# ASNA: An Intelligent Agent for Retrieving and Classifying News on the Basis of Emotion-Affinity

Shaikh Mostafa Al Masum[1], Md. Tawhidul Islam[2], Mitsuru Ishizuka[3]

[1,3]*Department of Information and Communication Engineering, University of Tokyo, Japan*
[2]*Biotechnology Research Institute, Macquarie University, Sydney, Australia*
*mostafa_masum@ieee.org, ishizuka@i.u-tokyo.ac.jp, mislam@micros.com*

## Abstract

*In this paper, we present a system Affect Sensitive News Agent (ASNA) developed as a news aggregator that fetches news employing several RSS news-feeds and auto-categorizes the news according to affect sensitivity. There are three main factors that distinguish our work from other similar ones. First, we have integrated the approach to sense affective information from news-texts by applying a cognitive theory of emotions known as the OCC model that none have ever considered for news classification. Second, instead of any machine learning algorithm, we used common-sense and current-affairs as our knowledgebase with a rule based approach to assess each line of text by assigning a numerical valence and finally, natural language processing (NLP) technologies are used to perform automated categorization of news stories on the basis of emotional affinity. Relying on these paradigms and content analysis technologies, we have developed a news-browser that can fetch the news from RSS news-feeds and categorizes the theme of the news according to eight emotion-types plus a neutral category for quicker and intuitive understanding.*

## 1. Introduction

The primary goal in developing the system described in this paper is to demonstrate the feasibility of categorizing news stories according to their emotional affinity using natural language processing techniques for quicker and intuitive understanding. The classification and synthesized retrieval of the large amount of news articles from the Web has been a topic attracting much research effort (e.g., [1], [2], [3], [5]) but none has ever considered to sense affective information from news-texts for grouping those on the basis of affective senses and largest drawback of these systems are that they are all based on static corpora of published news articles. We have followed a deep approach to synthesize news-text and classified those according to the concept of emotion types. To meet this objective we have developed a linguistic tool called SenseNet that employed common-sense and current-affairs knowledge to assign each line of text a numerical valence to be assessed by the rule-based implementation of the OCC [11] emotion-model.

### 1.1. Background

News categories include topics, industries, proper names, and geographic information. There are mainly two tasks involved in such classification. First, document indexing is needed in order to transform the natural language text into a numerical representation suitable for further processing. The second task is the actual classification. According to literatures, different techniques, both statistical and knowledge-based, have been followed to perform the above tasks.

To achieve qualitative and efficient categorization of news-text naive Bayesian method [4], support vector machines [1], decision trees [2], pattern matching [7] techniques are been employed but these are highly dependent of static corpora of previously published news articles and hence the sentences like *"...captured the gold medal at the summer Olympics..."* or *"...the battle on center court at Wimbledon...',* are classified as war/disorders based on the lexical affinity of the words *"capture"* and *"battle"* to *"war".* The categorization method discussed in [5] has produced high accuracy, consistency, and flexibility using both knowledge-based natural language processing techniques and statistical techniques. Ontology-based text categorization in which the domain ontologies are automatically acquired through morphological rules and statistical methods is also been implemented in [3]. A technique for personalized article classification exploiting user's awareness of a topic has showed better performance in order to classify articles in a 'per-user' manner [6].With the combination of taxonomy-based topic matching and personalized

word-list the technique expressed in [6] measures the distance between sense-derived keywords in the user's profile and words matched in the news feed to output user-focused categorization. News analysis system (NAS) [7] extracts stories from a newswire, parses the sentences of the story, and then maps the syntactic structures into a concept base. This process results in an index containing both general categories and specific details that matched the concept. Transforming each news-document into a vector of weights corresponding to an automatically chosen set of keywords and then applying either k-NN (nearest-neighborhood) [2] or cosine similarity method are used to compare the keyword vector of the news story to the feature vectors. Different threshold values are used for different categories to notice 91.4% success rate for news classifications.

