

# Improving Knowledge Discovery in Document Collections through Combining Text Retrieval and Link Analysis Techniques

Wei Jin<sup>1</sup>, Rohini K. Srihari<sup>1</sup>, Hung Hay Ho<sup>1</sup>, and Xin Wu<sup>2</sup>

<sup>1</sup>*Department of Computer Science & Engineering, University at buffalo, State University of New York, USA*

<sup>2</sup>*Department of Computer Science and Technology, University of Science and Technology of China, China*

*wjin2@buffalo.edu, rohini@cedar.buffalo.edu, hunggho@buffalo.edu, xinwu@mail.ustc.com.cn*

## Abstract

*In this paper, we present Concept Chain Queries (CCQ), a special case of text mining in document collections focusing on detecting links between two topics across text documents. We interpret such a query as finding the most meaningful evidence trails across documents that connect these two topics. We propose to use link-analysis techniques over the extracted features provided by Information Extraction Engine for finding new knowledge. A graphical text representation and mining model is proposed which combines information retrieval, association mining and link analysis techniques. We present experiments on different datasets that demonstrate the effectiveness of our algorithm. Specifically, the algorithm generates ranked concept chains and evidence trails where the key terms representing significant relationships between topics are ranked high<sup>1</sup>.*

## 1. Introduction

It is recognized that text information is growing at an astounding pace. These vast collections of publications offer an excellent opportunity for text mining, i.e., the automatic discovery of knowledge. In this paper we propose to combine information retrieval and link-analysis techniques over the extracted features provided by Information Extraction (IE) engine for finding new knowledge.

It is often the case that a document collection reveals interesting information other than what is explicitly stated, thus, the main theme of our paper is

based on the hypotheses that “The wealth of recorded knowledge is greater than the sum of its parts” [1], and interesting links and hidden information that connect facts, propositions or hypotheses can be formed by using some techniques to discover previously unknown logic connections among the existing information we have.

The goal of this paper is to sift through these extensive document collections and find such links. We present Concept Chain Query (CCQ), a special case of text mining focusing on handling questions like: “What is the relation of topics X and Y?” (e.g. two individuals). A traditional search involving, for example, two person names will attempt to find documents mentioning both of these individuals. Failing this, the search results in documents containing one of the names. This research focuses on a different interpretation of such a query: what is the best evidence trail across documents that connect these two topics assuming that one or more instances of both topics occur in the corpus, but not necessarily in the same document? We go one step further and require the response to be a set of text snippets extracted from multiple documents. This is in fact a cross-document summary of the plausible relationship between the two topics.

Figure 1 below shows an example for the concept chain query in the anti-terror domain. The snippets of text can lead to a plausible connection between *Nashiri* and *Nairobi attack* through concept *Jihad Mohammad Ali al Makki*. Note that the connection may not be contained in a single document, but could be found through the discovery of knowledge holding across documents.

The remainder of the paper is organized as follows. Section 2 discusses concept chain queries in more detail. Section 3 discusses related work. Section 4 presents the graphical text representation and mining

---

<sup>1</sup> This work is sponsored by NSF grant IIS-0325404 and a research grant from FAA 032-G-009.

model and section 5 described the algorithm for evidence trail generation. Section 6 discusses experiments based on processing two corpora from different domains, and is followed by conclusions.

Doc 1: *Nashiri* and *his cousin, Jihad Mohammad Ali al Makki*, returned to Afghanistan, probably in 1997, Nashiri again encountered Bin Ladin, still recruiting for "the coming battle with the United States." Nashiri joined al Qaeda and later was recognized as the chief of al Qaeda operations in and around the Arabian Peninsula.

Doc 2: In late 1998, al Qaeda decided mounting an attack against a U.S. vessel and *Jihad Mohammad Ali al Makki*, also known as *Azzam*, was a suicide bomber for the *Nairobi attack*.

**Figure 1. Concept chain and evidence trail**

## 2. Concept Chain Queries

Concept chain queries go beyond what traditional IR models handle [13]. A concept chain query involving topics A and B tries to find if (i) there is a direct connection (association) between them (e.g. co-occurrence in sentences), or (ii) if they can be connected by several intermediate concepts (paths). In the second case, the answer may not be contained in any particular individual document, but is the result of relating the contents of a small set of documents.

**Table 1. Comparison of concept chain queries with IR**

	IR	Concept Chain Queries
Query	keywords	topic pairs
Output	a ranked list of documents	ranked sets of concept chains and evidence trails
Ranking	(i) single document satisfies query (ii) documents ranked independently w.r.t. query	(i) single document may not satisfy query (ii) connection between documents critical

This problem can be thought of as generalizing the information retrieval (IR) task as shown in Table 1 which compares traditional IR to concept chain

queries. In IR systems, the input is a set of keywords, whereas here the query is for concept chains. In IR systems the output is a ranked list of documents, where each document is ranked independently w.r.t. the query. In concept chain queries, the output is ranked sets of concept chains and evidence trails (i.e. text snippets extracted from multiple documents) capturing the relationship between topics. The problem manifests itself in different IR applications including complex question-answering and summarization.

