# Discovering Event Evidence Amid Massive, Dynamic Datasets

Robert M. Patton
Oak Ridge National Laboratory
P.O. Box 2008 MS 6085
Oak Ridge, TN USA 37831
Ph: 1-865-576-3832

pattonrm@ornl.gov

Thomas E. Potok
Oak Ridge National Laboratory
P.O. Box 2008 MS 6085
Oak Ridge, TN USA 37831
Ph: 1-865-574-0834

potokte@ornl.gov

## ABSTRACT
Automated event extraction remains a very difficult challenge requiring information analysts to manually identify key events of interest within massive, dynamic data. Many techniques for extracting events rely on domain specific natural language processing or information retrieval techniques. As an alternative, this work focuses on detecting events based on identifying event characteristics of interest to an analyst. An evolutionary algorithm is developed as a proof of concept to demonstrate this approach. Initial results indicate that this approach represents a feasible approach to identifying critical event information in a massive data set with no apriori knowledgeof the data set.

## Categories and Subject Descriptors
I.7.0 [**Document and Text Processing**]: General

## General Terms
Algorithms, Design, Experimentation

## Keywords
Event detection, Events, Evolutionary Algorithm

## 1. INTRODUCTION
As global connectivity continues to grow, the worldwide relevance of formerly local events has become much more profound. Events that decades ago would be isolated news stories, now have the potential to create a cascading effect resulting in a global reaction. The development and expansion of the Internet as a news and information source provided much of the fuel for this cascading effect. News of events can now travel the globe at nearly the same rate at which the event itself unfolds. A recent example is that of the controversial Danish cartoons published in September 2005 [1]. Immediately, these cartoons created a strong reaction from the local Muslim community in Denmark. Within several months, this event led to a global reaction resulting in approximately 139 people dead, extensive property damage, and various bans on imports and exports between countries [1].

From the perspective of the military, intelligence, and national security, identifying such events before they create a global reaction presents many difficulties, but is of significant importance to national security in order to prepare a proper response. There are several major challenges to this. First, information concerning these events is often obscured, initially, within a massive amount of data. Consequently, this information is not likely to become "visible" until these events have created a global reaction. Second, since events may now have a global reach, information concerning these events may occur at any time within a 24-hour period, and must often be analyzed and responded to within a very short time frame. This creates a dynamic environment in which to track the latest "hot issues". Such dynamic environments pose significant challenges due to their unpredictability. Another challenge is that events have a temporally based reaction characteristic. Determining this characteristic is similar to detecting the shockwaves of an earthquake, and measuring their strength. Earthquakes vary in magnitude, and smaller earthquakes can often precede large earthquakes. In a similar way, events cause some reaction. These reactions vary in magnitude (local, regional, national, global, etc), and some events with large reactions can be preceded by events with smaller reactions (e.g., Danish cartoons). Some events cause an initial global reaction immediately. Some events cause only a local reaction that does not propagate. Unfortunately, there remains no simplistic or clearly defined approach for defining and measuring a reaction to an event. Finally, another challenge is that automation of event extraction from electronic sources remains very difficult [2][3]. Many of the difficulties are directing related to natural language processing. Unfortunately, thorough and accurate event extraction from massive data remains a laborious, manual process and is sometimes even impossible for humans to perform.

Given these challenges and the significance of a solution to the problem, the grand vision is to automatically detect events that may cause a global reaction. In an attempt to move towards this vision, this paper describes a novel approach to discovering evidence of events within massive data sets using an evolutionary algorithm. This work narrowly focuses on resolving a specific problem: Identify events of interest (if they exist) within a massive and dynamic data set without apriori knowledge of the contents of the data set and its characteristics (i.e., categories, clusters).

Section 2 provides some background information that provides the context under which this new approach was developed. Section 3 describes the design of an evolutionary algorithm (EA) for this approach. Section 4 describes some experimental testing to demonstrate the EA. Section 5 and 6 provide an analysis and conclusions, respectively.

