# Example Application
# of
# PolyAnalyst
# with
# IATA STEADES Data

*Prepared by:*

**Dr. Sergei Ananyan**
**President**
**Megaputer Intelligence Inc.**
**120 W. 7th Street, Suite 310**
**Bloomington, Indiana 47404**
**Tel: +1 (812) 330-0110**
**Fax: +1 (812) 330-0150**
**E-mail:  s.ananyan@megaputer.com**

**&**

**Mike Goodfellow**
**Assistant Manager, STEADES**
**International Air Transport Association (IATA)**
**800 Place Victoria, P.O. Box 113**
**Montreal, Canada H4Z 1M1**
**Tel: +1 (514) 874-0202 x-3210**
**Fax: +1 (514) 874-2661**
**E-mail: goodfellowmj@iata.org**

*In Conjunction with:*

**GAIN Working Group B, Analytical Methods and Tools**

**September 2004**

# Preface

This example application has been prepared by Megaputer Intelligence Inc. and the International Air Transport Association (IATA) in conjunction with the Global Aviation Information Network (GAIN) Working Group B (Analytical Methods and Tools) (WGB) as one of a number of such examples of the use of analytical methods and tools described in the "*Guide to Methods & Tools for Airline Flight Safety Analysis*".  The intent of these example applications is to illustrate how various tools can be applied within an airline flight safety department, and provide additional information on the use and features of the tool and the value of such analysis.  GAIN WG B hopes that these example applications will help increase the awareness of available methods and tools and assist the airlines as they consider which tools to incorporate into their flight safety analysis activities.

Each example application of an analytical method or tool is posted on the GAIN website (*www.GAINweb.org*).  Readers are encouraged to check the website periodically for a current list of example applications, as further examples will be added as they become available.

# PolyAnalyst

# 1 Introduction

## 1.1 OVERVIEW OF THE TOOL FUNCTIONALITY

In-depth analysis of historical incident reports can help airlines, air traffic controllers, aircraft manufacturers and regulating organizations further improve aviation safety. A typical airline accumulates up to six thousand minor safety incident reports per year, containing a mixture of structured and textual data. Over 80% of useful information about a safety event is stored in the text narrative describing the event, which is inevitable in the description of any complex system or situation. Today the analysis of this data is carried out manually, but analysts fail to cope timely and efficiently with quickly growing volumes of data. Correspondingly, there is an industry-wide need for analytical tools capable to perform intelligent automated analysis of huge volumes of mixed textual and structured safety data.

*PolyAnalyst*™ is a tool that discovers important patterns and trends hidden in large volumes of data and helps better manage and share discovered knowledge among safety officers and predict outcomes of future situations. PolyAnalyst offers a comprehensive data analysis solution - from data importing, cleaning and manipulation, to visualization, modeling, scoring and reporting. *PolyAnalyst* can access data stored in any major commercial databases, spreadsheets, or flat files, as well as in popular document formats. It offers a broad selection of semantic text analysis, clustering, prediction, and classification algorithms, link analysis, transaction analysis, and powerful visualization capabilities.

Results obtained with PolyAnalyst provide insights into happenings in different aviation processes, helping safety officers to:

- Perform automated categorization of events based on their text descriptions
- Find similar events in historical data
- Carry out clustering of event reports and deliver interactive visual representations of results
- Reveal hidden problems, patterns and trends
- Generate strategic overview charts across different parameters
- Identify bottlenecks in processes and highlight quality / supplier related issues

## 1.2 INTRODUCTION TO THE EXAMPLE APPLICATION

This case illustrates how PolyAnalyst was applied for the analysis of safety databases containing both structured and narrative data and how it can assist in expediting the process of identifying hidden trouble spots and help improve aviation safety.

PolyAnalyst was applied to the analysis of safety event reports in IATA's Safety Trend Evaluation, Analysis and Data Exchange System (STEADES). The STEADES database consists of pooled data representing incident reports submitted to IATA by a large number of international airlines. There are about 50,000 reports filed each year and this amount is growing. Today STEADES contains about 300,000 safety reports. The goal of the project was to carry out an in-depth investigation of TCAS II (Traffic Alert Collision Avoidance System) related events. About 9,800 TCAS events were utilized in the performed analysis.

## 2   Input Data

The STEADES data is a combination of both structured data and natural language text: an anonymized summary of each event provided by the contributing airline representative, the event date and location, type of the aircraft involved, flight phase, and the risk level assigned by the airline to the event. In addition to these attributes, there exists a system of descriptors for categorizing each event in accordance with a preset taxonomy to facilitate future retrieval and analysis of records.

