



# **Tweets Segmentation based on Popularity of Posted Tweets with Help of Hashtag**

Praveen M<sup>1</sup>, SK. Prasanth<sup>2</sup>

M.Tech Student, Dept. of CSE., VCE, JNTUH, Hyderabad, Telangana, India<sup>1</sup>

Associate Professor, Dept. of CSE., VCE, JNTUH, Hyderabad, Telangana, India<sup>2</sup>

**ABSTRACT:** Twitter has pulled in a huge number of clients to share and spread most state-of-the-art data, bringing about huge volumes of information created regular. Nonetheless, numerous applications in Information Retrieval (IR) and Natural Language Processing (NLP) endure extremely from the loud and short nature of tweets. In this paper, we propose a novel system for tweet division in a bunch mode, called HybridSeg. By part tweets into significant sections, the semantic or setting data is all around safeguarded and effortlessly removed by the downstream applications. HybridSeg finds the ideal division of a tweet by augmenting the whole of the stickiness scores of its hopeful fragments. The stickiness score considers the likelihood of a section being an expression in English (i.e., worldwide connection) and the likelihood of a section being an expression inside the clump of tweets (i.e., neighborhood setting). For the last mentioned, we propose what's more, assess two models to infer neighborhood setting by considering the semantic components and term-reliance in a cluster of tweets, separately. HybridSeg is likewise intended to iteratively gain from sure sections as pseudo input. Probes two tweet information sets demonstrate that tweet division quality is altogether enhanced by learning both worldwide and neighborhood settings contrasted and utilizing worldwide setting alone. Through investigation and examination, we demonstrate that nearby etymological components are more solid for learning neighborhood connection contrasted and term-reliance. As an application, we demonstrate that high precision is accomplished in named substance acknowledgment by applying section based grammatical feature (POS) labeling..

**KEYWORDS:** Twitter stream, tweet segmentation, named entity recognition, linguistic processing, Wikipedia

## **I. INTRODUCTION**

Mirco blogging destinations, for example, Twitter have reshaped the way individuals discover, share, and disperse auspicious data. Numerous associations have been accounted for to make what's more, screen focused on Twitter streams to gather and get it clients' conclusions. Focused on Twitter stream is for the most part built by separating tweets with predefined choice criteria (e.g., tweets distributed by clients from a geological locale, tweets that match one or more predefined watchwords). Because of its precious business estimation of opportune data from these tweets, it is basic to get it tweets' dialect for a substantial collection of downstream applications, for example, named element acknowledgment (NER) [1], [3], [4], occasion discovery and synopsis [5], [6], [7], supposition mining [8], [9], opinion examination [10], [11], and numerous others. Given the restricted length of a tweet (i.e., 140 characters) what's more, no confinements on its written work styles, tweets regularly contain syntactic blunders, incorrect spellings, and casual condensing. The mistake inclined and short nature of tweets frequently make the word-level dialect models for tweets less dependable. For instance, given a tweet "I call her, no answer. Her telephone clinched, she dancing," there is no sign to figure its genuine topic by slighting word request (i.e., pack of-word model).

The circumstance is further exacerbated with the constrained connection gave by the tweet. That is, more than one clarification for this tweet could be inferred by various perusers if the tweet is considered in confinement. Then again, in spite of the loud way of tweets, the center semantic data is all around safeguarded in tweets as named elements or semantic expressions. For instance, the rising phrase "she dancin" in the related tweets demonstrates that it is a key idea—it orders this tweet into the group of tweets discussing the melody "She Dancin", a pattern theme in Inlet Area in January 2013. In this paper, we concentrate on the assignment of tweet division. The objective of this assignment is to part a tweet into an arrangement of successive n-grams (n = 1p, each of which is known as a section. A section can be a

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 7, July 2016

named substance (e.g., a motion picture title "finding nemo"), a semantically significant data unit (e.g., "formally discharged"), or whatever other sorts of expressions which seem "more than by chance" [1].

