# Flight Safety Data Analysis with PolyAnalyst™

Case study

**Prepared by:**
Megaputer Intelligence
120 W. 7th Street, Suite 310
Bloomington, IN 47404

By
Vijay Kollepara
Sergei Ananyan

**Megaputer**
*Your Knowledge Partner*™

May 2003

# Table of contents

## Overview of Application

Aviation safety experts surmise that accidents are usually a culmination of a series of unsafe events that had gone unnoticed. For every accident and major event report that is thoroughly investigated, there can be 300 incident reports (Heinrich's triangle) that could have contained information about the impending accident. These reports can be in the form of pilot reports, maintenance reports or incident reports at different parts of Flight Operations.  Incident reports represent a combination of structured data and free form text narratives stored in a database.   The analysis of these reports is carried out is predominantly manually so far.

*PolyAnalyst™* is a comprehensive and user-friendly data and text mining system.  It offers a complete end-to-end data analysis solution - from data importing, cleaning and manipulation, to visualization, modeling, scoring and reporting.  *PolyAnalyst* can access data stored in any major commercial database and some proprietary data formats (Excel, SAS), as well as 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.  PolyAnalyst can be a convenient solution for easy, thorough and accurate analysis of aviation safety data.

Results obtained with PolyAnalyst provide key insights into happenings in different aviation processes, helping safety officers to
   a) Reveal hidden problem issues (irrespective of data type – structured or unstructured)
   b) Generate strategic overview charts for the management across different parameters
   c) Identify bottlenecks in processes and highlight quality / supplier related issues

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

## Input Data

The Aviation Safety Reporting System (ASRS) receives, processes, and analyzes reports of unsafe occurrences and hazardous situations that are voluntarily submitted by pilots, air traffic controllers, and others.  PolyAnalyst has been applied to analyzing ASRS data from NASDAC.  The chosen period covered October and November of 2001 and included 7,500 records and 61 attributes (including free form text in narrative fields).

The data was imported in PolyAnalyst using a built-in Data Import Wizard.   To ensure the most explicit interpretation of the results obtained from free text fields, user-made dictionaries of domain-specific synonyms, stop-words and abbreviation expansions were also imported in the system.
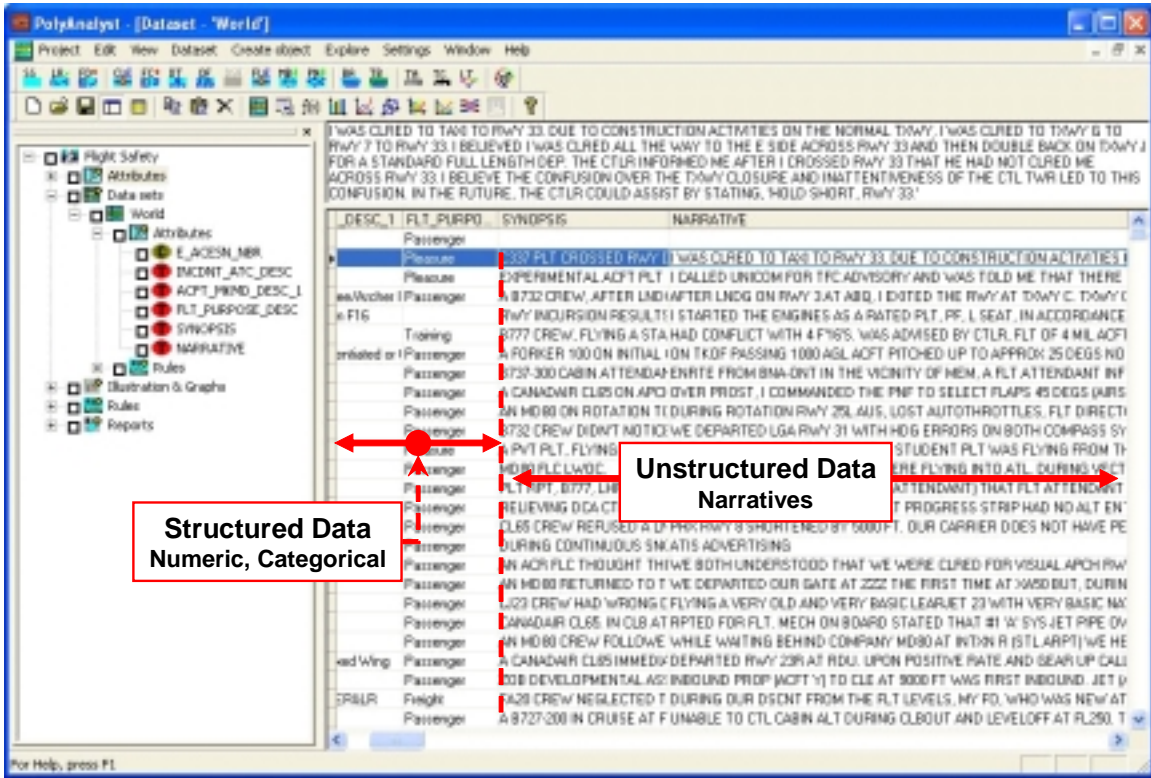
**Fig. 0. Snapshot of the investigated data.**