To ensure the most explicit interpretation of the results obtained from free text fields, user-defined dictionaries of domain-specific synonyms, stop-words and abbreviation expansions were imported in the system. Figure 1 shows the data as it appears after importing into the system.
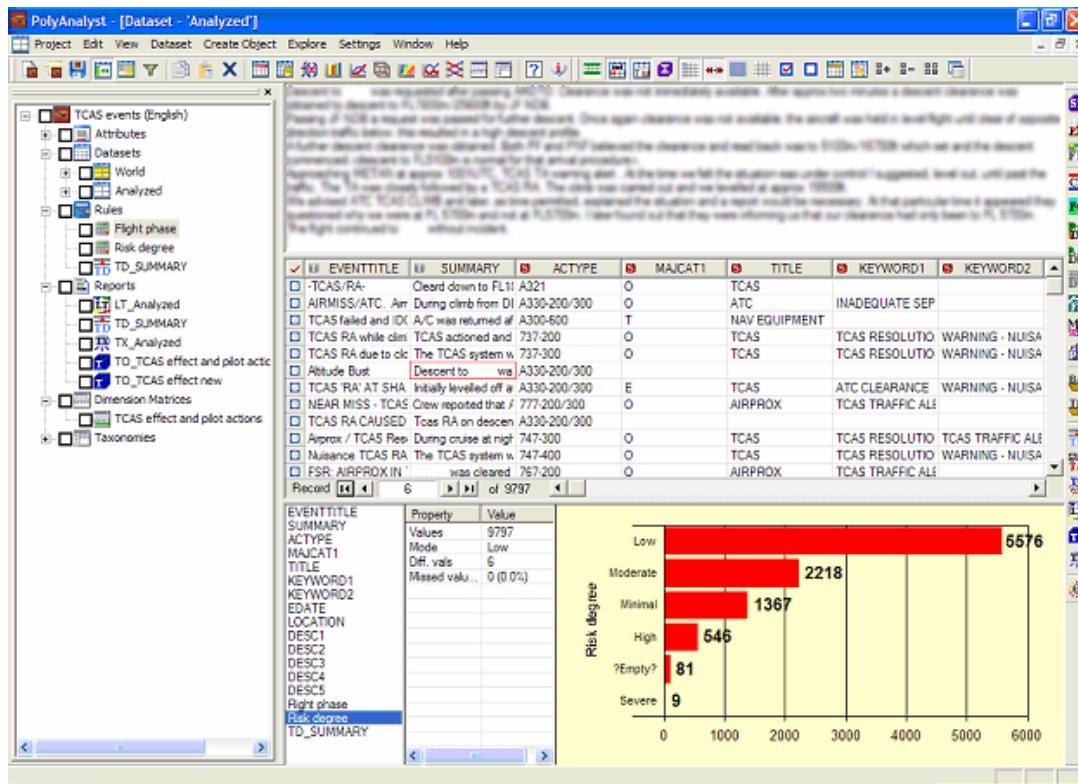


Figure 1 **STEADES data imported within PolyAnalyst. A histogram, in the bottom right pane, is an interactive representation of the distribution of the risk degree of events.**

No special formatting of original data is required for loading data in PolyAnalyst. The system offers a number of built-in data manipulation and cleansing tools, such as eliminating duplicate records, identifying almost identical fragments of text, and determining the language of narratives.

## 3   Analytical Process

PolyAnalyst provides a comprehensive set of tools for addressing major analytical tasks that safety officers are facing. It can be fine-tuned to a specific application domain.  A major portion of the user's involvement is in providing a direction for the analysis process and defining the main areas of interest.

The application of PolyAnalyst to the analysis of safety data at IATA provides the following main effects:

- Automating preliminary stages of data analysis, saving analysts' time for interpretation of the results (e.g. removing duplication and/or short summaries).
- Validating previously classified events. There may be the potential for shifting the responsibility of initial data categorization from airline representatives to computer algorithms monitored by IATA analysts in the future.
- Providing more elaborate and quick reporting capabilities.

The entire process of safety data analysis at IATA can be split in five major steps as depicted in Figure 2:

1) Data preprocessing
2) Report categorization
3) Problem areas ranking
4) Discovery of trends and patterns
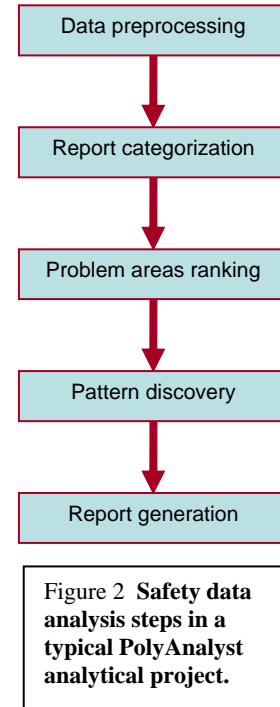5) Report generation

PolyAnalyst assists investigators with many steps of the data analysis process. The results are delivered in highly visual, interactive reports: easy to interpret, manage and share among the organization.



Figure 2 **Safety data analysis steps in a typical PolyAnalyst analytical project.**

# 4 Tool's Output

To perform automated categorization of incident reports, an analyst starts by defining patterns of terms for each taxonomy node determining which records are categorized to this node, iteratively perfects these node definitions, and then relies on the developed taxonomy to automatically categorize new safety reports. Once developed, a defined taxonomy should only require some further maintenance.

Patterns of terms can be defined by means of a simple but powerful Pattern Definition Language (PDL) developed by Megaputer Intelligence. This language allows the user to search for either exact terms, their morphological variations (through Stemming mechanism) or particular instances of terms (through built-in or user-added Semantic Dictionaries). PDL has a number of logical (AND, OR, XOR and NOT) and geometric (FOLLOW, SENTENCE, etc.) pattern definition operators helping the user to express any complex term pattern representing events to be monitored.