Several researches have been performed to analyze sentiment expressed through text. For example, *Sentiment!* [8], is a commercial application that reads news articles and shows if they are positive, negative or neutral claiming 85% accuracy against human analysts. Affective-News Theory [9] conceptualizes news as having (different) story structures; the inverted pyramid among others; certain structures meet intuitions on 'storyhood' by evoking specific emotional reactions (e.g. suspense or curiosity based on event and discourse structure) to different story structures in news. Approach mentioned in [10] used a sentiment analysis dictionary having 3,513 entries and instead of analyzing the favorability of the whole context each statement on favorability is extracted, and present them to the end users so that they can use the results according to their application requirements. But the system outputs -1 to indicate a negative sentiment for the sentence *"It's difficult to take a bad picture with this camera.",* whereas this is a positive statement for the camera.

## 1.2. Our Approach

We admit that analysis of favorable or unfavorable opinions or emotion-affinity is a task requiring emotional intelligence and deep understanding of the textual context, involving common-sense and domain knowledge as well as linguistic knowledge. The interpretation of opinions is usually debatable affair even for humans. However the system, ASNA, is an attempt towards this task. The approach of our system is quite straightforward and the step by step operation of the system is indicated in Figure 1. First a user chooses the sources of news according to his/her domain of interest. In this case we used RSS [14] feeds as the sources for the news. The justification of using RSS feeds are explained in section 2.1. After the news sources are selected, News Fetcher collects the news as tuples of news topic and brief story corresponding to the topic by parsing the results returned by the RSS feeds. Then the plain-text tuples are parsed by a language parser. We have implemented a deep parsing technique to output tuple(s) of Subject, Subject Type, Subject Attributes; Action, Action Status, Action Attributes; and Object, Object Type, Object Attributes for each line of text. The output of language parser is assessed by a linguistic tool SenseNet that we have developed employing WordNet [12] and ConceptNet [13]. SenseNet considers each tuple as a Sense and outputs a numerical value for each lexical-unit (e.g. sentence). Affect Sensing Engine then classifies the news-texts according to eight emotion-types namely, Happy, Sad, Hopeful, Fearful, Admirable, Shameful, Loveable, and Hatred plus a Neutral category. Finally a user can browse the news according to the emotion groups.
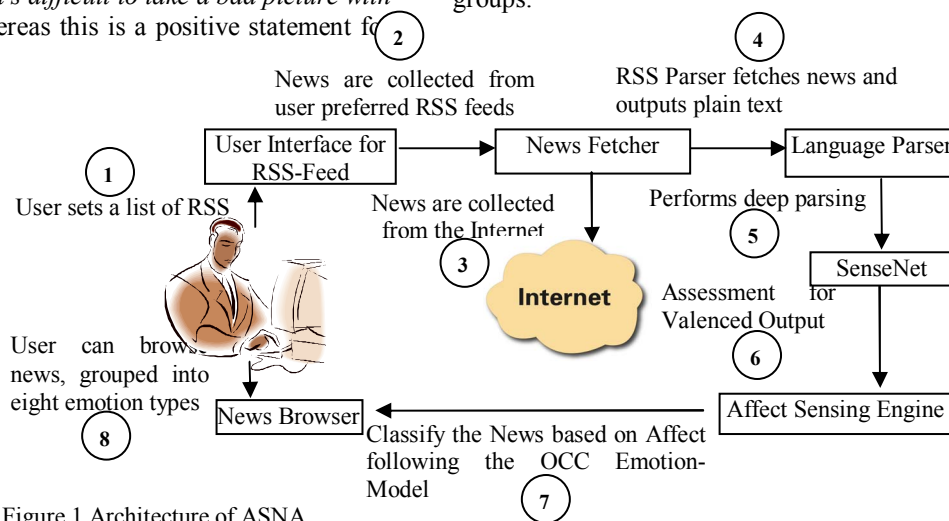


Figure 1 Architecture of ASNA

## 2. Implementation

The system ASNA consists of several modules namely, User Interface; News Fetcher; Language Parser; SenseNet; Affect Sensing Engine and News Browser. The Language Parser, SenseNet and Affect Sensing Engine are integrated as a server application written in Python. Others are written in C#.

