which consists of tweets and their transformed in-vocabulary counterparts. We demonstrate our results by extracting concepts from unconventionally written sentences and then passing them through our proposed module to convert them to standard format concepts<sup>7</sup> and their corresponding polarities from SenticNet.

## 4. Experiments

This section dives into experiments performed to develop a concept-level microtext normalization module. The experiments performed help in deciding the best set of parameters to achieve state of the art accuracy. We name the lexicon which we built from SenticNet as PhonSenticNet. It contains concepts and their related phonetic encoding which is extracted from Epi-tran [29].

### 4.1. Framework

The architecture of the proposed model is depicted in Figure 1. The framework classifies a sentence as OOV or IV using binary classifier. Following this, the OOV sentence is passed through the concept parser, and then the concepts are transformed to IPA using Epi-tran. The IPA of OOV concepts are matched to the PhonSenticNet, and then the IV concept is fetched. The corresponding polarity of the IV concept is retrieved from SenticNet. The detailed procedure is explained in the following subsections:

#### 4.1.1. Classification of microtext

In this subsection, we employ various binary classifiers to detect microtext so as to reduce the execution time of the overall algorithm. We observed that the execution time of polarity detection task was reduced by 20%. Different classifiers were trained on the two datasets namely NUS SMS data and Twitter dataset as shown in Table 2.

We use the term frequencyinverse document frequency (TF-IDF) [30] approach for the task of feature extraction from a given text. We first split the document into tokens assigning weights to them based on the frequency with which it shows up in the document along with how recurrent that term occurs in the entire corpora. We used this approach to train four different classifiers. The evaluation metrics such as Precision, Recall, F-measure and Accuracy have been enlisted in the Table 2.

#### 4.1.2. Soundex vs IPA

We compare our proposed IPA based method to Soundex [8] since only [8] have incorporated phonetic features to improve sentiment analysis. Soundex gives a lot of duplicate encoding, whereas IPA gives no duplicate encoding. Each concept is unique in phonetic subspace, thereby increasing the efficiency for microtext normalization at concept-level. The number of

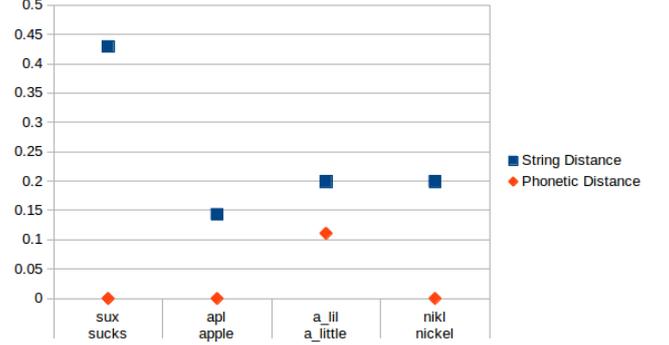


Figure 2: Visualization of string and phonetic distance

concepts present in the lexicon is 100000. The duplicates<sup>5</sup> due to Soundex encoding are 46080. This shows that using soundex for microtext normalization has some information loss at the concept-level. As a result of Soundex encoding, we have 46080 ambiguous concepts which affect the microtext normalization in real time. Hence, we propose to use IPA for all the phonetic based microtext normalization methods. The IPA based encoding has no redundancy, and thereby no information loss occurs during microtext normalization.

---

**Algorithm 1** Algorithm for microtext normalization using phonetic features

---

```

Sentence (S) = s1, s2, ..., sn
ci = concept-parser(S)
For each concept ci in Sn
closest-match-concept = PhonSenticNet(ci)
if closest_match(Ci, SenticNet) then
return concept polarity
else
return sentence polarity
end if
average over polarity of concepts for sentence polarity
EndFor
return sentence polarity

```

---



---

**Algorithm 2** Closest Match Algorithm

---

```

Concept (C) = c1, c2, ..., cm
For each concept ci in C
Sorensen (ci, Senticnet)
EndFor
return phonetically closest matching concept