Using an existing STEADES-BASIS system of descriptors, and Megaputer analysts defined nodes in a two level Event Type-Event ID taxonomy with tentative patterns of terms based on the provided definitions of categories, as well as common sense. A fragment of the resulting taxonomy with defined nodes is shown in Figure 3.
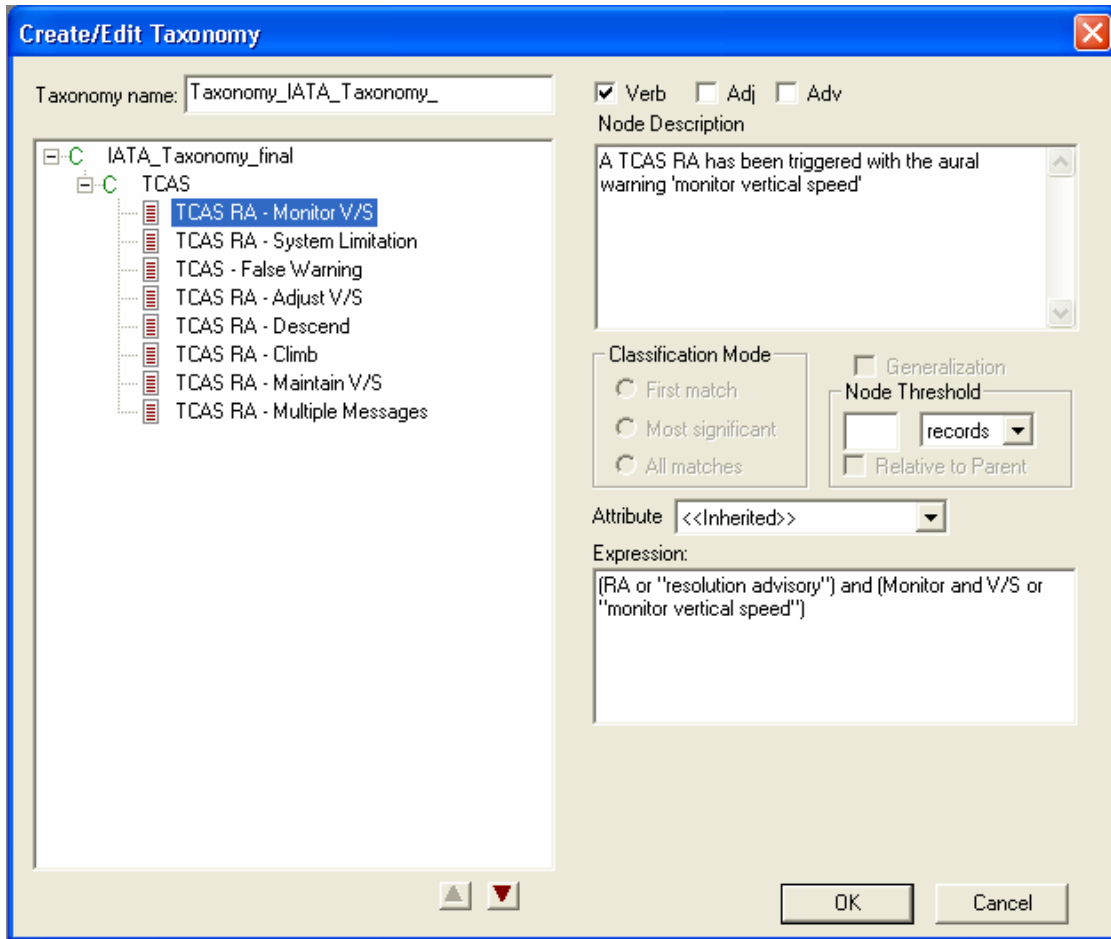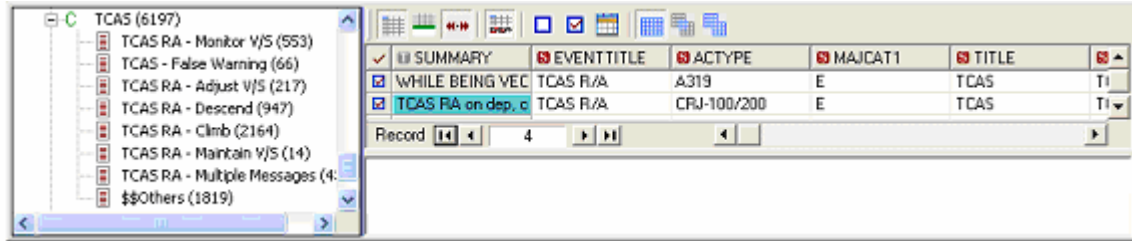
Figure 3 **PolyAnalyst Taxonomy Editor with imported IATA-British Airways taxonomy of event types and descriptors: each node is defined through the corresponding term pattern in PDL.**

For example, a pattern defining the "TCAS RA – Monitor V/S" requires records categorized to this node to contain a combination of one of the terms *RA* or *resolution advisory* and either the term pair *Monitor* and *V/S* or the term *monitor vertical speed*.

Applying the defined taxonomy to the IATA safety data resulted in matching the majority of records to one or more categories in the defined taxonomy. For example, the "TCAS RA - Monitor Vertical Speed" node caught 553 records, as can be seen in Figure 4. Terms matching the pattern are highlighted in different colors depending on the level of the taxonomy where they were encountered.

Figure 4  **The results of taxonomy-based categorization performed by PolyAnalyst: 553 records matched the definition of** *TCAS RA – Monitor V/S* **node.**

PolyAnalyst allows the user to easily browse through all relevant records with highlighted patterns of terms triggering categorization and to export the results in CSV or HTML files.

## 4.1  STRONGEST CORRELATIONS BETWEEN RISK DEGREES AND NARRATIVE TERMS

Link Chart performs calculating and visualizing pair-wise correlations between individual values of different attributes. Figure 5 reveals the strongest correlations between different risk degree values and terms extracted from event summaries through semantic text analysis. The heavier and more saturated color is the line on the graph, the stronger is the correlation between the corresponding objects on the graph.