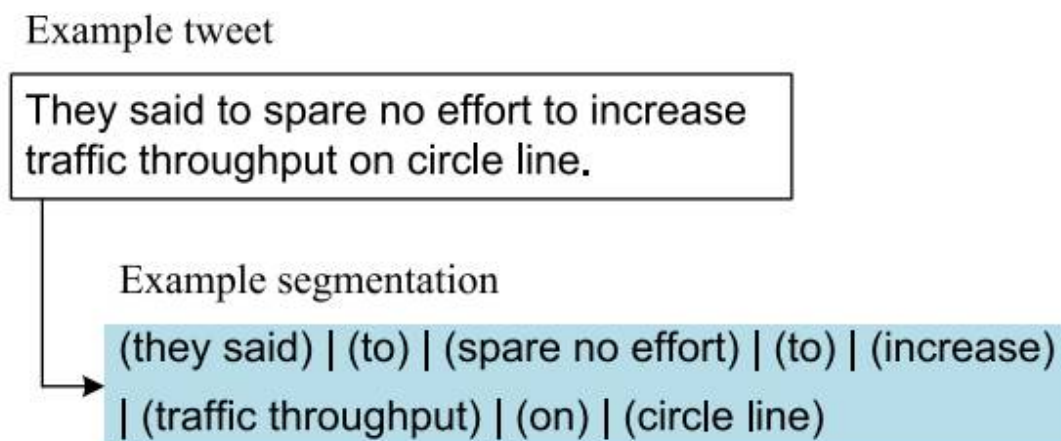


Fig. 1. Example of tweet segmentation.

Fig. 1 gives a case. In this case, a tweet "They said to save no push to increment activity throughput on circle line." is part into eight portions. Semantically significant sections "save no exertion", "activity throughput" and "circle line" are safeguarded. Since these portions safeguard semantic importance of the tweet more unequivocally than each of its constituent words does, the theme of this tweet can be better caught in the consequent preparing of this tweet. For example, this portion based representation could be utilized to improve the extraction of land area from tweets as a result of the portion "circle line" [12]. Truth be told, section based representation has demonstrated its viability over word-based representation in the errands of named element acknowledgment and occasion identification [1], [2], [13]. Note that, a named element is legitimate section; however a portion may not inexorably be a named element. In [6] the portion "korea versus greece" is distinguished for the occasion identified with the world container match amongst Korea and Greece.

To accomplish excellent tweet division, we propose a bland tweet division system, named HybridSeg. HybridSeg gains from both worldwide and nearby settings, and has the capacity of gaining from pseudo input. Worldwide connection. Tweets are posted for data sharing also, correspondence. The named substances and semantic expressions are all around protected in tweets. The worldwide setting gotten from Web pages (e.g., Microsoft Web N-Gram corpus) alternately Wikipedia subsequently helps distinguishing the important fragments in tweets. The technique understanding the proposed system that exclusively depends on worldwide setting is signified by HybridSegWeb. Neighborhood setting. Tweets are exceptionally time-delicate so that numerous developing expressions like "She Dancin" can't be found in outer learning bases. In any case, considering an expansive number of tweets distributed inside a brief span period (e. g., a day) containing the expression, it is not hard to perceive "She Dancin" as a substantial and significant portion. We along these lines examine two neighborhood connections, to be specific nearby etymological components and nearby collocation.

Watch that tweets from numerous official records of news offices, associations, what's more, promoters are likely elegantly composed. The very much saved etymological components in these tweets encourage named substance acknowledgment with high precision. Each named element is a substantial fragment. The technique using nearby etymological elements is signified by HybridSegNER. It acquires sure sections in view of the voting consequences of different off-the-rack NER instruments. Another technique using neighborhood collocation learning, meant by HybridSegNgram, is proposed based on the perception that numerous tweets distributed inside a brief day and age are



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 7, July 2016

about the same point. HybridSegNGram fragments tweets by evaluating the term-reliance inside a cluster of tweets. Pseudo criticism. The portions perceived taking into account neighborhood setting with high certainty serve as great input to separate more significant fragments. The gaining from pseudo criticism is led iteratively and the technique executing the iterative learning is named HybridSegIter. We direct broad trial investigation on HybridSeg1 on two tweet information sets and assess the nature of tweet division against physically commented on tweets. Our trial results demonstrate that HybridSegNER and HybridSegNGram, the two strategies fusing nearby connection in extra to worldwide connection, accomplish noteworthy change in division quality over HybridSegWeb, the strategy use worldwide connection alone.

Between the previous two techniques, HybridSegNER is less delicate to parameter settings than HybridSegNGram and accomplishes better division quality. With iterative gaining from pseudo criticism, HybridSegIter further enhances the division quality. As a use of tweet division, we propose what's more, assess two portion based NER calculations. Both calculations are unsupervised in nature and take tweet sections as information. One calculation misuses co-event of named substances in focused Twitter streams by applying arbitrary walk (RW) with the suspicion that named elements are more prone to co-happen together. The other calculation uses Grammatical feature (POS) labels of the constituent words in sections. The fragments that are prone to be a thing expression (NP) are considered as named substances. Our exploratory results demonstrate that (i) the nature of tweet division fundamentally influences the exactness of NER, and (ii) POS-based NER strategy beats RW-construct technique with respect to both information sets.