## 2. BACKGROUND

Before discussing the approach of detecting events, event characteristics must first be described and understood. There are various characteristics about events such as the temporally based reaction characteristic described earlier. For the purposes of this work, events are defined as having two primary characteristics: foundational and descriptive [4].

A foundational characteristic is a characteristic that, when changed, dramatically alters the description of the event. Descriptive characteristics are characteristics that, when changed, do not alter the description of the event dramatically. Their presence simply enhances the detail of the event, while their absence leaves the basics of the event intact. For example, consider the 2004 Madrid train bombings [5]. The title for this event provides sufficient wording to identify clearly this event. In this case, these words can be considered foundational characteristics of the event. To change any of these words would cause the very essence of this event to be dramatically different. However, upon further investigation of this event, there are other details such as the specific trains that were bombed, the time of day, etc. It is this information that would be considered descriptive characteristics. Changing these characteristics would not dramatically change the event. The event would still be defined as the 2004 Madrid train bombings. The foundational characteristics remain intact; however, the detail of the event is enhanced by the descriptive characteristics.

Foundational characteristics for different events can be generalized to several categories such as groups of people, countries, locations, infrastructure (e.g., buildings, bridges, roads), and actions (e.g., bombing, hijacking, smuggling). Descriptive characteristics can also be generalized to some degree; however, the variety of these characteristics makes this generalization more difficult.

Given these characteristics and the degree of difficulty in the problem space, this works narrowly focuses on the development of an automated method for event detection using the foundation characteristics of events that would be of interest to an analyst.

## 2.1 Analogous & Related Works

Below are listed two analogous works whose primary characteristics are that there is some "object" whose existence needs to be identified, but the direct detection of this object is very difficult.

In the field of sunspot detection [6], if a sunspot (i.e., an active region on the surface of the sun) is facing the earth, it can cause an effect on the earth's atmosphere. Such active regions on the sun emit different characteristics than other non-active regions of the sun. In the work of [6], the authors determined that there is a particular type of radiation that can be measured such that sunspots can be detected before they begin facing the earth. This

method of indirect detection provides an early warning of increased solar activity that may affect the earth.

Another analogous work is that of indirectly detecting extrasolar planets [7]. In their work, the authors describe several approaches for detecting extrasolar planets based on indirect measurements of the environment in which the planet exists. According to their work, the very existence of a planet will produce measurable effects on the stars that they orbit. Depending on the effects, it may even be possible to determine the size of the planet and distance from the star that it orbits.

These two analogous works developed indirect approaches that measure the environments to identify if there were any detectable effects that would be caused by the existence of the object of interest. In like manner, this work is focused on indirectly detecting the existence of an event (i.e., object of interest) based on the occurrence or frequency of words (i.e., environment) that may be affected by the existence of such an event.

In the field of natural language processing, there is substantial work in automated event detection, including related works [8][9][10]. Generally, these works focus on using clustering or categorization techniques for event detection. Unfortunately, for massive data sets, clustering techniques become computationally intensive and categorization techniques become plagued by the challenge of defining appropriate categories [12]. Unlike these approaches, this work does not depend on clustering or strict categorization techniques.

## 3. DESIGN

The problem at hand is fundamentally a search problem with a dynamic, unclear, and massive search space. As a search algorithm, the EA is ideally suited for this particular problem. To perform this search, there is a three-step process:

1. The user of the EA must define what foundational characteristics (FC) are of interest.

2. Documents in the data set that match with a particular FC are identified. In this case, a particular threshold that defines how well a document matches the FC defines a match.

3. The documents that match each FC are then compared between the different FC's to identify any connections between documents. In this case, exact matches of noun and proper noun phrases (names, organizations, locations, etc) define a connection.