For example, *High* risk degree is strongly correlated with the following terms (in decreasing strength order): *intruder*, *safety comment*, *same level*, *opposite traffic*, *safety* and *air traffic controller*. On the other hand, *Minimal* risk degree is strongly correlated with terms *flight*, *flight crew* and *notify*.
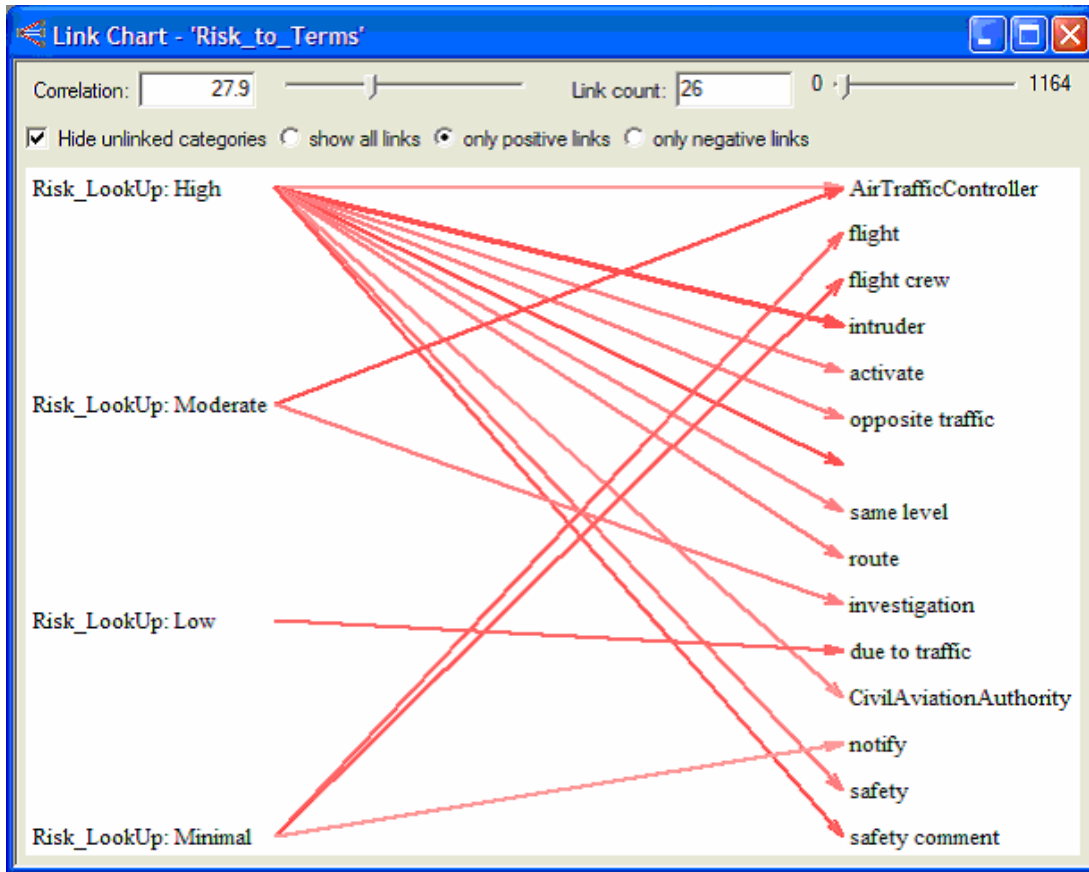
Figure 5 **Link Chart displays strong correlations between individual values of pairs of attributes.**

## 4.2 STABLE PATTERNS OF TERMS IN TEXT DESCRIPTIONS

Stable multi-dimensional patterns and correlations of terms can be discovered by the PolyAnalyst Link Terms engine. Patterns of terms occurring in a combination with a particular event can represent a characteristic signature of this event and can help trace nontrivial cause and consequence relations for this event.

The Link Terms diagram in Figure 6 displays several clusters of strongly correlated terms. Individual clusters are shown in different colors to facilitate simple visual identification. For example, the cluster shown in yellow indicates that the *minimum separation distance* to *intruder* is measured in *nautical miles*. The red cluster containing the *Ground Proximity Warning System* (GPWS) node suggests that, as expected, GPWS is strongly correlated with *radio altimeter*, *terrain*, and *warning triggered by traffic*, but this *warning* is often considered to be a *nuisance* when associated with a TCAS alert.
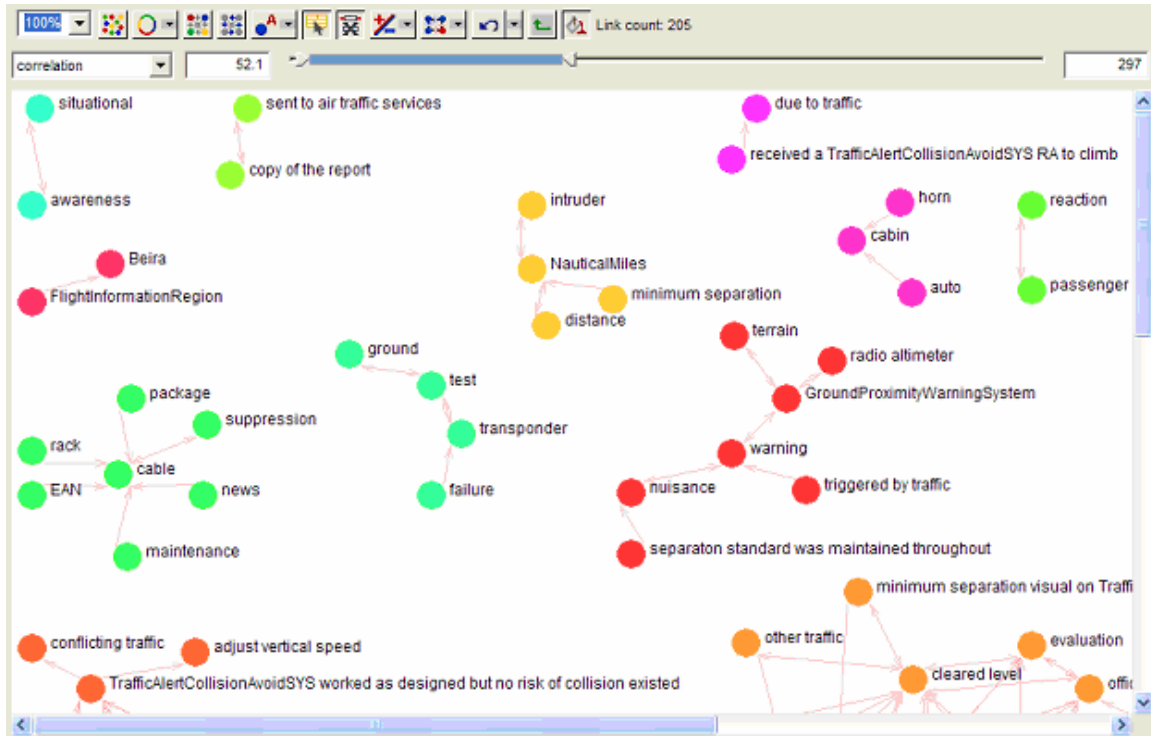
Figure 6 **Link Terms Diagram displays patterns of terms encountered in event text descriptions.**