```

---

#### 4.1.3. Microtext normalization of concepts using IPA

The analysis at concept-level is intended to infer the semantic and affective information associated with natural language opinions and, hence, to enable a comparative fine-grained feature-based sentiment analysis. Concept-based approaches to sentiment analysis focus on a semantic analysis of text through the use of web ontologies or semantic networks, which allow the aggregation of conceptual and affective information associated

<sup>5</sup>Repetition of a soundex encoding for greater than one

Table 2: Precision, Recall, F1 and Accuracy for each algorithm on different datasets

	NUS SMS Dataset								Twitter Dataset							
	Logistic - Regression		SGDC		SVC		Multinomial-NB		Logistic - Regression		SGDC		SVC		Multinomial-NB	
	IV	OOV	IV	OOV	IV	OOV	IV	OOV	IV	OOV	IV	OOV	IV	OOV	IV	OOV
Precision	0.91	0.95	0.84	0.98	0.87	0.97	0.89	0.97	0.71	0.69	0.63	0.72	0.74	0.67	0.81	0.68
Recall	0.95	0.90	0.99	0.81	0.98	0.85	0.97	0.87	0.68	0.71	0.80	0.52	0.64	0.77	0.61	0.85
F-measure	0.93	0.92	0.91	0.89	0.92	0.91	0.93	0.92	0.69	0.70	0.70	0.60	0.68	0.72	0.69	0.76
Accuracy	<b>0.9275</b>		0.89875		0.915		0.9225		0.6962		0.6605		0.7013		<b>0.7288</b>	

with natural language opinions. In order to normalize concepts found on the social media, we built a resource for concept-level phonetic encodings by using concepts from SenticNet. We used concept parser [28] to extract concepts from the input text. The concepts are then transformed to a subspace where they are represented by their phonetic encodings. Table 1 shows sample concepts with their respective soundex encodings and IPA from SenticNet 5.

Phonetic encoding transforms concepts from the string subspace into their phonetic subspace. This transformation eliminates the redundant concept encoding produced by Soundex. The input concept is then passed through this phonetic subspace (IPA used in PhonSenticnet) to find the most phonetically similar concept and then returns it. The algorithm 1 and 2 describe the procedures in detail.

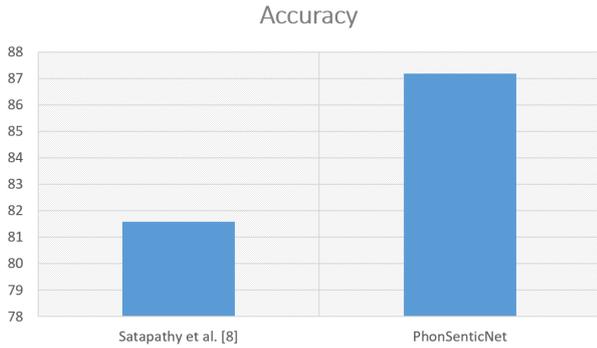


Figure 3: Accuracy for polarity detection

#### 4.1.4. Polarity Detection with SenticNet

While SenticNet 5 can be used as any other sentiment lexicon, e.g., concept matching or bag-of-concepts model, the right way to use the knowledge base for the task of polarity detection is in conjunction with sentic patterns, sentiment-specific linguistic patterns that infer polarity by allowing affective information to flow from concept to concept based on the dependency relation between clauses. The sentiment sources of such affective information are extracted from SenticNet 5 by firstly generalizing multiword expressions and words by means of conceptual primitives and, secondly, by extracting their polarity. We compare our polarity detection module results with [8]. The accuracy increases significantly by 6% as shown in Figure 3.

## 5. Discussion and Future Work

The proposed resource contains concepts from SenticNet and their phonetics by using Epitran which we name as PhonSenticNet. This resource is used as a lexicon for microtext normalization. The input sentence is broken down into concepts and

Text	Sentence Polarity before microtext normalization	Sentence Polarity after microtext normalization
I wil kil u	Neutral	Negative
m so hapy	Neutral	Positive
i dnt lyk reading	Positive	Negative
it is awesum 2 ride byk	Neutral	Positive

Table 3: Sample sentences before and after microtext normalization

then transformed into their phonetic encoding. The phonetic encoding is matched with the PhonSenticNet, the resource built in this work. Then the most similar matching concept and its corresponding polarity is returned as shown in the algorithm 1 and 2.