## 3. Related work

Much of the work in hypotheses generation in document collections makes use of an idea originated by Swanson in the 1980s called the "complementary structures in disjointed literatures" (CSD). Swanson and Smalheiser [4] found in their research that new interesting and unknown implicit information could be discovered if we study the linkage between textual records. They proposed a simple "A influences B, and B influences C, therefore A may influence C" model for detecting instances of CSD that is commonly referred as Swanson's ABC model. Using this technique, he found a connection implying patient benefit between fish oil and Raynaud's syndrome, two years before clinical trials established that the benefit was real. Small and Garfield [10] suggested a wider citation search and proposed that a path of several co-citation steps can connect two clusters of documents. Therefore, it will be meaningful to study the multistage paths between topics of interest when we are forming hypotheses.

CARTwheels [11] and Arrowsmith [12] are two knowledge discovery systems that address the similar problem of finding connections. CARTwheels builds connections among topics by "story-telling". Arrowsmith works at discovering connections by automatically using the content in a collection of scientific papers. It detects indirect relationships between topics in MEDLINE by finding common keywords between two document sets. More recently, Srinivasan [7] has used her approach to demonstrate the feasibility of her approach in text mining based on MeSH terms and UMLS semantic types and presented open and closed text mining algorithms that are built within the discovery framework established by Swanson and Smalheiser in biomedical domain. Our competitive baseline model is based on Srinivasan's closed text mining algorithm.

There has also been work on discovering connections between concepts across documents using social network graphs, where nodes represent

documents and links represent connections between documents. However much of the work has focused on special problems on social network analysis, such as detecting communities [8]. [9] is the work which is close to the problem being presented here. Detecting associations between people is modeled as finding a connection subgraph and a solution based on electricity analogues is presented. Here we pointed out several differences between our proposed model and the subgraph model [9]. The most notable difference is the reliance on URL links to establish connections between documents. The corpus we addressed is a general plain text corpus, and thus the associations are extracted based on content analysis. Second, our approach presents the paths individually instead of a connection subgraph including all paths together. This will allow greater user input in determining the best paths, including recency, novelty, semantic coherence, etc. Third, the approach presented here attempts to generate the evidence trail for each retrieved chain, whereas the connection subgraph approach does not. Finally, the connection subgraph solution only addresses named entities whereas our approach extends to general concepts.

Recently, [3] presents a model of using Stepping Stones and Pathways (SSP) to connect topics in document collections. Query results are networks of document groups representing topics, each group relating to and connecting (by documents) to other groups in the network. However, the starting point of this approach is to create an intermediate document set with the documents that are common to both document sets retrieved for the two topics respectively. Thus, if the intersection of both document sets is empty, the system could not find any connection, whereas our proposed solution handles the query without such a limitation and we assume that one or more instances of both topics occur in the corpus, but not necessarily in the same document. Second, we present the user an explanation of the retrieved chains, i.e. a cross-document summary of the plausible relationship between topics instead of document groups for each connection used in [3]. This may assist a user with the second dimension of the analysis process, i.e., when the user has to peruse the documents to figure out the nature of the relationship underlying a suggested chain. Third, as the authors also pointed out, the main disadvantage of the algorithm proposed in [3] is when adding a new stepping stone, there is no check for global coherence in the pathway. Since each update only affects the adjacent nodes, it is possible that while all local connections are valid, the meaning of all connections taken together will be difficult to understand. In our proposed approach, coherence is the

primary factor to be considered when formulating the algorithms of generating and ranking chains and evidence trails.

The main contributions of this paper are as follows:

- The goal of this work is to assist the user in identifying links that hold across documents, which is crucial for hypothesis generation. We assume that the knowledge we desire may not be contained in a single document, but can be found through the discovery of knowledge holding across multiple documents.
- We extended Swanson's ABC model [4, 7] to generate concept chains. This was achieved by considering multiple levels of interesting concepts instead of just one level as in the original technique. Specifically, our algorithm generates ranked concept chains of various lengths where the key terms representing significant relationships between topics are ranked high.
- The approach presented here attempts to generate an explanation of the chain through linking text snippets extracted from multiple documents with the relevant chain. This is in fact a cross-document summary of the plausible relationship between topics.
- Users can easily guide the discovery process through setting appropriate parameters. The generated results are easy to understand and interpret by the user.