## 4.3 CORRELATIONS BETWEEN STRUCTURED ATTRIBUTES AND NARRATIVES

The PolyAnalyst Link Analysis engine allows the user to visualize multidimensional correlations between values of structured attributes included in the analysis and strong patterns of terms in the corresponding event description summaries.

Figure 7 displays strong correlations between Risk Degree, Flight Phase and summary terms extracted from TCAS related events. One can observe that flight phases *Climb* and *Descent* have strong correlations with a phrase *monitor vertical speed*, while *Climb* is correlated in addition with *warning* and *rate of climb*. At the same time, *Moderate* risk degree is strongly correlated with terms *TA*, *heading*, *air traffic controller*, *opposite traffic*, *flight* and *action*.
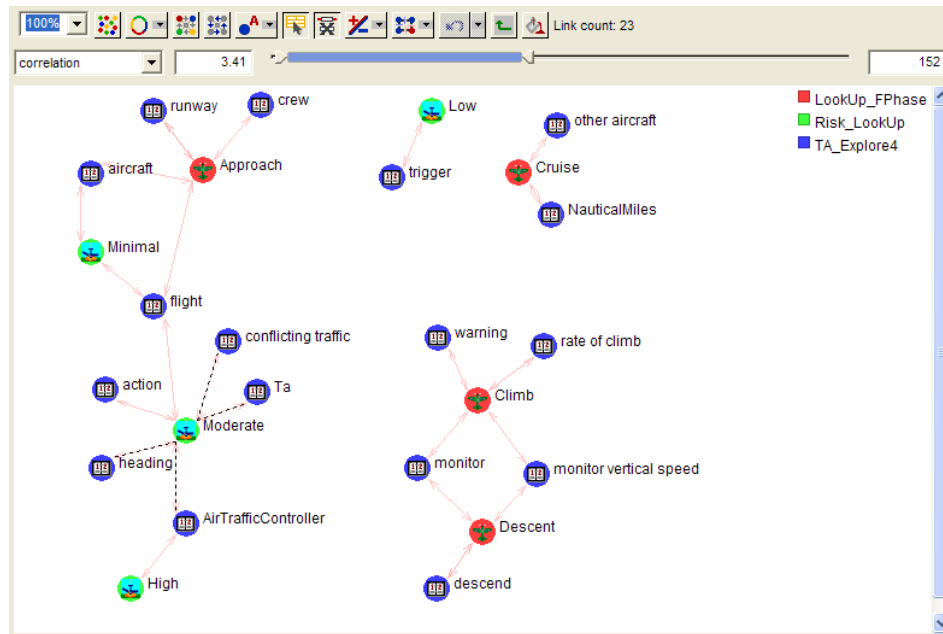
Figure 7 **Link Analysis Diagram helps visualize multi-dimensional correlations between values of structured attributes and patterns of terms extracted from text narratives.**

All PolyAnalyst Link Analysis engines support drill down capabilities. By clicking on a link of interest, the user displays all reports supporting the considered link with the corresponding terms highlighted in the narratives. The user can select how to perform the drill-down: selecting more than one link on a graph, the user can choose whether he wants to see an intersection or union of all records supporting the considered collection of links.