PolyAnalyst can directly access data from any major commercial database through standard OLE DB and ODBC protocols. Fig. 1 shows the data as it appears after importing into the system. The data contains both structured and narrative fields.

## Analytical Process

PolyAnalyst provides a comprehensive set of tools that can address many analytical tasks that safety officers are facing and can be easily fine-tuned to a specific application domain. A major portion of the user's involvement is in providing direction to the analysis process and defining their areas of interest. User-defined parameters for running analysis engines are entered in the corresponding dialog boxes.

Broadly, the process of gaining knowledge from narratives would involve two main steps: extraction and interpretation of knowledge. Further discussion will primarily concentrate on the analysis of unstructured portion of the data, as it often contains over 80% of useful information, while analysts lack efficient tools for the analysis of textual notes.

a) *Identify and extract all terms of interest occurring in narratives*

Fig. 2 illustrates simple steps performed by the user to run the PolyAnalyst Text Analysis (TA) engine to identify important concepts being discussed in the narratives. Text analysis can be carried out in two modes:

- Unsupervised TA Mode: In this mode, the TA engine extracts all important concepts occurring in the text.
- Supervised TA Mode: The user can guide the TA engine to only search and extract concepts of interest to them. For example, by defining the broad concept '*equipment*' the system can return a list of all equipment or device results like '*radar*', '*horn*', '*pump*' and '*rudder*' that have been discussed in the narrative.
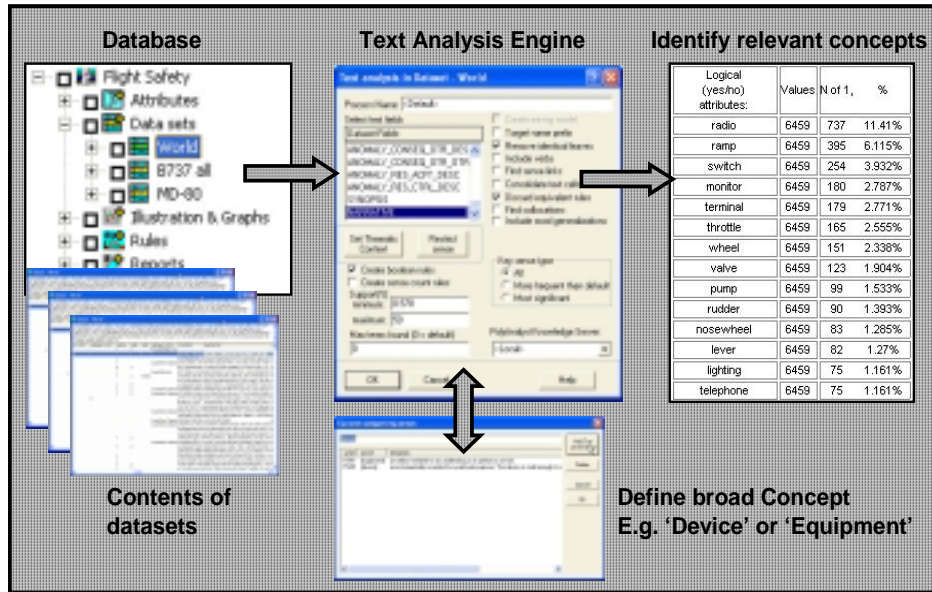
4

**Fig. 2. Automated text analysis exploration.**

Once the main terms are extracted, the user becomes able to either simply export the concepts to a Microsoft Excel sheet or conduct other advanced analysis and visualization within PolyAnalyst.

*b) Generate actionable reports for Management*

The system incorporates different visualization techniques to enable the user generate explicit and actionable results. Fig. 3 illustrates two visualizations graphs the user can employ to better understand patterns of terms and relations between them that had been identified in the previous step.