## II. RELATED WORK

Both tweet division and named substance acknowledgment are considered vital subtasks in NLP. Numerous current NLP systems vigorously depend on etymological elements, for example, POS labels of the encompassing words, word upper casing, trigger words (e.g., Mr., Dr.), and gazetteers. These semantic elements, together with compelling regulated learning calculations (e.g., hidden markov model (HMM)), accomplish great execution on formal content corpus [14], [15], [16]. In any case, these systems experience serious execution disintegration on tweets on account of the uproarious and short nature of the last mentioned. There has been a considerable measure of endeavors to join tweet's interesting qualities into the ordinary NLP systems. To enhance POS labeling on tweets, Ritter et al. train a POS tagger by utilizing CRF model with routine and tweet-particular elements [3].

Chestnut grouping is connected in their work to manage the poorly framed words. Gimple et al. consolidate tweet-particular elements including at-notice, hashtags, URLs, and feelings [17] with the assistance of another naming plan. In their methodology, they measure the confidence of uppercase words and apply phonetic standardization to badly shaped words to address conceivable curious compositions in tweets. It was accounted for to beat the state-of-the-craftsmanship Stanford POS tagger on tweets. Standardization of not well framed words in tweets has built up itself as an vital exploration issue [18]. A managed methodology is utilized in [18] to first recognize the badly shaped words. At that point, the right standardization of the badly framed word is chosen in light of various lexical similitude measures. Both administered and unsupervised methodologies have been proposed for named element acknowledgment in tweets.

T-NER, a part of the tweet-particular NLP structure in [3], first portions named elements utilizing a CRF model with orthographic, relevant, lexicon and tweet-particular elements. It then marks the named elements by applying Labeled-LDA with the outside learning base Freebase.2 The NER arrangement proposed in [4] is additionally taking into account a CRF model. It is a twostage expectation accumulation model. In the principal arrange, a KNN-based classifier is utilized to direct word-level classification, utilizing the comparative and as of late named tweets. In the second stage, those forecasts, alongside other phonetic components, are bolstered into a CRF model for better grained characterization. Chua et al. [19] propose to concentrate thing phrases from tweets utilizing an unsupervised methodology which is for the most part in light of POS labeling. Each separated thing phrase is a competitor named substance.

Our work is likewise identified with substance connecting (EL). EL is to distinguish the notice of a named substance and connection it to a passage in a learning base like Wikipedia [20], [21], [22], [23]. Traditionally, EL includes a NER framework took after by a connecting framework [20], [21]. As of late, Sil and Yates propose to join named substance

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 7, July 2016

acknowledgment and connecting into a joint model [23]. Correspondingly, Guo et al. propose an auxiliary SVM answer for all the while perceive specify and resolve the connecting [22]. While substance connecting plans to recognize the limit of a named substance and purpose its significance in light of an outer learning base, a commonplace NER framework recognizes substance notice just, similar to the work exhibited here. It is hard to make a reasonable examination between these two systems. Tweet division is adroitly like Chinese word division (CSW). Content in Chinese is a persistent arrangement of characters. Dividing the arrangement into important words is the initial phase in many applications.

Best in class CSW techniques are for the most part created utilizing administered learning methods like perceptron learning what's more, CRF model [24], [25], [26], [27], [28]. Both semantic what's more, vocabulary elements are utilized as a part of the administered learning in CSW. Tweets are to a great degree uproarious with incorrect spellings, casual shortened forms, and syntactic blunders. These unfriendly properties lead to countless specimens for applying a managed learning procedure. Here, we misuse the semantic data of outside learning bases and nearby settings to perceive significant fragments like named elements and semantic expressions in Tweets. Extremely as of late, comparable thought has likewise been investigated for CSW by Jiang et al. [28]. They propose to prune the inquiry space in CSW by misusing the regular comments in the Web. Their exploratory results show huge change by utilizing straightforward neighborhood highlights.