A user may select a number of links connecting Moderate risk degree value with strongly correlated narrative terms: *conflicting traffic*, *TA*, *heading* and *air traffic controller* and select the intersection drill-down mode. Figure 8 represents records that support all selected links at once.
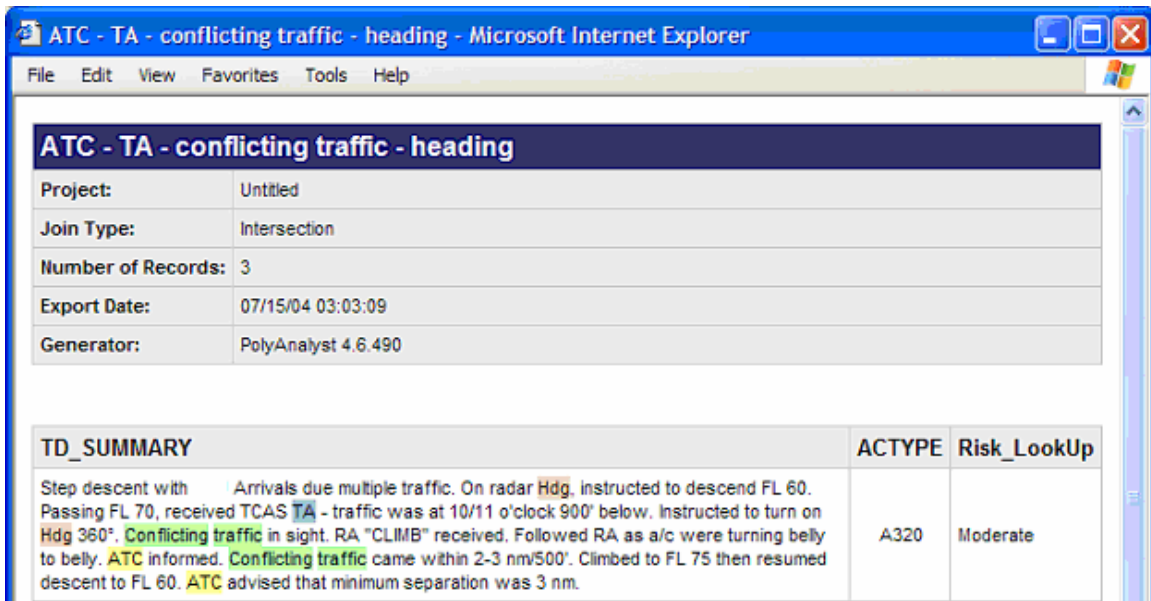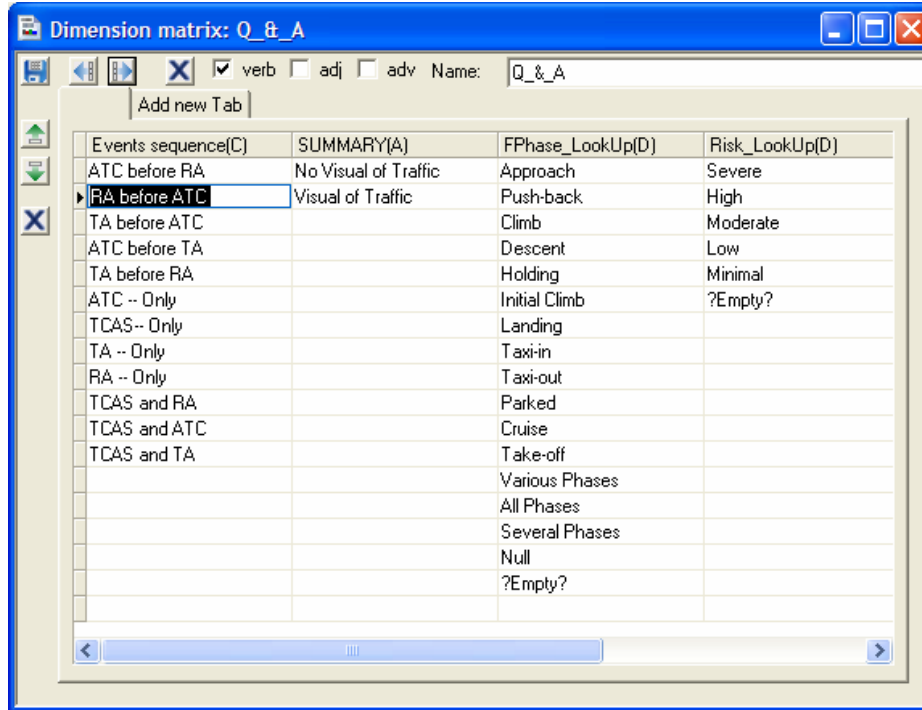
Figure 8  **Visual drill-down capability allows the user to select collections of original records supporting the selected link patterns.**

## 4.4  INTERACTIVE MULTI-DIMENSIONAL NARRATIVE INVESTIGATION – TEXT OLAP

Similarly to the analysis of structured data, defining dimensions of interest and developing a multi-dimensional cube holding data opens up new opportunities for the analysis of free form text. Manipulating the defined cube, the user can quickly slice and dice data across different dimensions, rotate data, and perform drill-down to original records with terms of interest highlighted.

To build an interactive Text OLAP report, the user defines a matrix of dimensions assigned to either textual or structured data fields. A dimension matrix utilized in this project is displayed in Figure 9.

Figure 9 **Dimension Matrix defines a set of dimensions utilized for interactive generation of analytical reports with PolyAnalyst Text OLAP.**

Values of individual cells in dimension matrices can be defined with the help of the same Pattern Definition Language (PDL) that was described when defining patterns for separate taxonomy nodes above. For example, the cell capturing RA events occurring prior to ATC can be defined through the following simple PDL expression: *follow*([RA], [ATC]).

Upon applying the developed dimension matrix to a collection of TCAS reports, one obtains a Text OLAP report, as depicted in Figure 10.

The first column creates a report on the sequence of events: for example, whether an ATC warning came prior to a Traffic Alert (TA) or Resolution Advisory (RA) warnings, or ATC was informed about the situation afterwards. Figure 12 illustrates that in the majority of cases (301 cases) the RA warning comes prior to an ATC communication, while in 264 cases ATC is mentioned prior to RA. Similarly, in 92 cases a TA precedes the RA warning. If the user decides to concentrate on these cases and clicks on the corresponding cell, she immediately sees the distribution of records stating that there was/wasn't a visual contact with the approaching traffic. It can be seen that in 45 of those cases where a TA preceded a RA or ATC alert, there was visual contact with the conflicting traffic, while in 20 cases there was no visual contact. Next, the user can click on the cell representing visual contact with traffic to see the distribution of flight phases where these TCAS events occurred. An immediate conclusion is that the majority of these events occurred during *Climb* (16 events) and *Descent* (11 events). Next, if a user selects the *Descent* and *Approach* cells to evaluate the final flight stages, the distribution of Risk Degree for the corresponding events is displayed. It can be seen that the majority of these events (9) had *Moderate* risk degree assigned to them.