1. We have taken Sorensen similarity to measure the distance. The Sorensen similarity shows how similar the two input texts are to one another, where 0 means similar and 1 means dissimilar as shown in Figure 2.
2. Figure 2 shows some of the distance metric between non-standard and their standard concepts. The similarity is shown at both string and phonetic level.
3. Previous paper [8] shows sentence-level sentiment analysis, whereas in this work we focus on concept-level microtext normalization.
4. It can be observed from Table 2 that the twitter dataset does not perform as good as the NUS SMS data. The reason behind it is, the twitter dataset contains acronyms like lol, rofl, etc instead of phonetic substitution. This also suggests how the way of writing differs in both messages and tweets.

Microtext is very much language-dependent: the same set of characters could have completely different meaning in different languages, e.g., ‘555’ is negative in Chinese language because the number ‘5’ is pronounced as ‘wu’ and ‘wuwuwu’ resembles a crying sound but positive in Thai since the number 5 is pronounced as ‘ha’ and three consecutive 5s correspond to the expression ‘hahaha’. Hence, we are working on its multilingual version [31]. The proposed work only works for the phonetic class of microtext analysis. Though, the acronyms still rely on the lexicon built in [8].

## 6. References

- [1] E. Cambria, S. Poria, A. Gelbukh, and M. Thelwall, “Sentiment analysis is a big suitcase,” *IEEE Intelligent Systems*, vol. 32, no. 6, pp. 74–80, 2017.
- [2] J. Read, “Using emoticons to reduce dependency in machine learning techniques for sentiment classification,” in *Proceedings of the ACL student research workshop*. Association for Computational Linguistics, 2005, pp. 43–48.
- [3] K. D. Rosa and J. Ellen, “Text classification methodologies applied to micro-text in military chat,” in *Proc. Eight International*