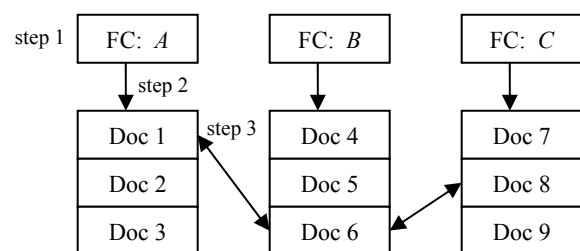Figure 1 shows a conceptual view of these steps.



**Figure 1. Conceptual View**

To illustrate this 3-step process, consider the following example. For step 1, the user has defined the FC's of interest as *A*, *B*, and *C* as shown in Figure 1. For step 2, documents that match *A* (described below) are then found in the data set. As an example, this could be documents *1*, *2*, and *3*. For step 3, documents that match with *A* are then compared to documents that match *B* and *C*. As an example, document *1* would be compared to documents *4 – 9*. In this comparison, if documents 1, 6, and 8 contained a high number of noun or proper noun phrases that were exact matches between these three documents, then we would say that the event described by *ABC* is supported by evidence contained in documents *1*, *6*, and *8*. These three documents would then be presented to an analyst. This is the overall concept of the way that this approach works.

In defining the FC's of interest, the user must develop a taxonomy of words that defines an FC as clearly as possible. For this current work, a simple list of words is used to define a particular FC; however, a more formal ontology could be used. For example, if the FC of interest is "earthquake", then the user may use the following words to define more clearly this concept:

- Earthquake(s)
- Quake(s)
- Tremor(s)
- Fault line(s)
- Richter scale

For this work, this listing of words would be represented using a Vector Space Model (VSM) [11]. Using the TF-ICF approach [12], documents in a dataset would then be compared to this VSM of the concept "earthquake". Documents that had a high similarity (based on a pre-defined threshold) would then be identified as being "earthquake" documents. By doing this, documents that simply mention the word "earthquake" one time without any other earthquake-related terms occurring would be ignored. As a result, documents can be more accurately identified as pertaining to the concept of interest. This would complete step 2 in the process defined previously.

For step 3, the selected documents that match the FC's of interest are then processed by an entity extraction module that identifies entities such as people, locations, organizations, etc. For this particular system, the LingPipe entity extractor was used [13]. In addition, noun pair phrases were extracted as well based on the Stanford Part of Speech tagger [14]. Examples of noun pair phrases would be "quake victims" or "rescue operations". After extracting entity and noun phrases for each document, the various documents for each FC are then compared based on these entity and noun phrases. If document *1* representing *A* and document *6* representing *B* both contained the noun phrase "rescue operations", then the system would classify that as a match between those 2 documents during the step 3 processing. This would complete step 3 in the process defined previously.

The following sections will describe implementation details concerning the specific encoding, fitness function, and other EA operators to implement this approach.

## 3.1 Genetic Encoding
To represent properly the solution domain in the EA, a structured EA encoding was used [15][16]. This encoding relies on a two level "virtual" structure of the genes as shown in Figure 2. The first (or higher) level is referred to as "control" genes. These genes activate or deactivate the second (or lower) level of genes. Depending on the problem to be solved, there can be *C* number of higher-level genes. The lower level genes are used to encode the actual solution parameters of interest. Depending on the problem to be solved, there can be *N* number of lower level genes.
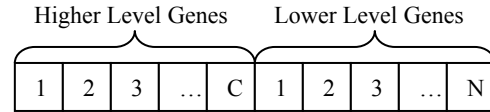
Figure 2. DNA encoding

The higher-level genes are generally encoded as binary with 1 meaning "activate" and 0 meaning "deactivate". The lower level genes are encoded as real numbers. In addition to this, each higher-level gene can have a 1-to-1 or 1-to-Many mapping to the lower level genes.