## 4. A concept-association-graph based approach

Different from previous approaches [4, 7] that do not represent the associations as graph, the solution being pursued here takes the view that the concepts and associations in a particular domain can be represented as a network (i.e. concept association graph (*CAG*) constructed below) with the nodes representing concepts and edges between them representing associations. Then based on *CAG*, both the existing explicit (direct) associations and also potential new associations (hypotheses) could be discovered.

### 4.1. Data preparation

An Information Extraction (IE) engine, Semantex [2] was used to extract the features from document collections. Semantex tags named entities, common relationships associated with person and organization,

as well as providing subject-verb-object (SVO) relationships. All named entities, as well as any noun or noun phrases participating in SVO relationships were extracted as concepts.

After filtering phase, the left concepts were then mapped to their respective ontology categories in an automated fashion. Table 2 illustrates some of the concepts that are extracted by the system; Concepts in the table such as *visa applicant*, *CIA report* indicate the data-driven manner in which these concepts are extracted.

**Table 2. Portion of terrorism ontology**

Ontology Node	Examples (Instances)
Organization	Counterterrorism center Aviation unit
Human Action	Assignment, assumption, attack
City	Manila, Nairobi, Rome
Person	Visa applicant, Arab volunteer
Document	CIA report, conspirator statement
Country	Albania, Sudan, Afghanistan

#### 4.2. Concept association graph construction

We build the concept association graph (*CAG*) by first identifying frequent itemsets in the text corpus. Discovery of frequent itemsets is a well-understood data mining problem, arising in the *market basket* association rule problems [6]. To achieve the fine granularity required by concept chain queries, each passage (e.g. sentence in our experiment) is viewed as a market basket of concepts. Observe that two concepts which occur in the same sentence seem more correlated than two terms which occur far apart in a document, it would be worthwhile to factor in the distance between two concepts in mining their associations.

All the resulting frequent items are represented as vertices in the weighted graph  $CAG = \langle V, E \rangle$ . The adjacency relationship between vertices is decided by the obtained frequent 2-itemsets, i.e., there exists an edge between  $v_i \in V$  and  $v_j \in V$  ( $i \neq j$ ) if  $(v_i, v_j)$  constitutes a frequent 2-itemset.

The weight of the edge connecting  $v_i$  and  $v_j$  represents the strength of the relationship between them. The traditional association rule confidence may

not be an appropriate measure of correlation strength: Confidence overemphasizes common items as consequents and rare items as antecedents. The consequent in such cases rarely adds much meaning and interest to the user. In our case, for example, concept chains going through concepts that are very generic, e.g., United States, would not be considered as interesting as a chain that goes through a vertex labeled *bomb shelter*.

Instead of using confidence, we used mutual information [12] as the definition of edge weights in the *CAG*:

$$\log \frac{P(x, y)}{P(x)P(y)} \quad (1)$$

This measure computes the ratio of the actual probability of a frequent 2-itemset occurring in a sentence to the probability of the items occurring together by chance. This measure also emphasizes relatively rare items that generally occur together and deemphasizes common items, and hence leads to more interesting chains in our case.

In terms of text mining and knowledge discovery, the representation also satisfies the following properties:

- **Correlation to documents:** The representation relates concepts and associations to the underlying text snippets in documents. Without this support, the *CAG* is not complete since users need to validate conclusions by looking at actual documents. For this purpose, we added additional information in the *CAG* called **Sentence Clusters:** This data structure links the sentences in the document collection with the relevant association information represented in the *CAG*. For example, the edge between  $v_i$  and  $v_j$  is related to a sentence cluster composed of sentences in which  $v_i$  and  $v_j$  co-occur. The resulting sentence clusters offer a view of the document collection which is highly characterized by the presence of concept associations and play an important role in generating a cross-document summary as an explanation of the relationship between topics, which we formalize in section 4.
- **Path finding:** The representation supports the ability to generate the most likely paths connecting two topics of interest, emanating from a given topic etc.
- **Quantification of information:** The representation enables the generation of a

measure reflecting the amount of information in a particular document, or path.

- **Interactive Search:** The user can intervene, guide the discovery process and the interactive interface is provided.

### 4.3. Concept chain queries

This section presents the formulation and solution of concept chain queries (*CCQ*) based on the constructed *CAG*. Given two topics of interest, *CCQ* tries to find if (i) there is a direct connection (association) between them, or (ii) if they can be connected by several intermediate concepts (paths). Note that finding direct links between two topics is trivial, in the following we mainly focus on discovering and ranking indirect connections between topics. We consider this problem as finding optimized transitive associations at various lengths between topics in the *CAG*. A transitive association has the form:

$$c_1 \rightarrow c_2 \rightarrow \dots \rightarrow c_n \text{ (transitive strength)}$$

Where  $c_1$  and  $c_n$  are the source topic and destination topic respectively.