### III. PROPOSED ALGORITHM

The proposed HybridSeg system fragments tweets in cluster mode. Tweets from a focused on Twitter stream are assembled into clusters by their distribution time utilizing a settled time interim (e.g., a day). Every cluster of tweets are then divided by HybridSeg by and large. Given a tweet  $t$  from cluster  $T$ , the issue of tweet division is to part the " words in  $t$   $w_1 w_2 \dots w'$  into  $m$

" back to back sections,  $t = s_1 s_2 \dots s_m$ , where every fragment  $s_i$  contains one or more words. We define the tweet division issue as an enhancement issue to expand the aggregate of stickiness scores of the  $m$  sections, appeared in Fig. 2.3 A high stickiness score of portion  $s$  demonstrates that it is an expression which seems "more than by chance", and further part it could break the right word collocation or the semantic importance of the expression.

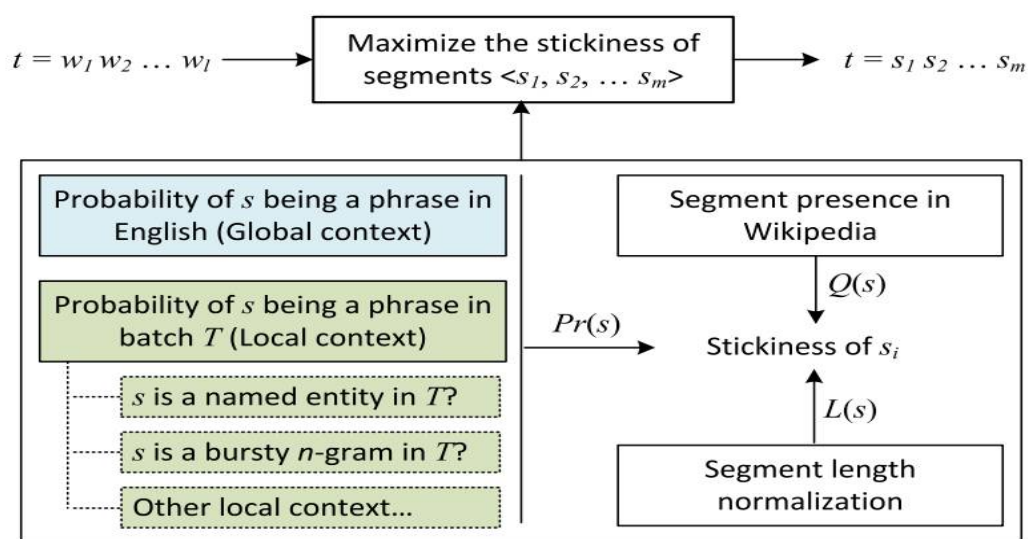


Fig. 2. HybridSeg framework without learning from pseudo feedback.



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 7, July 2016

Formally, let  $C_{s,p}$  mean the stickiness capacity of fragment  $s$ . The ideal division can be determined by utilizing dynamic programming with a period multifaceted nature (rf. [1] for point of interest). As appeared in Fig. 2, the stickiness capacity of a section takes in three variables: (i) length standardization  $L_{s,p}$ , (ii) the section's nearness in Wikipedia  $Q_{s,p}$ , and (iii) the section's phrasings, or the likelihood of  $s$  being a phrase in light of worldwide and nearby settings. The stickiness of  $s$ ,  $p$ , is formally characterized in Eq. (2), which catches the three variables: Length standardization. As the key of tweet division is to extricate significant expressions, longer fragments are favoured for safeguarding all the more topically particular implications.

Leave  $js_j$  alone number of words in portion  $s$ . The standardized fragment respectably reduces the punishment on long fragments. Nearness in Wikipedia. In our system, Wikipedia serves as an outer word reference of substantial names or expressions. Give  $wiki_{s,p}$  and  $wikia_{s,p}$  a chance to be the quantity of Wikipedia passages where  $s$  shows up in any structure and  $s$  shows up as Tweets are viewed as uproarious with heaps of casual shortened forms furthermore, syntactic mistakes. In any case, tweets are posted predominantly for data sharing and correspondence among numerous reasons.

**Perception 1.** Word collocations of named substances and basic expressions in English are very much protected in Tweets. Numerous named elements and regular expressions are safeguarded in tweets for data sharing and spread. In this sense,  $Pr_{s,p}$  can be evaluated by checking a section's appearances in an extensive English corpus (i.e., worldwide connection). In our execution, we swing to Microsoft Web N-Gram corpus [31]. This N-Gram corpus is gotten from all web records filed by Microsoft Bing in the EN-US market. It gives a decent gauge of the measurements of generally utilized expressions as a part of English.