Note that we followed only a single path for diving in the results – out of a large number of possible paths providing immediate answers to other important questions analysts might have.

An interactive navigation bar shows the analyst the current navigation path with the corresponding populations of intermediate nodes. The text viewer at the bottom of the window allows the user to browse through individual text records with the corresponding patterns of terms highlighted.
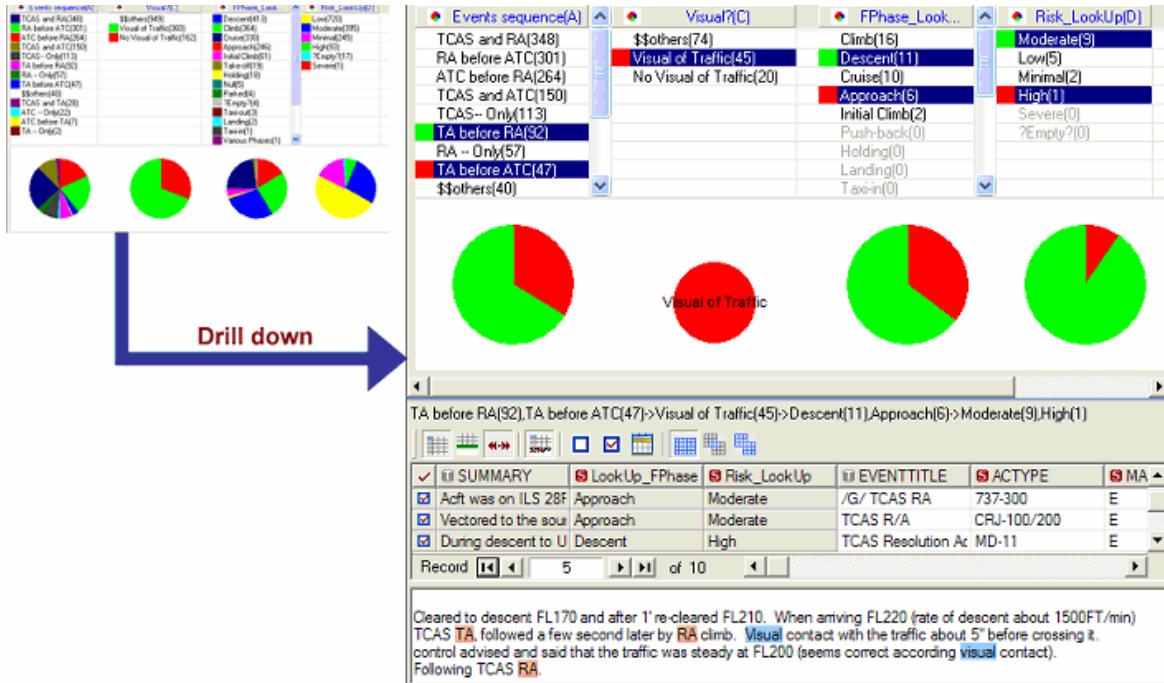


Figure 10 **PolyAnalyst Text OLAP report view displaying events where TA was received prior to RA, conflicting traffic was seen visually, the problem occurred during the descent and was assigned moderate risk level.**

An interactive grid in the middle of the window allows the user to export a CSV or HTML report capturing interesting findings. The user can export an HTML report listing all records supporting TCAS TA warnings followed by RAs, which occurred during the *Descent* and were assigned *Moderate* risk.

If one needs to see a distribution of all visual contact versus no visual contact events first rather than splitting them by different sequences of ATC and TCAS TA and RA warnings, one simply needs to right click on the second column in the dimension matrix and move it to the left.

## 4.5  FIND SIMILAR

In order to better classify a new safety event, Flight Safety Officers sometimes need to locate similar events in past reports. Due to large volumes of historical safety data, analysts often cannot necessarily find similar reports in a timely manner. PolyAnalyst offers an efficient mechanism for carrying out this task automatically: the Find Similar engine allows the user to retrieve safety reports similar to the report being considered. Patterns of terms making the retrieved reports similar to the new report are highlighted and the measure of relevance for these reports is calculated. Similarity of reports is calculated with the help of a Case Based Reasoning algorithm.

Figure 12 presents the summary of a historical safety event that the system identified as most similar to an arbitrary selected report depicted in Figure 11. One can observe that these summaries are considered to be similar because they contain the same terms *wake turbulence*, *experienced* and *glidescope*.
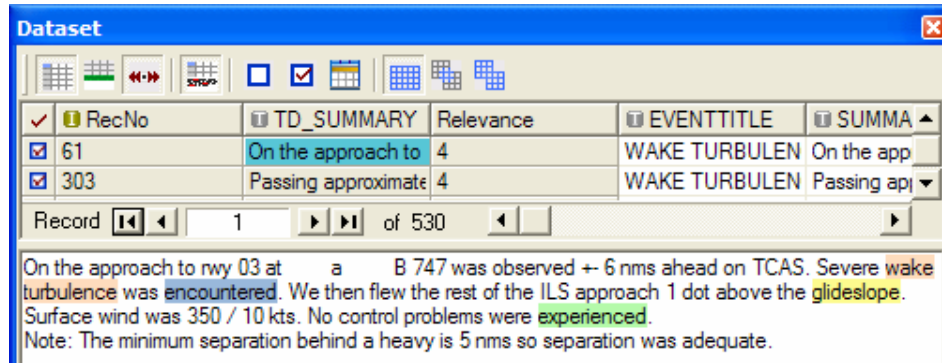
Figure 11 **Arbitrary safety report selected by the user. PolyAnalyst Find Similar engine identified a historical record most similar to the selected record, as illustrated below.**
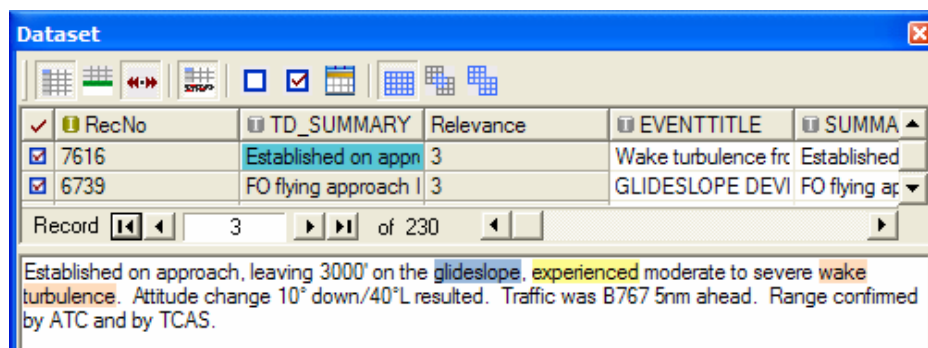