Considering the user's need, the output takes the following form: the retrieved chains are first grouped by *association degree*, which is indicated by the length of the links, and then ranked by *association strength* (estimated potential) within each group. In other words, we would like to automatically return ranked lists of concept chains, one list for each length between the shortest connection length and the specified maximum length *max\_len*. For the chains of the same length, they will be ranked based on the 'goodness' function  $g$  defined below. Users can easily guide this discovery process through setting appropriate parameters such as the desired maximum length of the chains and the number of top ranked chains to be returned for each possible length. The generated result offers a view of top  $K$  connections at various lengths connecting topics.

The function  $g$  represents the "goodness" of the solution concept chain *CC*. The definition of  $g$  is to produce concept chains that capture salient aspects of relationships between topics in document collections. For that purpose, we tried a variety of ways as the definition of the goodness function  $g(CC)$ . Finally, through experiments, we find directly using the product of the weights of its constituent edges as the

goodness function  $g$  in the *CAG* achieved the best evaluation result.

Formally, given the source topic  $c_1$  and destination topic  $c_n$ , a concept chain  $CC_{1 \rightarrow n}$  is basically a sequence of concepts  $c_1 \rightarrow c_2 \rightarrow c_3, \dots, c_{n-1} \rightarrow c_n$ , the transitive strength of a path from  $c_1$  to  $c_n$  made up of the links  $\{(c_1, c_2), \dots, (c_{n-1}, c_n)\}$ , denoted by  $g(c_1, c_n)$ , is given by the following equation:

$$g(c_1, c_n) = \prod_{i=1}^{n-1} (w(c_i, c_{i+1})) \quad (2)$$

where  $w(c_i, c_{i+1})$  is the weight of the edge between  $c_i$  and  $c_{i+1}$  in the *CAG*. We call a path an optimal path from  $c_1$  to  $c_n$  with path length  $l$  if its transitive strength is maximal among all the possible paths with length  $l$ .

### 4.4. Generating and ranking concept chains

This section presents the algorithm of generating and ranking concept chains. First, the *CAG* may undergo the optional cleaning phase depending on the user's need. We use both *support* and *mutual information* (edge weight) as filters. Very high support itemsets are almost always significant, as are high mutual information itemsets. The size of the desired *CAG* can be adjusted using the parameters *support* and *edge confidence* (which will filter the edges whose weights are below a user specified threshold).

Then our optimization problem is now to find a chain that maximizes our measure for each possible length. For this problem, the traditional dynamic programming algorithm could solve it in  $O(ek)$  time and  $O(vk)$  space, where  $v$  is the number of vertices in the input graph;  $e$  is the number of edges and  $k$  is the maximum length of any allowed  $c_1$ - $c_n$  path.

However, there are often cases that the user might be interested in exploring more potential chains instead of just one optimal chain. For this purpose, we have adapted the traditional dynamic programming algorithm into finding  $T$ -best chains for any given length. The following search strategy is employed considering the efficiency of finding  $T$ -best paths from the starting point  $s$  to ending point  $e$  with maximum path length  $L$ . For every intermediate concept  $v$  in the candidate path, only  $T$ -best paths are kept from  $s$  to  $v$ , and then staring at  $v$ , proceed to the next length level, i.e. explore all the neighboring vertices of  $v$ , and keep applying the same principle for the neighbors of  $v$  until reaching the ending topic  $e$  or surpassing the length  $L$ . The correctness of this algorithm can be briefly described as follows: let  $P_{s,e}$  denotes the path from  $s$  to  $e$ ,  $O_{s,v}(t)$  represents the  $T$ -best paths from  $s$  to  $v$  and  $O_{s,e,v}(t)$  denotes the  $T$ -best paths from  $s$  to  $e$  through the

intermediate vertex  $v$ , then for any intermediate vertex  $v$ , we show that the initial segment  $P_{s,v}$  in  $O_{s,e,v}(t)$  must be contained in  $O_{s,v}(t)$ .

## 5. Evidence trail generation

This stage is to collect evidence that supports the generated chain in the previous step. Evidence trails may assist a user with the second dimension of the analysis process, i.e. when the user has to peruse the documents to figure out the nature of the relationship underlying a suggested potential chain. We go one step further and require the response to be an optimal set of sentences extracted from multiple documents in which the concept chain occurs. This is in fact a cross-document summary of the plausible relationship between topics.

We consider evidence trail generation as chain-focused sentence retrieval from the document collection. Recall that given the source topic  $c_i$  and destination topic  $c_j$ , a retrieved concept chain  $cc_{i \rightarrow j}$  is in the form of  $c_i \rightarrow c_{i+1} \rightarrow, \dots, c_{j-1} \rightarrow c_j$ . As mentioned earlier, one sentence cluster is linked to each individual segment of the chain (e.g.  $c_i \rightarrow c_{i+1}$ ) during the phase of the *CAG* construction.