**Perception 2.** Numerous tweets contain helpful etymological highlights. Albeit numerous tweets contain untrustworthy semantic components like incorrect spellings and untrustworthy capitalizations [3], there exist tweets formed in legitimate English. For instance, tweets distributed by authority records of news offices, associations, and publicists are regularly elegantly composed. The semantic components in these tweets empower named element acknowledgment with moderately high precision.

**Perception 3.** Tweets in a focused on stream are not topically free to each other inside a period window.

## IV. SAMPLE CODE

```
<% @page contentType="text/javascript" pageEncoding="UTF-8"% >
var body= document.getElementById("body");
<% java.util.ArrayList widgetids= new java.util.ArrayList();
if(request.getParameter("wid")==null || request.getParameter("wid").equals("All"))
{
widgetids.add("613942941946527744");
widgetids.add("613943438036242433");
widgetids.add("613943846225862656");
widgetids.add("613944265220030464");
widgetids.add("613944907837747200");
widgetids.add("613945561771700224");

widgetids.add("613945983643156480");
widgetids.add("613947458519871489");
widgetids.add("613947740792328192");
widgetids.add("613948185451495424");
widgetids.add("613948485700751360");
widgetids.add("613948910059393024");
```



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 7, July 2016

```
widgetids.add("613950796518944770");
widgetids.add("613952410902335488");
widgetids.add("613953429803659265");
widgetids.add("614078054382440450");

}
else
{
    widgetids.add(request.getParameter("wid"));
}

int ws=widgetids.size();
for(int i=1;i<=ws; i++)
{
    String cl= "odd";
    if(i%2==0)
    {

        cl="even";
    }
    %>
    var v="<h1><%=widgetids.get(i-1)%></h1><div class='<%=cl%>' id='<%=widgetids.get(i-1)%>' </div><a
target='_blank' href='viewOne.jsp?wid=<%=widgetids.get(i-1)%>'>View This Only</a>";
    body.innerHTML = body.innerHTML + v+"<hr>";
    var config<%=i%> = {
    "id": '<%=widgetids.get(i-1)%>',
    "domId": '<%=widgetids.get(i-1)%>',
    "maxTweets": 100,
    "enableLinks": true
    };
    twitterFetcher.fetch(config<%=i%>);
    <%
    }
    %>
```

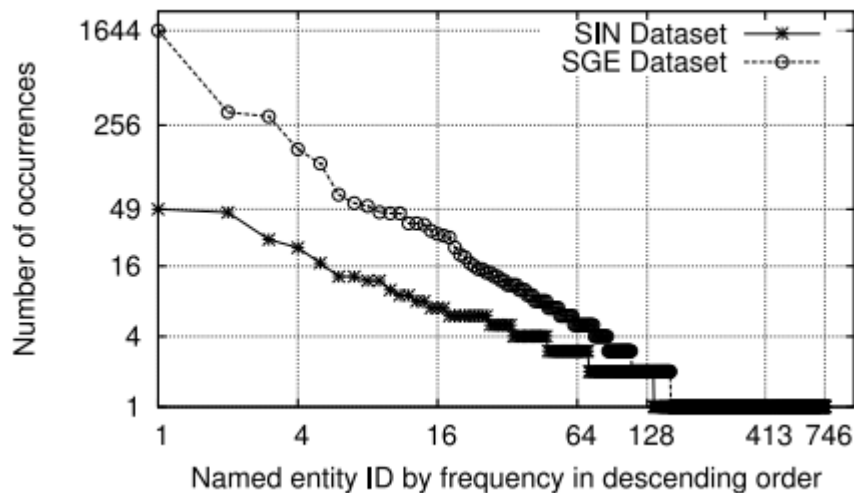
## V. SIMULATION RESULTS

**Tweet Data Sets.** We used two tweet data sets in our experiments: SIN and SGE. The two data sets were used for simulating two targeted Twitter streams. The former was a stream consisting of tweets from users in a specific geographical region (i.e., Singapore in this case), and the latter was a stream consisting of tweets matching some predefined keywords and hashtags for a major event (i.e., Singapore General Election 2011). We randomly selected 5;000 tweets published on one random day in each tweet collection. Named entities were annotated by using BILOU schema [4], [14]. After discarding retweets and tweets with inconsistent annotations, 4;422 tweets from SIN and 3;328 tweets from SGE are used for evaluation. The agreement of annotation on tweet level is 81 and 62 percent for SIN and SGE respectively. The relatively low agreement for SGE is mainly due to the strategy of handling concepts of GRC and SMC, which refer to different types of electoral divisions in Singapore. 8 Annotators did not reach a consensus on whether a GRC/SMC should be labeled as a location name (e.g., “aljunied grc” versus “aljunied”). Table 2 reports the statistics of the annotated NEs in the two data sets where fgs denotes the number of occurrences (or frequency) of named entity s (which is also a valid segment) in the annotated ground truth G. Fig.3 plots the NEs’ frequency distribution..