### 2.1. RSS-Feeds

Most of the systems for news categorization primarily target to cluster news according to specific domains (e.g. sports, war, business, technology etc.), but this problem has been solved explicitly by RSS [14] technology. The RSS 1.0, 2.0 and ATOM standards (see [23] for detail) include categorical information for each news item, which enables a more elegant way to determine the domain of a news item than matching keywords against the item's text body. But still, the problem of intelligent filtering of information exists. In general, one can subscribe to a web-site's (e.g. MyYahoo!) RSS feed using a desktop news aggregator to get the news of one's domain of interest. If, for example, 10 RSS news feeds are subscribed by a user, and each news feed delivers 10 news items per day on average, then the user will have to filter through 100 news items in total per day. So one of the problems of desktop news aggregation is the issue of information overload and hence grouping news by studying the relationship between natural language and affective information by a theory of cognitive appraisal for emotion might be worthwhile in this case.

### 2.2. News Fetcher

In this scenario, a user either selects from a list of provided news-feeds or can add others according to his/her preference towards specific news domains. Otherwise the system uses a default list of news-feeds as the source of news to fetch. The system then requests the RSS news feed to provide the topic of the news along with a brief story (usually 1 or 2 lines) corresponding to it from the provider's web server and receive an XML-like data. This data of RSS feed is then parsed by an RSS-feed parser, which extracts category tags from the news items. The detail parsing techniques of RSS-feeds is not in the scope of this paper and hence it is not discussed here.

### 2.3 Language Parser

We are using the *Machinese Syntax* [15] program to obtain XML-formatted shallow-parsed information for an input sentence for further processing. As an example, for the input sentence, *"Two members of Tonga's royal family were killed when a teenager racing her car crashed into their vehicle."*, we obtain XML-like syntactical information from the parser, which is further processed to output as a tuple of Subject, Subject Type, Subject Attributes; Action, Action Status, Action Attribute; Object, Object Type and Object Attribute, as indicated in Figure 2. Since a tuple is initiated with an occurrence of a verb in the sentence, we may obtain multiple tuples if deep-parser encounters multiple verbs in a sentence. A tuple encodes information about *"who is associated with what and how"*. The output given in Figure 2 has three such tuples.

[[['Subject:' 'member', 'Subject Type:' 'Person', 'Subject Attrib:' ['quantity: two', 'N GEN SG: tonga', 'A ABS: royal', 'N NOM: family']]
['Action:' 'kill', 'Action Status:' 'Past Particle', 'Action Attrib:' ['time: when']],
['Concept:' '', 'Concept Type:' '', 'Concept Attrib:'']]],
[['Subject: ' 'teenager', 'Subject Type:' 'Person', 'Subject Attrib: [ ]],
['Action:' 'race', 'Action Status:' 'Continuous', 'Action Attrib:'[]],
['Concept:' 'car', 'Concept Type:' 'N NOM', 'Concept Attrib:' ['PRON PERS GEN SG3:she']]],
[['Subject:' 'car', 'Subject Type:' 'Other 3rd', 'Subject Attrib:' [ ]],
['Action:' 'crash', 'Action Status:' 'past', 'Action Attrib:' ['goal: vehicle']],
['Concept:' 'vehicle', 'Concept Type:' 'N NOM', 'Concept Attrib:' ['PRON PERS GEN PL3: they']]]

Figure 2: Output of deep-parse

### 2.4 SenseNet

In a linguistic context, as e.g. in WordNet [20], a word sense is a given meaning of a word based on the context. Unlike WordNet, by the term "sense" used in SenseNet, we mean a lexical tuple, formed by 'a subject or agent', 'a verb or action', 'an object or concept' and associated 'adjectives or attributes' and each sense is assigned a value that we call sense-valence. SenseNet employs two lexical resources namely, WordNet and ConceptNet [21]. A sentence may contain several such senses. For instance, Figure 2 indicates three "senses" for the input sentence.

