

# Unapparent Information Revelation

## A Concept Chain Graph Approach

Dr. R. Srihari, Sudarshan Lamkhede, Anmol Bhasin

### The Problem

#### Connecting the Dots...

"The importance of an all-source analysis cannot be overstated. Without it, it is not possible to connect the dots". No one component holds all the relevant information...

The 9/11 Commission Report

#### Unapparent Information Revelation

A document collection often reveals interesting information other than what is explicitly stated because these documents are generated by many authors working independently and at various times. Snippets of information from individual documents by themselves may seem inconsequential. However, when put together synergistically the synthesized document may reveal more information. We refer to this phenomenon as Unapparent Information Revelation (UIR). This paradigm mines documents to expose latent links that connect facts in an interesting way. It also aims at substantiating hypotheses with corpus support and leveraging links discovered.

#### UIR in perspective...

Query	IR	UIR
Keywords	Complex Query	Complex Query
ranked document list	ranked document list	ranked document sets
single document	single document	Single document cannot satisfy query
Documents ranked independently wrt query	Documents ranked independently wrt query	Connections between documents is important

#### Instances

- Text Mining for Intelligence Analysis
  - what are plausible connections between landing-gear and weather phenomena?
- Unintended information revelation
  - releasing information which is inconsistent with overall goals, objectives and operation of an organization
- Detect UIR scenarios
  - User-specified patterns
  - Quantify information represented by concept chains, document sets

### The Approach

#### Concept Chain Graphs (CCG)

##### CCG Architecture

###### Document Layer:

This layer contains documents and links (e.g. hyperlinks) between them.

###### Ontology Layer:

The Domain Ontology is roped in the architecture for Semantic Inference

###### Concepts and Associations Layer:

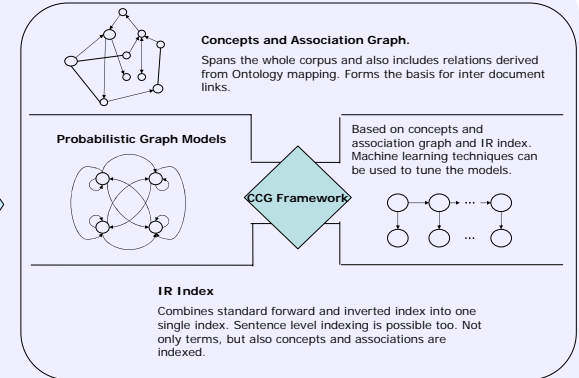
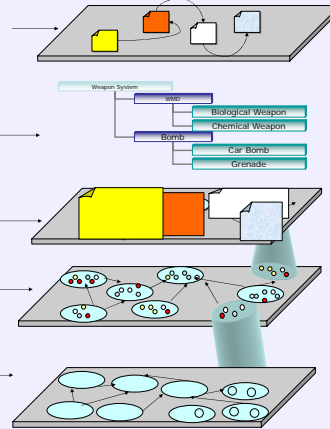
Consists of concepts and associations coming from corpus or ontology mapping.

###### Instances Layer:

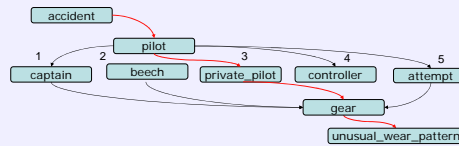
Tracks instances of concepts and associations detected in the corpus back to documents. Also maintains instance specific information last offsets and type.

###### Index Terms Layer:

Consists of index terms and hits.



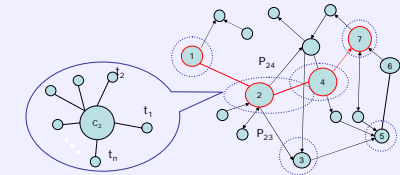
##### Evidence Trail



"A factor related to the accident was the pilot in commands..." [20001212X19333]. "The pilot was a private pilot with single and multi-engine land ratings..." [20001212X18707]. "Unable to fully extend one main landing gear, the private pilot retracted the other gear..." [20001212X20127]. "Other gears in the accessory case exhibited unusual wear patterns..." [20001204X00074]

##### Concept Chains

- Uses a weighted automaton defined on top of concept and associations graph.
- Represents a plausible connection between two given concepts.
- Transition probabilities based on context of the concepts.
- Context of a concept can be modeled in different ways depending on what user is looking for.

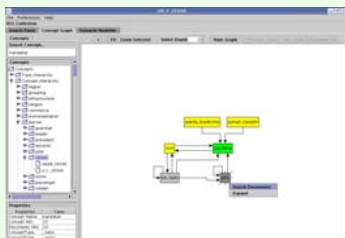


Context of C2

$P_{23}$ ,  $P_{24}$  are the probabilities of transition from C2 to C3 and C4 respectively.  $P_{23} < P_{24}$  as C2 and C4 have somewhat overlapping context.

### Applications

#### Visualization



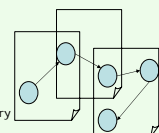
#### Document Set Retrieval

A minimal set of documents that completely "covers" a concept chain. (In general it is the minimal set of documents that answers all parts of a query.)

Presented to the user as the evidence of a retrieved concept chain.

Ranked according to the probabilities of generating the query chain.

Uses a Hidden Markov Model defined on CCG. The transition and emission probabilities can be learnt.

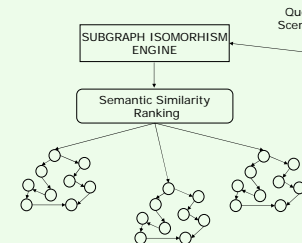


#### Cross Document Summarization



- Select the most "central" nodes from CCG based on degree centrality or information centrality.
- Select important sentences from the document set.
- An abstract can be created using the query concept chain as a theme for summary generation.

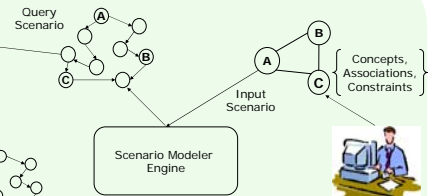
#### Graph Mining



Extracting Isomorphic scenarios to a query scenario is targeted at unearthing all instances of the pattern of interest based on the corpus.

The ranking is based on similarity of arcs of the graphs. Various similarity measures based on LSA and Dependency Tree Kernels can be applied.

#### Scenario Modeling



A scenario depicting the Input Scenario is constructed on the basis of the document corpus using Concept Chains.

The constraints used for the generation of scenarios can be extended to specifying concept and association types and chain selection criteria