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 7, July 2016





# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 7, July 2016

Go to Post Page 67       Facebook Stories on Twitter: "In celebration of today's ruling and Pride Month, watch Justin Kamimoto's story #LoveWins #PrideConnectsUs http://t.co/3NaY689CyJ"
Go to Post Page 166       Facebook on Twitter: "Meet the little voices inside your head with the Inside Out Sticker Pack. Download it here: http://t.co/Ux71MUavvf http://t.co/UghcNANckB"
Go to Post Page 95       Facebook on Twitter: "Today's Q&A with Mark live from Facebook HQ will begin at 4:30 PT. Tune into the live stream here: https://t.co/fghLNY52u9"
Go to Post Page 261       Facebook on Twitter: "Soak up the happiness with SpongeBob SquarePants and Friends! Download the sticker pack here: http://t.co/D3JYBICLxs http://t.co/Ky0AqV0Cs0"
Go to Post Page 668       Facebook on Twitter: "Thousands of Nepal earthquake survivors need our help. Donate now and Facebook will match your donation. http://t.co/KuSiClamx"
Go to Post Page 242       Facebook on Twitter: "Mark Zuckerberg is answering questions on his Timeline in an hour-long live Q&A happening now: https://t.co/kd0CVz8RqW"
Go to Post Page 122       Facebook on Twitter: "In 15 minutes, watch the keynote from the second day of #F8 live here: https://t.co/V6A2R3eJxZ"
Go to Post Page 187       Facebook on Twitter: "In 15 minutes, tune into the live stream of today's keynote from #F8 here: https://t.co/V6A2R3eJxZ"
Go to Post Page 183       Facebook on Twitter: "Join Jimmy Kimmel for a live Q&A on his Page now: https://t.co/R2pQhVn4U"
Go to Post Page 152       Facebook on Twitter: "Today's Q&A with Mark live from Barcelona, Spain will begin at 9am PT. Tune into the live stream here: https://t.co/EMudjm9sQV"
Go to Post Page 231       Facebook on Twitter: "Watch Mark Zuckerberg's keynote live from Mobile World Congress now: http://t.co/Dh92NkxQgY"
Go to Post Page 516       Facebook on Twitter: "Join actor/director Hank Azaria for a Q&A live on his Page now: https://t.co/hHuDSYJ55Z"
Go to Post Page 660       Facebook on Twitter: "Love is for sharing. #WhatsLoveIn4Words http://t.co/CivXVP5j7T"
Go to Post Page 844       Facebook on Twitter: "Celebrate your friends and all they do with the Friendship Sticker Pack. Download it here: http://t.co/JZLOEsyYk http://t.co/uPFH2HfoT"
Go to Post Page 626       Facebook on Twitter: "Meet The Press is answering questions about tonight's State of the Union live: https://t.co/OJXa8yRzaS"
Go to Post Page 542       Facebook on Twitter: "Join Emmy Rossum, star of Shameless on Showtime, for a live Q&A now: https://t.co/UOUJ5xzq8P"
Go to Post Page 567       Facebook on Twitter: "Today's Q&A with Mark Zuckerberg live from Bogotá, Colombia will begin at 1 pm PT. Tune into the live stream: http://t.co/264dX26jZ6"
Go to Post Page 531       Facebook on Twitter: "Join Moisés Naim, author of The End of Power, the first book in A Year of Books for a live Q&A now: https://t.co/vCFzIkeZE"
Go to Post Page 803       Facebook on Twitter: "Cozy up with the Home for the Holidays Sticker Pack, available for download here: http://t.co/Dv0BaWp6SY http://t.co/yQmgB8T3OJ"
Go to Post Page 1627       Facebook on Twitter: "Stickered for Messenger is available for download now in the App Store: http://t.co/RQGCG8XQzm and the Play Store: http://t.co/R9UCIheKa4"