#### 2.4.1 Knowledge-Base of SenseNet

*2.4.1.1 Action and Adjective Polarity* A group of students and volunteers have manually counted the positive and negative senses of each word of our customized list of verbs and adjectives according to the contextual understanding of each sense appeared in WordNet 2.1; and thus we maintain a database of scored verbs and adjectives. For example, a verb's score is stored as following tuple-like format:

verb-word [*Positive Sense Count, Negative Sense Count, Prospective Value, Praiseworthy Value,*

*Polarity Value*]. An excerpt from database is given to illustrate the idea.

| appear | 6 | 1 | 3.571 | 4.286 | 3.929 |
|---|---|---|---|---|---|
| applaud | 2 | 0 | 5.000 | 5.000 | 5.000 |
| appreciate | 5 | 0 | 5.000 | 5.000 | 5.000 |
| approve | 2 | 0 | 5.000 | 5.000 | 5.000 |
| arrest | 1 | 3 | -2.500 | 1.250 | -0.625 |

Table 1: An excerpt from verb database

The formula used to calculate the values (scale of -5 to 5) are; for each word,

*Polarity Value* = Average (((*Positive Sense Count - Negative Sense Count*) / Total Sense Count) * 5.0)

*Prospective Value*= Average ((*Positive Sense Count* / Total Sense Count) * 5.0)

*Praiseworthy Val* = Average (*Polarity Value + Prospective Value*)

The adjectives are also assigned a similar type of numerical value based on the count of senses found in the WordNet. At present we scored 723 verbs, 205 phrasal verbs, 237 adjectives related to shape, time, sound, taste/touch, condition, appearance and 711 adjectives related to emotional affinity.

***2.4.1.2 Domain-knowledge*** we have also developed a knowledge-base of current-affairs and stored as a database of scored named-entities. For example, an entity's score is stored as following tuple-like format: Named-entity [*Type, Role, General-Sentiment*], the field *Type* indicates whether the entity indicates a person (living body), company, or an object and the *Role* stores a keyword to represent the concept of the entity. The *General-Sentiment* field indicates either a negative (-1) or positive (+1) impression towards the named-entity. An excerpt from database is given to illustrate the idea.

| Bin Laden | Person | Militia | -1 |
|---|---|---|---|
| Discovery | Object | Skyrocket | 1 |
| George W. Bush | Person | President | -1 |
| Harold Pinter | Person | Scientist | 1 |
| IBM | Company | Electronics | 1 |
| Katrina | Object | Cyclone | -1 |
| Kofi Annan | Person | Official | 1 |
| Microsoft | Company | Software | 1 |
| NASA | Company | Research | 1 |

Table 2: An excerpt from News-Domain Knowledge

The value for *General Sentiment* is a subject to personal-view or opinion. But in general we assigned negative values for those entities that are usually associated with wars, crime or negative concept. Collection of such entities and scoring is still in a surveying and reviewing phase. At present the system has 2000 such named entities that serve as the knowledge-base of the current-affairs.

**2.4.2 Assumptions for SenseNet**

The rules and algorithms of SenseNet are based on the following assumptions.