The Snake Chart is utilized for providing a comparative overview of concepts across different business entities. The Link Terms engine conducts 'n-dimensional' correlation analysis and visual layout of the results to help reveal close associations and patterns of terms in the investigated data.
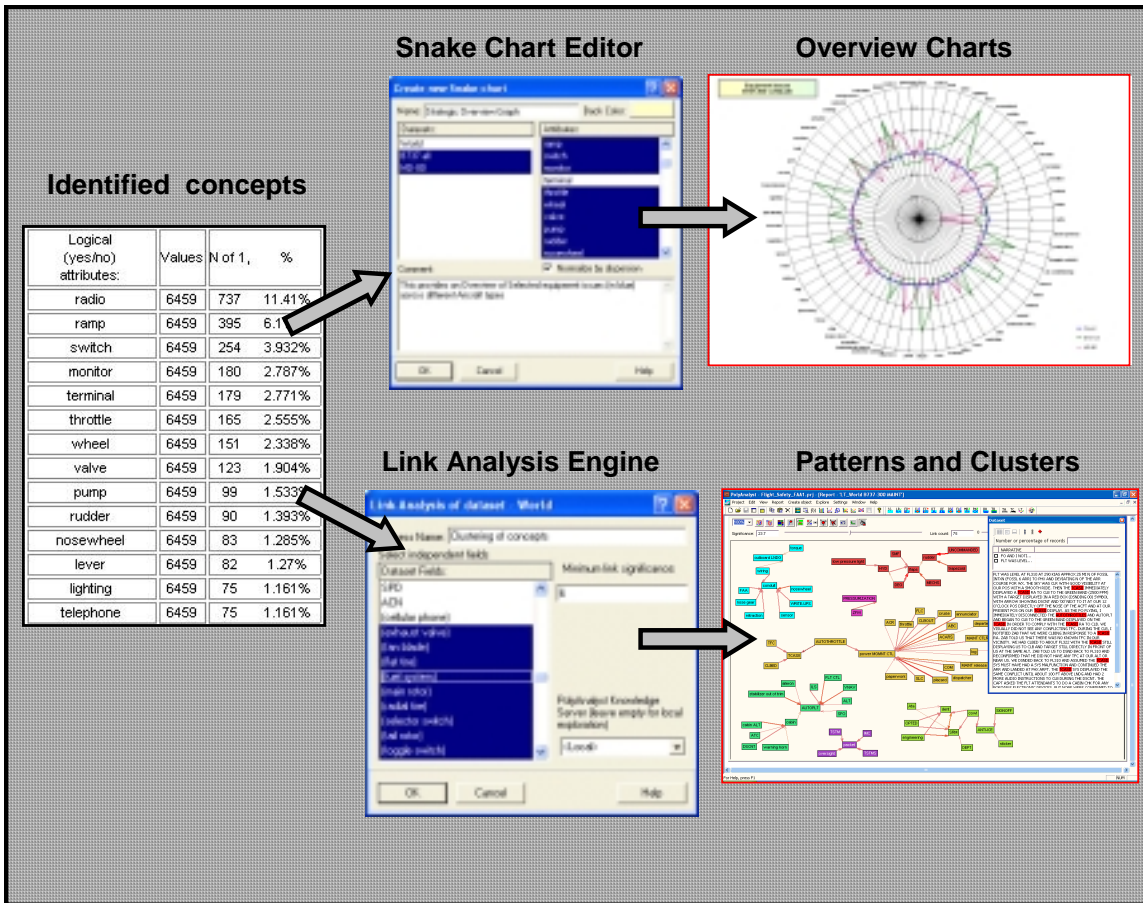
**Fig. 3. Possible steps of typical data analysis scenarios.**

In addition to the above two processes, PolyAnalyst provides numerous other analysis and visualization capabilities and scenarios based on the needs and desires of the analyst. Overall, it offers sixteen different analytical engines and a few dozen visualization techniques that can be used either independently or sequentially to derive new knowledge from data. This broad range of analytical engines also allows the user to conduct the analysis irrespective of the type of data (Numeric, Boolean, Categorical or Textual).

## *Tool's Output*

### 1) Intelligent Text Analysis

Flight Safety Officers (FSO) would like to quickly learn major concerns as perceived by the pilots. Of course, one can try to extract all text concepts occurring in the narratives and then manually browse through them. However, a far more powerful mechanism would be to search for concepts within a specific domain, and let the text analysis engine do the hard work of identifying and reporting back only relevant concepts.