Top 3 Tweet    Top Tweet

Fig.4 TOP 3 Segmented Tweets Screen

Go to Post Page 67       Facebook Stories on Twitter: "In celebration of today's ruling and Pride Month, watch Justin Kamimoto's story #LoveWins #PrideConnectsUs http://t.co/3NaY689CyJ"
Go to Post Page 166       Facebook on Twitter: "Meet the little voices inside your head with the Inside Out Sticker Pack. Download it here: http://t.co/Ux71MUavvf http://t.co/UghcNANckB"
Go to Post Page 95       Facebook on Twitter: "Today's Q&A with Mark live from Facebook HQ will begin at 4:30 PT. Tune into the live stream here: https://t.co/fghLNY52u9"
Go to Post Page 261       Facebook on Twitter: "Soak up the happiness with SpongeBob SquarePants and Friends! Download the sticker pack here: http://t.co/D3JYBICLxs http://t.co/Ky0AqV0Cs0"
Go to Post Page 668       Facebook on Twitter: "Thousands of Nepal earthquake survivors need our help. Donate now and Facebook will match your donation. http://t.co/KuSiClamx"
Go to Post Page 242       Facebook on Twitter: "Mark Zuckerberg is answering questions on his Timeline in an hour-long live Q&A happening now: https://t.co/kd0CVz8RqW"
Go to Post Page 122       Facebook on Twitter: "In 15 minutes, watch the keynote from the second day of #F8 live here: https://t.co/V6A2R3eJxZ"
Go to Post Page 187       Facebook on Twitter: "In 15 minutes, tune into the live stream of today's keynote from #F8 here: https://t.co/V6A2R3eJxZ"
Go to Post Page 183       Facebook on Twitter: "Join Jimmy Kimmel for a live Q&A on his Page now: https://t.co/R2pQhVn4U"
Go to Post Page 152       Facebook on Twitter: "Today's Q&A with Mark live from Barcelona, Spain will begin at 9am PT. Tune into the live stream here: https://t.co/EMudjm9sQV"
Go to Post Page 231       Facebook on Twitter: "Watch Mark Zuckerberg's keynote live from Mobile World Congress now: http://t.co/Dh92NkxQgY"
Go to Post Page 516       Facebook on Twitter: "Join actor/director Hank Azaria for a Q&A live on his Page now: https://t.co/hHuDSYJ55Z"
Go to Post Page 660       Facebook on Twitter: "Love is for sharing. #WhatsLoveIn4Words http://t.co/CivXVP5j7T"
Go to Post Page 844       Facebook on Twitter: "Celebrate your friends and all they do with the Friendship Sticker Pack. Download it here: http://t.co/JZLOEsyYk http://t.co/uPFH2HfoT"
Go to Post Page 626       Facebook on Twitter: "Meet The Press is answering questions about tonight's State of the Union live: https://t.co/OJXa8yRzaS"
Go to Post Page 542       Facebook on Twitter: "Join Emmy Rossum, star of Shameless on Showtime, for a live Q&A now: https://t.co/UOUJ5xzq8P"
Go to Post Page 567       Facebook on Twitter: "Today's Q&A with Mark Zuckerberg live from Bogotá, Colombia will begin at 1 pm PT. Tune into the live stream: http://t.co/264dX26jZ6"
Go to Post Page 531       Facebook on Twitter: "Join Moisés Naim, author of The End of Power, the first book in A Year of Books for a live Q&A now: https://t.co/vCFzIkeZE"
Go to Post Page 803       Facebook on Twitter: "Cozy up with the Home for the Holidays Sticker Pack, available for download here: http://t.co/Dv0BaWp6SY http://t.co/yQmgB8T3OJ"
Go to Post Page 1627       Facebook on Twitter: "Stickered for Messenger is available for download now in the App Store: http://t.co/RQGCG8XQzm and the Play Store: http://t.co/R9UCIheKa4"

Top 3 Tweet    Top Tweet

Fig.5 TOP 1 Segmented Tweets Screen





# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 7, July 2016

## VI. CONCLUSION AND FUTURE WORK

In this paper, we introduce the HybridSeg structure which sections tweets into significant expressions called fragments utilizing both worldwide and nearby connection. Through our system, we show that nearby phonetic elements are more solid than term-reliance in managing the division process. This discovering opens open doors for instruments created for formal content to be connected to tweets which are accepted to be substantially more loud than formal content. Tweet division safeguards the semantic importance of tweets, which therefore advantages numerous downstream applications, e.g., named substance acknowledgment. Through examinations, we demonstrate that fragment based named substance acknowledgment techniques accomplishes much better exactness than the word-based option. We distinguish two headings for our future examination. One is to facilitate enhance the division quality by considering more nearby elements. The other is to investigate the viability of the division based representation for undertakings like tweets rundown, seek, hashtag suggestion, and so forth..