For the problem space defined in this work, the lower level genes are used to represent different foundational characteristics of interest to a person. The higher-level genes are then used to activate or deactivate various foundational characteristics to more accurate identify specific events. There is a 1-to-1 mapping of the higher-level genes to the lower level genes. For example, the encoding scheme would then consist of six genes if there were three categories of the foundation characteristics (Country, Action, and Infrastructure). Gene 4 would represent the Country characteristic. Gene 5 would represent the Action characteristic. Gene 6 would represent the Infrastructure characteristic. Genes 1 – 3 would represent the control genes for genes 4 – 6, respectively. By doing this, events may be simply characterized with just genes 4 and 6, but not 5 if gene 2 were set to 0. This allows a much wider range of possible solutions using these foundational characteristics because it includes all possible combinations of events defined by only 2 genes as well as all possible combinations of events defined by 3 genes.

## 3.2 Fitness Function
The fitness function evaluates each individual to determine how well that individual provides "evidence" to support the event that is encoded by its DNA and how unique that individual is within the EA population. In this particular work, "evidence" is defined by two criteria:

- The existence of documents that match a particular FC
- The number of noun and proper noun phrases that match between documents of different FCs.

To measure the first criteria, a search is performed on the data set to identify documents that match with an "activated" FC (i.e., a lower level gene whose higher level gene is set to *1*). If an FC fails to match with any documents in the data set, then the "activated" FC is said to be a non-contributing gene and the first criteria is not met. As a result, the individual in the EA population is penalized for having an "activated" FC with no supporting documents. For this work, the penalty was defined as shown in Equation 1 where *n* is the number of non-contributing genes in the individual. A scaling factor of 1/10 was used to provide a stronger penalty. If the FC does match with documents in the data set, then the first criterion is met, and no penalty or reward is given. As an example, consider individual i with 3

lower level genes. Suppose that gene k represents the FC "Earthquake". If there are no documents that match this FC in the data set, then the number of non-contributing genes would be 1, and the penalty would 1/10.

$$penalty = \frac{1}{10 \times n}$$

**Equation 1. Penalty function**

To measure the second criteria and given that there are $N$ lower level genes, a comparison is performed between each document $i$ of FC $k$ $(k < N)$ to each document $j$ of FC $m$ $(m < N)$ to count the number of matches between noun and proper noun phrases that were extracted as described earlier. The individual's fitness will be higher with a higher number of matches (not counting any penalties for non-contributing genes). The individual's fitness will be weaker with a lower number of matches. Once this count is performed, it is then multiplied by the penalty function shown in Equation 1. This final calculation is the individual's fitness prior to niching.

For an information analyst, it would be valuable to identify multiple events within a given data set. However, the traditional EA is designed to converge to a single, dominant individual. To compensate for this, a simple fitness sharing niching scheme was used [17]. For this work, if $N$ individuals contained exactly the same DNA encoding, then each individual's fitness value was divided by $N$.

## 3.3 Operators

For the selection operator, a tournament selection process was used with a tournament size of 3. For crossover, a 1-point crossover operator was used with a crossover rate of 0.7. For each pair of individuals, the crossover operator randomly selects a particular point at which to perform the crossover.

For mutation, a 1-point mutation operator was used and was controlled by an adaptive mutation rate defined by the function shown in Equation 2.

$$MRate(j) = \left( MR_{REF} - MR_{REF} * \frac{GSD(j)}{Max_{GSD}} \right)^{FSD}$$

**Equation 2. Adaptive Mutation Rate**

This function incorporates measures of phenotypic and genotypic diversity as well as known reference points of useful mutation rates. In this function, MRate represents the mutation rate for a particular gene locus $j$. The value $MR_{REF}$ represents a known reference point for an acceptable mutation rate. There has been a considerable amount of effort from a variety of domains in which EAs are used to find the "best" mutation rate. While there is no clear answer that is best suited across all domains, it is still valuable to have some known boundaries that provide useful results. As a result, this value provides the practitioner with the ability to influence the mutation rate around values that are known to work well for a particular domain. As the EA begins to converge, this value will strongly influence the bounds of values returned by this function. The GSD function represents the standard deviation of the gene values across the population for a given gene locus $j$. For example, values for gene locus 3 across all individuals may have a standard deviation of 2.4. This would