Our goal is to build a chain-focused sentence retrieval mechanism and the output, an optimal set of sentences relevant to the corresponding chain, is returned to the user as a chain-focused summary. This is accomplished by selecting one sentence that has the highest importance score (to be described shortly) from each sentence cluster and the order of sentences in the summary follows the order of associations appearing in the corresponding chain.

The algorithm for choosing the most relevant sentence from each sentence cluster w.r.t. the chain is motivated by [5]. [5] introduced a random walk model for finding answers to complex questions. This model is based on the idea that answers can be found by scoring each sentence against a complex question and the sentence rank is produced by a mixture model that combines an estimation of a sentence’s relevance to a question as well as the similarity between sentences. We adapted this model into our needs, i.e., given each sentence cluster and the corresponding chain, how can we select the most representative sentence from each cluster such that the generated summary best explains the corresponding chain?

### 5.1. Relevance to the chain

To measure the relevance of a sentence in the cluster w.r.t. the query (i.e. the concept chain), we first compute concept IDF’s by the following formula:

$$idf_c = \log\left(\frac{N}{s_c}\right) \quad (3)$$

Where  $N$  is the total number of sentences in the cluster, and  $s_c$  is the number of sentences that the concept  $c$  appears in.

Then the relevance of a sentence  $s$  to the query  $q$  is calculated by:

$$rel(s, q) = \sum_{c \in q} \log(f_{c,s} + 1) \times idf_c \quad (4)$$

Where  $f_{c,s}$  is the frequency of concept  $c$  appears in  $s$ . Note that we use the whole concept chain as the query instead of the corresponding segment of the chain when measuring the relevance of a sentence in the cluster, which aims to estimate the potential contribution of a sentence as an indicator of the whole chain. From experimental results, this measure is shown to be very successful in selecting sentences that maximize coherence and coverage for inclusion in the final summary w.r.t. the query.

### 5.2. The mixture model

As described in [5], we used a mixture model to capture any inter-sentence information in a cluster.  $val(s, q)$ , the score of a sentence  $s$  w.r.t. query  $q$ , is computed as the sum of its relevance to the query itself and the similarity to the other sentences in the cluster:

$$val(s, q) = d \frac{rel(s, q)}{\sum_{t \in C} rel(t, q)} + (1-d) \sum_{u \in C} \frac{sim(s, u)}{\sum_{t \in C} sim(t, u)} val(u, q) \quad (5)$$

Where  $C$  is the set of all sentences in the cluster. The value of  $d$ , which we call the “chain bias”, is a parameter to adjust the weights of the two terms in the equation. Higher  $d$  gives more importance to the relevance to the chain than the sentence similarity in the cluster. The widely used cosine measure was adopted as the similarity between two sentences in a cluster:

$$sim(u, v) = \frac{\sum_{c \in u, v} f_{c,u} \times f_{c,v} \times (idf_c)^2}{\sqrt{\sum_{u_i \in u} (f_{u_i, u} idf_{u_i})^2} \cdot \sqrt{\sum_{v_j \in v} (f_{v_j, v} idf_{v_j})^2}} \quad (6)$$

Next, the  $val$  vector can be written as  $Val = [dX + (1-d)Y]^T \times Val$ . A stochastic matrix  $X$  is considered such that all the elements in the  $i^{th}$  column are

proportional to  $\text{rel}(i, q)$ . In order for matrix  $\mathbf{X}$  to be stochastic, normalize it so that row sums add up to 1. Similarly, a stochastic matrix  $\mathbf{Y}$  is defined such that  $\mathbf{Y}(i, j)$  is proportional to  $\text{sim}(i, j)$ . The square matrix  $\mathbf{E} = [d\mathbf{X} + (1-d)\mathbf{Y}]$  is also stochastic and defines a Markov chain where each entry  $\mathbf{E}(i, j)$  specifies the transition probability from state  $i$  to state  $j$  in the corresponding Markov chain. The  $\text{val}$  vector is the stationary distribution of the Markov chain. With probability  $d$ , a transition is made from the current state (sentence) to the states that are similar to the query. With probability  $(1-d)$ , a transition is made to the states that are similar to the current state lexically.

Then the summary of the plausible relationship underlying a suggested chain can be easily generated. For each sentence cluster corresponding to each individual segment of the chain, we select the sentence that has the highest  $\text{val}(s, q)$  score for inclusion in the final summary. The order of selected sentences in the summary follows the order of associations appearing in the corresponding chain.

## 6. Experiments

In this section we evaluate the performance of the various components in our system. Our experiments were designed to answer the following questions:

- How good is the proposed concept chain query in helping the user identify the potential relationship between two topics of interest?
- Do our generation and ranking algorithms capture salient aspects of relationships between topics?
- How well does the evidence trail generation algorithm perform? Does the generated evidence trail make sense?