*Assumption 1: A concept or named-entity has a valence.* SenseNet maintains a growing list of concepts scored with the help of ConceptNet. A named-entity can be represented by its' type and valence can be calculated by considering the valence of the role and general sentiment. For example, the sentence "*Nearly a year after Katrina flooded New Orleans, the city still does not have a plan for rebuilding",* the valence of 'Katrina' is set according to the concept-valence of "Cyclone" (-4.433) and moreover the general sentiment (-1) validates the negative polarity of the assigned valence, on the contrary for the entity *"George W. Bush"* or *"Bush"* SenseNet will first get the valence of "President" (3.42), which is a positive value but the general sentiment value (-1) arises a contradiction to the polarity of the obtained valence. For such cases SenseNet considers such entities ambiguous and sets the value to 0 to indicate the valence of such entities as neutral.

*Assumption 2: An action or verb is the core of a sense-unit accompanied by a concept and/or a subject/actor and/or adjective/attributes.* The smallest unit of the SenseNet processing element is the sense-unit and the core element is a verb. A valid sense-unit must have a verb and a concept associated with that verb. If a verb has a missing concept, a positive concept is imagined to form the sense-unit for that verb. So a 'sense' may be formed by a sense-unit with or without a subject and associated attributes.

*Assumption 3: A sense-unit outputs either a negative, positive or neutral valence.* For the input, *'President Bush called the space shuttle Discovery on Tuesday to wish the astronauts well, congratulate them on their space walks and invite them to the White House.'* The sense-units are: [call, Discovery], [wish, astronauts], [congratulate, them] and [invite, them]. The rules to assign the polarity sign of sense-unit are:

- Neg. Verb + Pos. Concept → Neg. Polarity (e.g. quit job)
- Neg. Verb + Neg. Concept→ Pos. Polarity (e.g. quit drug)
- Pos. Verb + Pos. Concept→ Pos. Polarity (e.g. buy car)
- Pos. Verb + Neg. Concept → Neg. Polarity (e.g. buy gun, encourage terrorist)

The valence is calculated by adding the scores of both verb and concept.

*Assumption 4: Intensifier and Modifier-* An adjective deals with intensity of the sense-valence and concept-valence of actor may modify the polarity of a sense-valence. As examples, *"The hurricane of the season has formed",* and *"The first hurricane of the season has formed"*; if the intensity of negative sense of the sense-unit ("form-hurricane") for the first sentence is neutral, but the intensity of the negative sense for the second one is higher because of concept-intensifier 'first'. Similarly the intensity of the positive sense of the

sentence *"President Bush has a straightforward message for Russian leader"* is higher than that of the sentence *"President Bush has a message for Russian leader"* for the word 'straightforward'.

*Assumption 5: Valence of an 'Abstract-Concept' for an input concept can be assigned by the action(s) valences that possibly are performed by that concept.* It is tedious to enlist all the key-concepts and Abstract-Concept because the list might be too long. If a concept is not found in the database, ConceptNet 2.0's function *DisplayNode()* is employed and it returns all the possible semantically connected entities that ConceptNet has found for the concept. We then make two groups of semantic relations; in the first group we collect all the entries for the relations like 'IsA', 'DefinedAs', 'MadeOf', 'PartOf' and the second group enlists the entries for the relations like, 'CapableOf', 'UsedFor', 'CapableOfReceivingAction'. The first list is again searched for any matching concept in the list. If it fails, from the second list which is actually a list of verbs, the first 5 unique verbs or actions are matched with the verb list and an average score for those verbs is retuned as the concept-valence.

*Assumption 6: The average value of sense-valences of a sentence,* S, *is the sense-degree of that sentence.* If a sentence, S has N many senses, the sense-degree of the sentence, S is assigned as:

| Sense-Degree(S) | = average (abs (sense1_valence) + abs (sense2_valence) +…… abs (senseN_valence))

The polarity sign of the sense-degree is set according to the sign of the sense-valence which value is the maximum among the sense-valences of that sentence. These assumptions are explained in the next section with an example.

### 2.4.3 SenseNet Processing

SenseNet processes each sense according to the rules and algorithms stated in the previous section. How valences are assigned to the input sentence, indicated in Figure 2 is discussed below.