be the value returned by GSD($3$). However, if all individuals in the population have exactly the same value for gene locus 3, then the standard deviation would be 0. The value $Max_{GSD}$ represents the maximum standard deviation observed across all genes loci. As the EA begins to converge, some gene values across all individuals will begin to converge while others are still very much divergent. This value represents the most divergent standard deviation of the genes in the individual. In combination with the GSD value, these two variables represent the genotypic diversity by measuring if particular gene values are converging. Finally, the value FSD represents the standard deviation calculated for the fitness values. When the EA first begins, this value will often be high. As the EA progresses and begins to converge, this value will become lower. Consequently, the values returned by the mutation function will often be zero or near zero at the beginning. It is not until the EA starts to converge that FSD becomes lower, thereby increasing the values returned by the mutation function. The FSD variable represents the phenotypic diversity by measuring if the population fitness values are converging.

At a macro-level, this mutation function will begin to return larger values as the fitness values begin to converge. This convergence is measured by the standard deviation of the fitness values (FSD variable). As the fitness values begin to converge, the gene values will then begin to converge. The genotypic convergence is measured by the standard deviation of the gene values and ultimately represented by the maximum standard deviation observed ($Max_{GSD}$). When the gene values begin converging, the mutation function will begin returning even larger values. This behavior allows the EA to converge as needed, but helps keep the EA from completely converging to a single value.

## 4. TESTING

To evaluate this approach, the test objective was made to be as simple as possible: Given a large set of documents, identify a single event of interest (EOI). A simple objective was established so that the behavior of the EA could be more easily observed and understood. Future work will investigate objectives that are more complex.

For the initial testing, the single EOI was defined as an earthquake that occurred on January 26, 2001 in the city of Bhuj, India. Ten newspaper articles were defined as "evidence" of this event. These documents came from a variety of sources and describe different aspects of this event. Once these documents were identified, they were then embedded into a collection of 990 documents, bringing the total document set to 1,000. These documents were categorized into 7 different categories as shown in Table 1.

**Table 1. Test Data Category Distribution**

| Number of Documents | Category | Percentage |
|---|---|---|
| 366 | Basketball | 36.6% |
| 240 | Financial News | 24.0% |
| 162 | Biological Weapon | 16.2% |
| 98 | Soccer | 9.8% |
| 75 | Dirty Bomb | 7.5% |
| 49 | Gas Prices | 4.9% |
| 10 | Earthquake Disaster | 1.0% |

For the initial testing, the focus category is the "Earthquake Disaster" category. This category was purposely made to be the smallest category in an effort to determine how well the EA would retrieve these documents given the specified FC's.

Once the test data had been established, a set of FC's were developed as shown in Table 2.

**Table 2. Foundational Characteristics of Interest**

| Country | Natural Disaster | Infrastructure |
|---|---|---|
| Germany | Earthquake | Power plants |
| India | Flood | Roads |
| Japan | Tornado | Bridges |
| United States | Hurricane | Rail |

An event may defined by any combination of 2 or 3 of these FC's. For example, an event may be defined by {Germany, Flood, Roads} or {United States, Hurricane}, but not be defined by {Earthquake, Tornado} or {Roads, Rail} or {Japan, India}. The event must be defined by distinct FC types. As described previously, each of the FC's listed were defined by specific terms that would help clearly identify the concept for which it represented.

In comparing the FC's and the test data set, the data set contained a large variety of documents that would match well with the chosen countries. A smaller variety of documents would match well with the chosen natural disasters, and an even smaller number of documents would match with the chose infrastructures. In doing this, the goal was to observe the performance of the GA given FC's that ranged from very general (countries) to very specific (natural disasters and infrastructure). Of specific importance is to determine how the EA would filter noise from the data of interest by identifying specific connections between the FC's. Once the test data and test FC's were identified, several runs of the EA were performed using various parameter settings. Significant results of these tests are discussed in the next section.