### 6.1. Experimental corpus

Unfortunately no standard data sets are available for quantitatively evaluating concept chains. Thus, a major task in evaluation was constructing an evaluation data set. For the experiments we used two corpora from different domains: one is the 9/11 counterterrorism corpus. This involves processing a large open source document collection pertaining to the 9/11 attack, including the publicly available 9/11 commission report; the other dataset is provided by the U.S. National Transportation Safety Board (NTSB). The narratives tables of this dataset including the aviation accident reports coming from the year 2001 to 2003 were used. All the documents were preprocessed by the Semantex information extraction engine [2], and

the concepts are extracted and mapped to the domain ontology as described in section 4. A thesaurus file is used for matching synonyms and acronyms in the texts.

### 6.2. Goodness of chains – Case studies

Figure 2 shows the output by our system for two representative test cases. Figure 2(a) shows the generated concept chains in the 9-11 corpus linking two named entities ‘Bush’ and ‘Bin Ladin’. In this example there are strong links, at various lengths, as expected.

```
<Paths      first-word="bush"      second-word=
"bin_ladin">
<Paths length=1>
-<Path id="1" value="bush bin_ladin">
<Paths length=2>
-<Path id="2" value="bush attack bin_ladin">
-<Path id="3" value="bush al-qaeda bin_ladin">
-<Path id="4" value="bush Afghanistan bin_ladin">
<Paths length=3>
-<Path id="5" value = "bush Afghanistan al-qaeda
bin_ladin">
-<Path id="6" value = "bush administration attack
bin_ladin"> .....
```

**Figure 2(a). Concept chains connecting ‘Bush’ and ‘Bin Ladin’ ( $\text{chain\_len}=3$ ;  $\text{chain\_num}=3$ )**

```
<Paths      first-word="landing"    second-word=
"accident">
<Paths length=1>
-<Path id="1" value="landing accident">
<Paths length=2>
-<Path id="2" value="landing full_stop_landing
accident">
-<Path id="3" value="landing flight_control
accident">
-<Path id="4" value="landing landing_gear
accident">
<Paths length=3>
-<Path id="5" value="landing airport gust accident"
>
-<Path id="6" value="landing snow_covered field
accident"> .....
```

**Figure 2(b). Concept chains connecting ‘Landing’ and ‘Accident’ ( $\text{chain\_len}=3$ ;  $\text{chain\_num}=3$ )**

Figure 2(b) above shows another application of concept chain queries in the NTSB dataset. Here the purpose of concept chain queries is to discover trends and patterns that might be useful for improving aviation safety and preventing future accidents. Suppose the user is interested in discovering potential safety hazards during the LANDING phase of flight, then the query may be issued using topics ‘Landing’ and ‘Accident’ and the generated output is shown. Different from the 9-11 corpus including many named entities, in the NTSB dataset, the majority of concepts appearing in the corpus are general concepts, and thus the concept chain query goes beyond named entities and establishes correlations between named entities and/or general concepts across documents.

Note that in both examples the system outputs ranked concept chains where the key terms representing strong and informative relationships between topics are ranked high. In the following sections we will describe an evaluation set we synthesized and the experimental results are presented.

### 6.3. Evaluation set

The previous section shows the proposed approach handles the queries with a high precision at various lengths. Now we will evaluate the recall of the concept chains that the system generated. We synthesized an evaluation set as follows. We selected chains of lengths ranging from 1 to 4 in terms of number of associations. The chains were selected as follows:

- We chose various pairs of topics: in the counterterrorism corpus, the topics were mostly named entities, that is, the end points of the chains were named entities (although intervening concepts were not required to be named entities); In the NTSB dataset, the topics were mostly general concepts, such as ‘landing roll’ and ‘on-ground collision’.
- For each topic pair, the relevant paragraphs for either topic respectively were then manually inspected: we selected those where there was a logical connection between the two topics.
- After achieving agreement among all annotators, we then generated the concept chains for these topic pairs (and paragraphs) as evaluation data.

The criteria used to select the target chains and evaluate the retrieved chains are as follows:

- Cause-Effect link: this type of link is mostly found in the NTSB dataset, e.g., *loss of*

*engine power* may cause a *forced-landing* event.

- Entity-to-Entity link: correlation between entities;
- Entity-to-Event link: e.g., co-participation in one event;
- Entity-to-Person link: e.g., two persons linked by the same entity *al Qaeda*;
- Co index (identity): refers to the alias of person name. For example, *Jihad Mohammad Ali al Makki* and *Azzam* refer to the same person as shown in Fig. 1;
- Co reference: If person A lives in a house, and person B lives in a house, they are related by the concept house, this is an example where co reference produces noisy results: since the house in instance A is not the same as the house in instance B. However there are often cases when the co reference criterion does not need to be strictly adhered to. In our experiment, we relaxed this condition, for example, two individuals may share the fact that they are both involved with educational institutions, even though not the same institution.