For sense1, sense-unit (kill, positive-concept), is formed since it does not contain a concept. SenseNet looks up the verb list for the score of 'kill' and gets the value -3.667 and Sense1 has no adjective attributes. So intensity is set to neutral and Actor type being 'Person' compels SenseNet to resolve the concept-valence for the actor ('Member'), and assigns a positive value 3.625, according to assumption 5. ConceptNet 2.1 server returns two lists, *Possible_concept_list* and *Possible_action_list* for the concept 'Member' as explained in the previous section. In this case SenseNet first tries with the list, *Possible_concept_list* and it fails to assign a value. So the second list, *Possible_action_list*, is processed and from the second list SenseNet returned the value 3.625 by averaging the scores of the verbs ('pay'; 'attract'; 'impress'; 'attend' and 'marry') found in the scored-list. Similarly the

value 3.333 is assigned as the concept-valence for the concept 'family'. The average-score of 'family' and 'member' (3.479) is set as the valence of the subject and intensity of subject's valence is set high for the adjective 'Royal' which is used to decide emotional intensity. Since sense1 does not have any accompanied concept, Valence of sense1 = abs (action valence) +5.00 = 8.667, and according to the formula, negative action with a positive concept (assumed in this case) give a negative sense, the polarity sign of the above value is set negative. Moreover SenseNet assigned the actor's valence with a positive polarity. So the resultant polarity of this sense-valence (8.667+3.479=12.146) is set negative (i.e. '-12.146'). For sense2, the sense-unit, (race; car), is assessed with the actor 'teenager'. SenseNet outputs +10.280 as the valence for Sense2. Similarly SenseNet assigns '-10.98' for the Sense3 and for the sentence the value of sense-degree is assigned as: abs ((-12.146) + abs (+10.280) + abs (-10.980))/3 = 11.135. The polarity is set to negative because the sign of the highest sense-valence is negative. Finally the sense-degree of the sentence becomes -11.135 which is further used to decide for the specific kind of negative emotion-type. SenseNet also takes care of the negation by reversing the polarity sign of the sense-valence.

### 2.5 Affect Sensing

The system assigns value to several OCC-emotion model inspired linguistic variables. The variables calculated for each sense are, *Action_ Name, Action_Polarity, Action_Status, Agent_ Type, Prospective_Val, Praiseworthy_Val, Sense_ Degree, Action_Likingness, Action_Deservingness, Effort_for_Action*. For instance, the variable, *Event Likingness*, is set by considering the polarity of the event and whether any determiner or adjective or adverb used to emphasize the event etc. Taking the average values of the corresponding variables *Emotion Type* and *Emotion Valence* are determined for each sentence or paragraph associated with each news-title. In this case RSS-feed parser usually returns 1 or 2 sentences associated with each news-title and hence an emotion-type is assigned by assessing those sentences. We have implemented rules for 8 emotion types following the OCC model using the aforementioned variables. Some of the rules are listed below due to space limitation.

"Happy" is true if *Sense_Degree* > 5.0, *Action_Polarity* >0.0, *Action_Status* = 'Past' or 'Present', *Agent_Type* = 'Person' or 'Company', *Action_likingness*>=2.0, *Action_Deservingness*>= 2.0 and *Effort_for_Action*>=2.0

"Hope" is true if *Sense_Degree* > 5.0, *Action_Polarity* >0.0, *Action_Status* = 'Present' or 'Future', *Agent_Type* = 'Person' or 'Company',

*Action_likingness*>=2.0, *Action_Deservingness*>= 3.0 and *Prospective_Val>=3.0*

"Love" is true if *Sense_Degree > 5.0*, *Action_Polarity >0.0*, *Action_Status* = 'Present' or 'Past', *Agent_Type = any*, *Action_likingness*>=3.0, *Action_Deservingness*>=3.0, *Praiseworthy_Val>=3.0* and *Prospective_Val>=3.0*

## 2.6 News Browser

The news browser finally enlists the news according to emotion-types and a user can browse news thereby. Figure 3 shows a snap-shot of the emotion sensitive news browser having 9 buttons. Clicking on any button shows a list of news summary corresponding to a specific emotion-affinity. An avatar reads out the news summary and a user can also view the full story of the news on this browser by clicking either on the headline or the image associated with the news.