# 5. RESULTS & CONCLUSIONS

Initial test results show that the EA successfully found the EOI and successfully described this EOI with the FC's {India, Earthquake}. The outcome of the EA was a set of 10 documents that matched exactly the 10 documents that supported the EOI. Table 3 shows the phrases that matched between the documents to help support {India, Earthquake} as a valid event description.

**Table 3. Phrase matches between documents**

| Category | Phrases |
|---|---|
| Nouns | quake victims<br>fresh tremors<br>relief operations<br>sniffer dogs<br>disaster management<br>MI-26 helicopters |
| Proper Nouns | Pakistan High Commissioner in New Delhi<br>Indian Air Force<br>Kutch district<br>Kandla port |
| Locations | Ahmedabad<br>Gujarat<br>Bhuj |
| Organizations | None found |
| People | None found |

An important outcome of this effort was the discovery that this approach can help reveal to the analyst not only the documents that would support the existence of an event, but specifically what evidence in these documents supports this proposition. In this particular test case, the EA not only revealed these 10 documents from a set of 1,000 documents, but also specific words such as "disaster management", "Kutch district", and "Bhuj". These phrases were not part of the original FC descriptions, but are phrases that may provide new insight into specific details of interest to the analyst.

In addition to this result, there were some other interesting behavioral characteristics of the EA in this approach. As with most EA applications, there are varieties of parameters that must be established such as the crossover rate, population size, etc. Actual values used for the parameters can significantly affect the behavior of the EA and the result that is achieved. In this particular work, the threshold parameter for step 2 described in section 3 proved to be yet another parameter that can significantly alter the outcome of the EA. For this particular work, a threshold value of 0.07 was used. With this value, the correct EOI was found. However, upon additional testing, it was observed that changing this value to 0.065 would dramatically affect the outcome, and another event defined by the FC's {Japan, Power Plants} was identified. This event was previously unknown in this test set, and was initially discarded as an incorrect outcome. Upon further analysis however, this second event identified by the EA exposed a set of 15 documents discussing and describing a political situation involving North Korea, Japan, and the United States that centered on the development of nuclear power plants by the North Koreans. Many of these documents described Japan's response and role in the political situation. Consequently, the EA had found another valid event. Table 4 shows the phrases that matched between the documents to help support {Japan, Power Plants} as a valid event description.

**Table 4. Phrase matches between documents**

| Category | Phrases |
|---|---|
| Nouns | economic crisis<br>grave situation<br>nuclear weapons<br>peaceful purposes<br>power plants<br>power generation<br>nuclear reactors<br>graphite-moderated reactor |
| Proper Nouns | Clinton administration<br>Bush administration<br>DPRK<br>Soviet Union<br>North Korea<br>P'yongyang<br>UN Security Council<br>NPT |
| Locations | United States<br>Japan<br>Korea<br>People's Republic of Korea |
| Organizations | None found |
| People | None found |

In comparing these two events, it was determined that the threshold value for step 2 was affecting the "coupling" of the

documents that supported the events. For the event {India, Earthquake}, the 10 documents supporting this event were strongly coupled. While each document was different and from a different source, they all described specific details of the event that provided a clear understanding of the event and the publication of these documents spanned a very short period of time (4 days). For the event {Japan, Power Plants}, the 15 documents were more loosely coupled. The publication of these documents spanned a much larger period (years 1998 to 2003), and described various aspects of a political "situation" rather than a specific "event". As a result, it was hypothesized that a higher threshold value for step 2 of the process would identify strongly coupled documents supporting a particular event, while lower threshold values would identify loosely coupled documents supporting a particular event. Additional testing supported this hypothesis. In one test, the threshold was set to 0.12 resulting in the event {India, Earthquake} being identified. However, only 2 documents of the 10 known documents were identified as supporting this event. These 2 documents were very strongly coupled. Raising the threshold value to 0.15 resulted in no events being identified. Lowering the threshold value to 0.01 resulted in an invalid event being identified. As a result of this relationship between the threshold and the document coupling, the current niching scheme could not support the identification of both of these events in a single run of the EA. Only one of the two could be identified at a given time depending on how the threshold was set. An improved approach is needed to identify both of these events at the same time.