## REFERENCES

1. C. Li, J. Weng, Q. He, Y. Yao, A. Datta, A. Sun, and B.-S. Lee, "Twiner: Named entity recognition in targeted twitter stream," in Proc. 35th Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval, 2012, pp. 721–730.
2. C. Li, A. Sun, J. Weng, and Q. He, "Exploiting hybrid contexts for tweet segmentation," in Proc. 36th Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval, 2013, pp. 523–532.
3. A. Ritter, S. Clark, Mausam, and O. Etzioni, "Named entity recognition in tweets: An experimental study," in Proc. Conf. Empirical Methods Natural Language Process., 2011, pp. 1524–1534.
4. X. Liu, S. Zhang, F. Wei, and M. Zhou, "Recognizing named entities in tweets," in Proc. 49th Annu. Meeting Assoc. Comput. Linguistics: Human Language Technol., 2011, pp. 359–367.
5. X. Liu, X. Zhou, Z. Fu, F. Wei, and M. Zhou, "Extracting social events for tweets using a factor graph," in Proc. AAAI Conf. Artif. Intell., 2012, pp. 1692–1698.
6. A. Cui, M. Zhang, Y. Liu, S. Ma, and K. Zhang, "Discover breaking events with popular hashtags in twitter," in Proc. 21st ACM Int. Conf. Inf. Knowl. Manage., 2012, pp. 1794–1798.
7. A. Ritter, Mausam, O. Etzioni, and S. Clark, "Open domain event extraction from twitter," in Proc. 18th ACM SIGKDD Int. Conf. Knowledge Discovery Data Mining, 2012, pp. 1104–1112.
8. X. Meng, F. Wei, X. Liu, M. Zhou, S. Li, and H. Wang, "Entitycentric topic-oriented opinion summarization in twitter," in Proc. 18th ACM SIGKDD Int. Conf. Knowledge Discovery Data Mining, 2012, pp. 379–387.
9. [9] Z. Luo, M. Osborne, and T. Wang, "Opinion retrieval in twitter," in Proc. Int. AAAI Conf. Weblogs Social Media, 2012, pp. 507–510.
10. [10] X. Wang, F. Wei, X. Liu, M. Zhou, and M. Zhang, "Topic sentiment analysis in twitter: a graph-based hashtag sentiment classification approach," in Proc. 20th ACM Int. Conf. Inf. Knowl. Manage., 2011, pp. 1031–1040.
11. [11] K.-L. Liu, W.-J. Li, and M. Guo, "Emoticon smoothed language models for twitter sentiment analysis," in Proc. AAAI Conf. Artif. Intell., 2012, pp. 1678–1684.
12. [12] S. Hosseini, S. Unankard, X. Zhou, and S. W. Sadiq, "Location oriented phrase detection in microblogs," in Proc. 19th Int. Conf. Database Syst. Adv. Appl., 2014, pp. 495–509.
13. [13] C. Li, A. Sun, and A. Datta, "Twevent: segment-based event detection from tweets," in Proc. 21st ACM Int. Conf. Inf. Knowl. Manage., 2012, pp. 155–164.
14. [14] L. Ratinov and D. Roth, "Design challenges and misconceptions in named entity recognition," in Proc. 13th Conf. Comput. Natural Language Learn., 2009, pp. 147–155.
15. [15] J. R. Finkel, T. Grenager, and C. Manning, "Incorporating nonlocal information into information extraction systems by Gibbs sampling," in Proc. 43rd Annu. Meeting Assoc. Comput. Linguistics, 2005, pp. 363–370.

## BIOGRAPHY

**Mr Praveen M**, is M.Tech student in Vardhaman College of Engineering , Hyderabad, Telangana, India.

**Mr Sk. Prasanth**, is working as an Associate Professor in Vardhaman College of Engineering , Hyderabad, Telangana, India. He is pursuing Ph.D from JNTUH. His interesting subjects are C & Data Structures, Network Security, Operating System , Linux Programming, Computer Organization, Theory Of Computation, Computer Networks, Mobile Computing, Cloud Computing.