- Conference on Machine Learning and Applications, Miami, 2009, pp. 710–714.
- [4] Z. Xue, D. Yin, and B. D. Davison, “Normalizing Microtext,” *Analyzing Microtext*, pp. 74–79, 2011.
  - [5] F. Liu, F. Weng, B. Wang, and Y. Liu, “Insertion, deletion, or substitution? Normalizing text messages without pre-categorization nor supervision,” *ACL-HLT 2011 - Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, vol. 2, pp. 71–76, 2011.
  - [6] R. Khoury, “Microtext Normalization using Probably-Phonetically-Similar Word Discovery,” in *Wireless and Mobile Computing, Networking and Communications (WiMob), 2015 IEEE 11th International Conference on.*, 2015, pp. 392–399.
  - [7] C. J. Hutto and E. Gilbert, “VADER: A parsimonious rule-based model for sentiment analysis of social media text,” in *Eighth International AAAI Conference on Weblogs and Social Media*, 2014, pp. 216–225.
  - [8] R. Satapathy, C. Guerreiro, I. Chaturvedi, and E. Cambria, “Phonetic-based microtext normalization for twitter sentiment analysis,” in *2017 IEEE International Conference on Data Mining Workshops (ICDMW)*. IEEE, 2017, pp. 407–413.
  - [9] Z. Li and D. Yarowsky, “Unsupervised translation induction for chinese abbreviations using monolingual corpora,” in *In Proceedings of ACL/HLT*, 2008.
  - [10] B. Han and T. Baldwin, “Lexical normalisation of short text messages: Mkn sens a# twitter,” in *ACL*, 2011, pp. 368–378.
  - [11] K. W. Church and W. A. Gale, “Probability scoring for spelling correction,” *Statistics and Computing*, vol. 1, no. 2, pp. 93–103, 1991.
  - [12] E. Brill and R. C. Moore, “An improved error model for noisy channel spelling correction,” in *Proceedings of the 38th Annual Meeting on Association for Computational Linguistics*, 2000, pp. 286–293.
  - [13] M. Li, Y. Zhang, M. Zhu, and M. Zhou, “Exploring distributional similarity based models for query spelling correction,” in *ACL*, 2006, pp. 1025–1032.
  - [14] D. L. Pennell and Y. Liu, “A Character-Level Machine Translation Approach for Normalization of SMS Abbreviations,” in *IJCNLP*, 2011, pp. 974–982.
  - [15] K. Toutanova and R. C. Moore, “Pronunciation modeling for improved spelling correction,” in *ACL*, 2002, pp. 144–151.
  - [16] R. Sproat, A. W. Black, S. Chen, S. Kumar, M. Ostendorf, and C. Richards, “Normalization of non-standard words,” *Computer speech & language*, vol. 15, no. 3, pp. 287–333, 2001.
  - [17] P. Cook and S. Stevenson, “An unsupervised model for text message normalization,” in *Proceedings of the workshop on computational approaches to linguistic creativity*, 2009, pp. 71–78.
  - [18] M. Choudhury, R. Saraf, V. Jain, S. Sarkar, and A. Basu, “Investigation and modeling of the structure of texting language,” *International Journal of Document Analysis and Recognition*, vol. 10, no. 3-4, pp. 157–174, 2007.
  - [19] C. Kobus, F. Yvon, and G. é. Damnati, “Normalizing SMS: are two metaphors better than one?” in *Proceedings of the 22nd International Conference on Computational Linguistics*, vol. 1. Association for Computational Linguistics, 2008, pp. 441–448.
  - [20] A. Aw, M. Zhang, J. Xiao, and J. Su, “A phrase-based statistical model for SMS text normalization,” in *ACL*, 2006, pp. 33–40.
  - [21] M. Kaufmann and J. Kalita, “Syntactic normalization of Twitter messages,” *natural language processing, Kharagpur, India*, 2010.
  - [22] D. L. Pennell and Y. Liu, “Normalization of informal text,” *Computer Speech & Language*, vol. 28, no. 1, pp. 256–277, 2014.
  - [23] R. Beaufort, S. Roekhaut, L.-A. I. Cougnon, and C. d. Fairon, “A hybrid rule/model-based finite-state framework for normalizing SMS messages,” in *ACL*. Association for Computational Linguistics, 2010, pp. 770–779.
  - [24] P. Wang and H. T. Ng, “A beam-search decoder for normalization of social media text with application to machine translation,” in *HLT-NAACL*, 2013, pp. 471–481.
  - [25] A. Brown, *Singapore English in a nutshell: An alphabetical description of its features*. Federal Publications, 1999.
  - [26] M. Warschauer, “The internet and linguistic pluralism,” *Silicon literacies: Communication, innovation and education in the electronic age*, pp. 62–74, 2002.
  - [27] E. Cambria, S. Poria, D. Hazarika, and K. Kwok, “SenticNet 5: Discovering conceptual primitives for sentiment analysis by means of context embeddings,” in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018, pp. 1795–1802.
  - [28] S. Poria, B. Agarwal, A. Gelbukh, A. Hussain, and N. Howard, “Dependency-based semantic parsing for concept-level text analysis,” in *International Conference on Intelligent Text Processing and Computational Linguistics*. Springer, 2014, pp. 113–127.
  - [29] D. R. Mortensen, S. Dalmia, and P. Littell, “Epitran: Precision G2P for Many Languages,” in *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*. Paris, France: European Language Resources Association (ELRA), May 2018, pp. 7–12.
  - [30] J. Ramos *et al.*, “Using tf-idf to determine word relevance in document queries,” in *Proceedings of the first instructional conference on machine learning*, vol. 242, 2003, pp. 133–142.
  - [31] D. Vilares, H. Peng, R. Satapathy, and E. Cambria, “Babelsenticnet: A commonsense reasoning framework for multilingual sentiment analysis,” in *2018 IEEE Symposium Series on Computational Intelligence (SSCI)*. IEEE, 2018, pp. 1292–1298.