Figure 12 **Report identified to be most similar to the selected report. Highlighted are collections of terms that make the two reports similar.**

Upon discovering patterns and trends of interest in the data, analysts may need to share their findings with other colleagues. In the case of IATA STEADES analysis, the results of safety data analysis are delivered to safety officers in all member organizations.

Recognizing the need for results sharing, PolyAnalyst offers a number of simple means for generating business reports capturing the important findings throughout the project. Currently PolyAnalyst provides the following main mechanisms for adding results to business reports:

1) All graphical objects can be right clicked and saved in one of several graphical formats (BMP, JPG, PNG or WMF).
2) All objects representing results of the analysis support the drill-down capability for arriving at a collection of records supporting the obtained results.
3) The results of drill-down can be saved in HTML format.
4) Is this supposed to be another mechanism or should it be disregarded?

A combination of these mechanisms allows the user to conveniently build business reports with graphical elements, as illustrated in Figure 13.
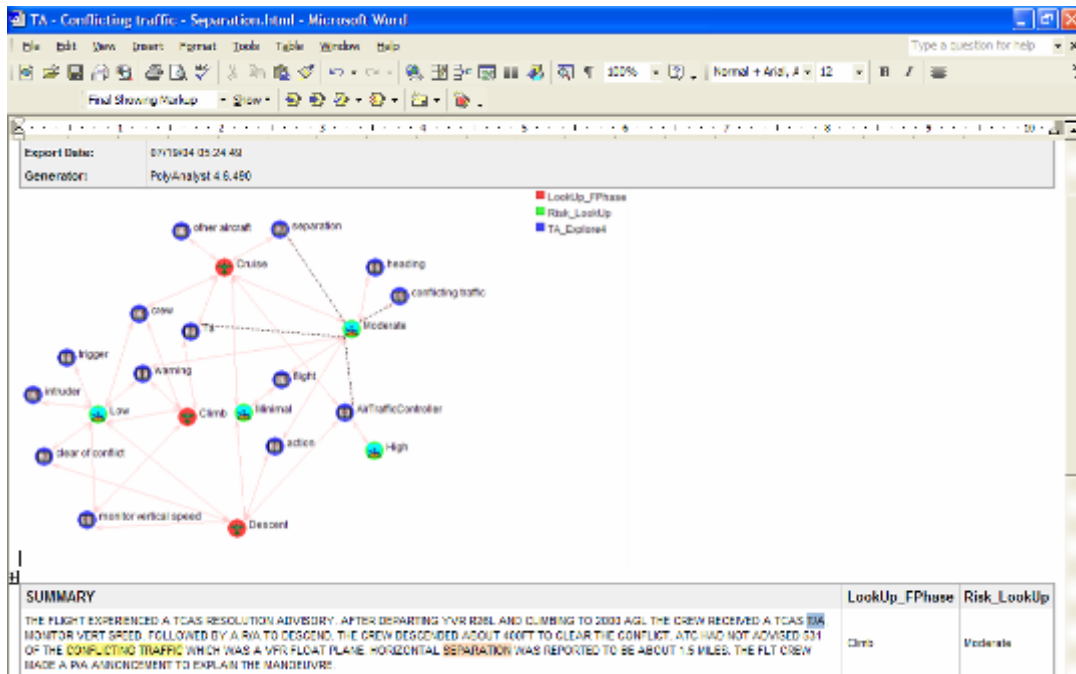
Figure 13 **Sample PolyAnalyst business report in Microsoft Word format.**

# 5 Application of the Results of the Analysis

This case illustrated how PolyAnalyst can be utilized for:

1) Monitoring data for known issues and assisting in categorizing safety events to predefined classes
2) Discovering and visualizing patterns from the analysis of raw data
3) Quickly finding similar reports in historical data
4) Summarizing the obtained results in easy to interpret reports.

Overall, the results obtained by PolyAnalyst can be further investigated and manipulated within the system and exported to reports, while the predictive models can be scheduled for online execution or applied to data in the original database to predict outcomes of future situations.

The case demonstrated that value could be generated through:

- Using PolyAnalyst to extend the analytical capabilities beyond an existing classification system
- Efficient use of analysts' time for many tasks
- Potential of automating repetitive processes and reducing manual processing
- Quick, intelligent analysis of textual data in certain cases
- Consistent and comprehensive use of both structured and unstructured data

In summary, PolyAnalyst could be helpful in aviation safety by the early and accurate detection of problem areas; patterns and trends based on the analysis of incident reports data. An evaluation copy of PolyAnalyst can be downloaded from www.megaputer.com.