### 6.4. Evaluation result

The above process resulted in approximately 50 query pairs in each dataset. As a baseline model, we have implemented a text mining model based on traditional information retrieval techniques. To the best of our knowledge, we are not aware of any previous work that has proposed a general model of discovering associations between topics flexibly at multiple levels, which is crucial for hypothesis generation. The baseline algorithm was based on the ‘closed’ discovery algorithm proposed in [7] but because of different corpora’ properties, we adapted this algorithm to meet our needs using sentence level co-occurrence, which achieved fine granularity required by concept chain queries in our document collections.

The experiment processed all the query pairs and for each query pair, top 5 chains for each length between the shortest connection length and the maximum allowable length (i.e. 4) were returned. The baseline model only returned the chains of length 1 and 2 for each query pair. The evaluation looked at (i) whether the target chains were found among the top 5 choices by the retrieval models; (ii) the average rank of the target chains, and (iii) the average recall achieved by both models at various lengths. The following figures present the comparative results we obtained on

executing concept chain queries from the evaluation set.

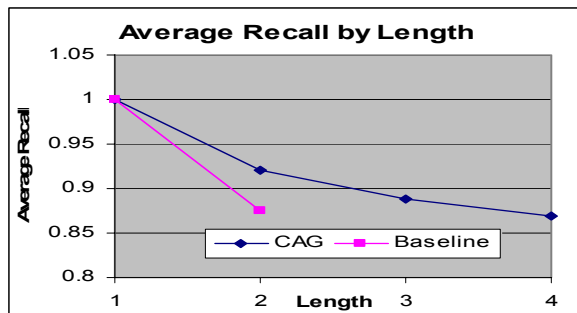


Figure 3(a). Average recall by length (CAG: Concept Association Graph)

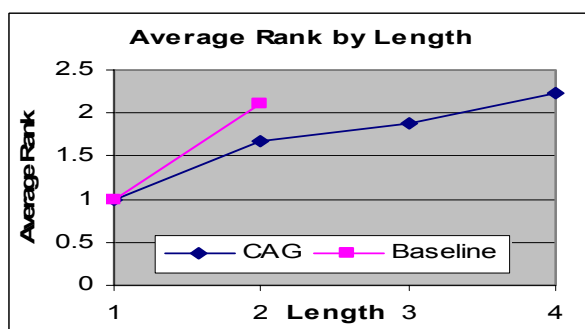


Figure 3(b). Average rank by length (CAG: Concept Association Graph)

As a post analysis, the proposed *CAG* model handles a majority of the queries with a considerably high precision at various lengths. Through inspecting the XML documents generated for each query pair, we got good coverage of the links we were looking for. It also outperforms the baseline model, in terms of average rank of the detected target chains, thereby manifesting the goodness of our chain generation and ranking algorithms, and the model performs consistently well in the NTSB corpus.

For the missed truth chains, the reason could be that the model didn't consider anaphors of person's names. The anaphora "He" or "She" and the alias of person name refer to the person we are interested in but the process couldn't find it. Another reason could be that due to stringent association extraction (sentence co-occurrence links), some of these links may have been missed. We are looking into these issues further.

### 6.5. Evidence trail evaluation

This section describes how our system performs in evidence trail generation. The goal is to see whether

the generated evidence trails explain the relationship underlying a suggested chain between topics. Coherence is the main criterion that needs to be satisfied. It is not sufficient for individual links between concepts to make sense: the chain (and evidence trails) as a whole must make sense also. Note that in formulating the algorithm for evidence trail generation, coherence is considered by selecting the sentence from each cluster which is highly relevant to the whole chain instead of each individual segment of the chain, which attempts to maximize contextual similarity between the whole chain and the final summary. We used a high chain bias  $d=0.85$  (nearly choices make little difference) and as expected, it achieved a satisfying result in terms of coherence and coverage. An example of such a scenario is illustrated in Table 3 below, which shows the top 3 sentences from the cluster corresponding to the partial chain 'al-qaeda → attack' w.r.t. the whole chain 'bush → Afghanistan → al-qaeda → attack → bin\_ladin'. Note that the sentences in a cluster which are more relevant to the whole chain would be favored.

Table 3. Top ranked sentences from the cluster corresponding to 'al qaeda → attack'

Rank	Sentence Index (docID, senID)	Sentence
1	docID='10.2 .txt' senID='84'	President Bush instructed Rumsfeld and Shelton to develop further the Camp David military plan to attack the Taliban and al Qaeda if the Taliban rejected the ultimatum.
2	docID='SS1 5.2.txt' senID='65'	What is clear is that these plots were major benchmarks in the evolving Islamist threat to the United States and foreshadowed later attacks that were indisputably carried out by al Qaeda under Bin Ladin's direction.
3	docID='10.3 .txt' senID='85'	In Phase Two, air strikes and Special Operations attacks would hit key al Qaeda and Taliban targets.