Figure 3: Affect Sensitive News Browser

## 3 Conclusion

The linguistic approach to affect-sensing from news would strengthen human-computer interaction with fun. We plan to implement a user interface to set user-specific preferences (e.g. personal opinion about particular entities) that might help the system to perform better based on certain user-centric model. Basically, we have found two types of systems; one classifies news according to taxonomical categories and the other realizes news-topics as story-events to assess sentiment (positive or negative or neutral) and limited emotional reactions (suspense or curiosity). But none of those ever considered classifying news-articles into broad range of emotion categories. Hence the system would help the news readers to conceptualize and sense news-articles in a quick and easy manner.

## References

[1] Maria, N., and M. J. Silva, "Theme-based Retrieval of Web News", *Proceedings of the Third International Workshop on the Web and Databases*, Springer, Athens, Greece, 2001, pp. 26-37.

[2] Bacan, H., I. S. Pandzic, and D. Gulija, "Automated News Item Categorization", *Proceedings of the 19th Annual Conference of The Japanese Society for Artificial Intelligence*, Springer-Verlag, Kitakyushu, Japan, 2005, pp. 251-256.

[3] Wu, S.H. and W.L. Hsu, "SOAT: A Semi-Automatic Domain Ontology Acquisition Tool from Chinese Corpus", *Proceedings of the 19th International Conference on Computational Linguistics- Volume 2*, Association for Computational Linguistics, Taipei, Taiwan, 2002, pp. 1-5.

[4] Sebastiani, F., "Machine Learning in Automated Text Categorization", *ACM Computing Surveys*, Vol 34, No 1, March 2002, pp. 1-47.

[5] Jacobs, P., "Joining statistics with NLP for text categorization", *Proceedings of the Third Conference on Applied Natural Language Processing*, Association for Computational Linguistics, Morristown, NJ, 1992, pp. 178-185.

[6] Antonellis, I., C. Bouras, and V. Poulopoulos, "Personalized news categorization through scalable text classification", *Proceedings of the 8th Asia-Pacific Web Conf*, Springer, Harbin, China, 2006, pp. 391-401.

[7] Kuhns, R. J., "A News Analysis System", *Proceedings of the 12th conference on Computational linguistics*, Association for Computational Linguistics, Budapest, Hungry,1988, pp. 351–355.

[8]Corpora Software, UK, http://www.corporasoftware.com

[9] Knobloch, S., "Affective News- Effects of Discourse Structure in Narratives on Suspense, Curiosity, and Enjoyment While Reading News and Novels", *Communication Research*, Vol 31, No 3, June 2004, pp. 259-287.

[10] Nasukawa, T., and J. Yi, "Sentiment Analysis: Capturing Favorability Using Natural Language Processing", *Proceedings of the 2nd international conference on Knowledge capture*, ACM Press, Sanibel Island, FL, 2003, pp. 70-77.

[11] Ortony, A., G. L. Clore, and A. Collins, *The Cognitive Structure of Emotions*, Cambridge University Press, 1988.

[12] Fellbaum C., (Ed.). *WordNet: An Electronic Lexical Databases*, MIT Press, Cambridge, Massachusetts, 1999.

[13] Liu, H., and Singh, P., "ConceptNet: A Practical Commonsense Reasoning Toolkit", *BT Technology Journal*, Vol 22, No 4, Oct. 2004, pp. 211-226.

[14] RSS, Available WWW: http://www.rssboard.org/

[15]Connexor Oy, http://www.connexor.com/connexor/