## 6. FUTURE WORK

Future work will investigate the identification of multiple events rather than a single, dominant event. The current niching scheme does not sufficiently identify multiple events. An improved niching scheme or alternative approach is needed to identify more accurately multiple events at the same time. In addition, further work is needed to better clarify and understand the importance of document coupling in supporting the existence of an event. Currently, evaluation of document coupling characteristics is a subjective matter. A more objective evaluation would enhance automated capabilities to recognize more coherently events of interest.

## 7. REFERENCES

[1] Wikipedia: Danish cartoons, Current March 2007, http://en.wikipedia.org/wiki/Jyllands-Posten_Muhammad_cartoons_controversy

[2] Huttunen, S., Yangarber, R., and Grishman, R. "Complexity of Event Structure in IE Scenarios", in *Proc of the 19th International Conf. on Computational Linguistics*, August 2002.

[3] Turmo, J., Ageno, A., and Catala, N. "Adaptive Information Extraction", *ACM Computing Surveys*, Vol. 38, No. 2, July 2006.

[4] Patton, R.M. *Application of Intelligent Method for Improved Testing and Evaluation of Simulation Systems Software.* Ph.D. Thesis, University of Central Florida, Orlando, FL, 2002.

[5] Wikipedia: Madrid train bombings, Current March 2007, http://en.wikipedia.org/wiki/11_March_2004_Madrid_attacks

[6] Bertaux, J.L., et. al. "Monitoring solar activity on the far side of the Sun from sky reflected Lyman alpha radiation", *Geophysical Research Letters*, Vol. 27, No. 9, pages 1331-1334, May 2000.

[7] Marcy, G.W. and Butler, R. P., "Detection of Extrasolar Giant Planets", *Annual Review of Astronomy and Astrophysics* Vol 36: 57-97, September 1998.

[8] Yang, Y., Pierce, T., and Carbonell, J. "A Study on Retrospective and On-Line Event Detection", in *Proc. of the 21st annual international ACM SIGIR*, August 1998.

[9] Allan, J., Papka, R., and Lavrenko, V. "On-line New Event Detection and Tracking", in *Proc. of the 21st annual international ACM SIGIR*, August 1998.

[10] Yang, Y., Ault, T., Pierce, T., and Lattimer, C. "Improving text categorization methods for event tracking", in *Proc. of the 23rd annual international ACM SIGIR*, July 2000.

[11] Salton, G., Wong, A., and Yang, C. S., "A Vector Space Model for Automatic Indexing", *Communications of the ACM,* Vol. 18, No. 11, pages 613–620, 1975.

[12] Reed, J., et al. "TF-ICF: A New Term Weighting Scheme for Clustering Dynamic Data Streams" in *Proc. of the 5th International Conference on Machine Learning and Applications (ICMLA'06)*. Orlando, FL., 2006.

[13] LingPipe, Current March 2007, http://www.alias-i.com/lingpipe/

[14] Stanford Log-linear Part-Of-Speech Tagger, Current March 2007, http://nlp.stanford.edu/software/tagger.shtml

[15] Dasgupta, D. and McGregor, D.R. "Species adaptation to nonstationary environments: A structured genetic algorithm" *Presented at Artificial Life III workshop*, Santa Fe, New Mexico, June 1992.

[16] Dasgupta, D. and McGregor, D.R. "Nonstationary Function Optimization using the Structured Genetic Algorithm" *In Proc. of Parallel Problem Solving From Nature Conference*, Brussels, Belgium, September 1992.

[17] Mahfoud, S.W. *Niching Methods for Genetic Algorithms.* Masters Thesis, Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, Illinois, 1995.