The final summary of the whole chain is presented below. It can be seen that the evidence trails extracted

from multiple documents identified the specific nature of the relationship underlying the suggested chain.

<docID='10.2.txt' senID='82' chain='bush→Afghanistan'> Bush and his advisers discussed new legal authorities for covert action in Afghanistan, including the administration's first Memorandum of Notification on Bin Ladin.

<docID='6.4.txt' senID='135' chain='Afghanistan→al qaeda'> CIA briefing slides described al Qaeda as the "most dangerous group we face," citing its "leadership, experience, resources, safe haven in Afghanistan, [and] focus on attacking US.

<docID='10.2.txt' senID='84' chain='al qaeda → attack'> President Bush instructed Rumsfeld and Shelton to develop further the Camp David military plan to attack the Taliban and al Qaeda if the Taliban rejected the ultimatum.

<docID='SS15.2.txt' senID='65' chain=' attack → bin ladin'>

What is clear is that these plots were major benchmarks in the evolving Islamist threat to the United States and foreshadowed later attacks that were indisputably carried out by al Qaeda under Bin Ladin's direction.

**Figure 4. Evidence trails of the concept chain 'bush → Afghanistan → al qaeda → attack → bin ladin'**

## 7. Conclusions

This paper focuses on detecting links between two topics of interest. Different from traditional search, we interpret such a query as finding the most meaningful evidence trails across documents that best explain the relationships between topics. A graphical text representation and mining model is proposed which combines information retrieval, association mining and link analysis techniques. In addition to generating concept chains connecting topics, the evidence trail extracted from multiple documents as an explanation of the suggested chain is also provided, which may assist users in further investigations.

In future research in addition to the specific points raised in the paper, we are also researching extensions of concept chains to concept graph queries. This will enable users to quickly generate hypotheses graphs which are specific to a corpus. These matched

instances can then be used to look for other, similar scenarios.

## 8. References

- [1] Davies, R. "The creation of new knowledge by information retrieval and classification." *Journal of Documentation*, Location 45(4), 1989, pp. 271-301.
- [2] Srihari, R. K., Li, W., Niu, C. and Cornell, T. "InfoXtract: A Customizable Intermediate Level Information Extraction Engine." *Journal of Natural Language Engineering*, 2006, pp. 1-26.
- [3] Das-Neves, F., Fox, E. A. and Yu, X. Connecting Topics in Document Collections with Stepping Stones and Pathways. *CIKM'05*, ACM Press, 2005, pp. 91-98.
- [4] Swanson, D. R. "Complementary Structures in Disjoint Science Literatures". In *proceedings of the 14<sup>th</sup> ACM SIGIR Conference*, ACM Press, Chicago, IL, 1991, pp 280-289.
- [5] Otterbacher, J., Erkan, G. and Radev, D. "Using random walks for question-focused sentence retrieval." In *Proceedings of HLT/EMNLP*, 2005, pp. 915-922.
- [6] Agrawal, R., Mannila, H., Srikant, R., Toivonen, H. and Verkamo, A. I. "Fast discovery of association rules." *Advance in knowledge discovery and data mining*, AAAI Press/The MIT Press, pp. 307-328, 2004.
- [7] Srinivasan, P. "Text Mining: Generating Hypothesis from Medline." *JASIST*, Vol.55, pp. 396-413 (2004).
- [8] Gibson, D., Kleinberg, J. and Raghavan, P. Inferring web communities from link topology. In *Proc. of the Ninth ACM Conference on Hypertext and Hypermedia*, pp. 225-234, 1998.
- [9] Faloutsos, C., McCurley, K. S. and Tomkins, A. Fast discovery of connection subgraphs. In *Proceeding of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'04)*, pp.118-127, Seattle, US, 2004.
- [10] Small, H. and Garfield, E. Verification of results that logically related non-interactive literatures are potential source of new knowledge. *Journal of the American Society for Information Science* 40 (3), pp. 152, 1989.
- [11] Ramakrishan, N., Kumar, D., Mishra, B., et al. "Turning CARTwheels: An alternating algorithm for mining redescrptions." *KDD'04*, pp. 266- 275, 2004.
- [12] Swanson, D., Smalheiser, N. and Bookstein, A. "Information discovery from complementary literatures: categorizing viruses as potential weapons." *JASIST* 52(10), pp. 797-812, 2001.
- [13] Srihari, R. K., Lamkhede, S. and Bhasin, A. "Unapparent Information Revelation: A Concept Chain Graph Approach," *Proc. ACM Conf. Information and Knowledge Management (CIKM '05)*, ACM Press, Bremen, Ger., Nov. 2005, pp. 329-330.