An example of this type of analysis (also called supervised analysis) would be to identify only equipment related issues mentioned by pilots in their narratives. Being instructed to focus on specific concepts ('*equipment*' and '*device*' in this project), the Text Analysis (TA) engine sifts through the entire Narrative portion of the database and automatically returns back concepts like '*radio*', '*switch*', '*brakes*', and '*nosewheel*'. Fig. 4 pictorially represents this process of intelligently extracting chosen categories of concepts (in this case, *Equipment*).
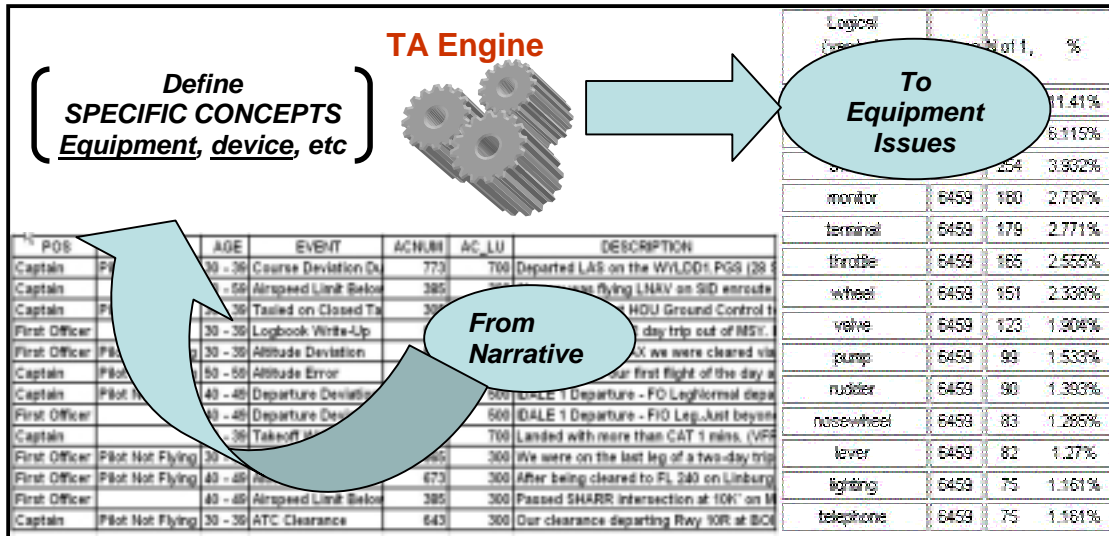


**Figure 4: Process of searching for specific concepts.**

Note that the system is smart enough to understand the '*equipment*' query and then identify all related words and phrases. The system enables specific user-desired charts and visualizations as shown in Fig. 5.
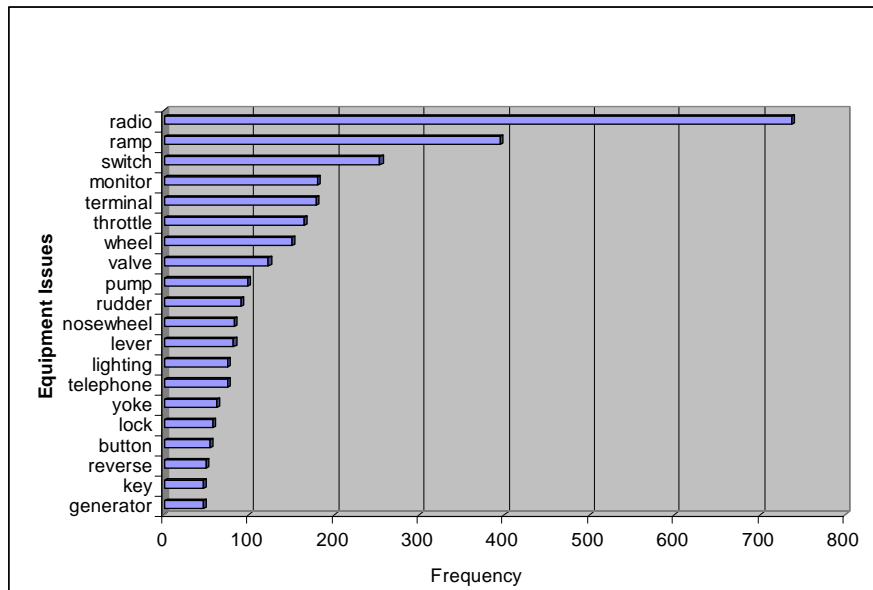


**Figure 5: Equipment related issues mentioned by Pilots.**

## 2) Gain Strategic Insight

Airline management seeks better understanding of how the pilot concerns are varying across the company business assets, divisions and other categories. PolyAnalyst Snake Charts, utilized in combination with supervised text analysis, can deliver such understanding by producing insightful views of the investigated issues.

For example, Fig. 6 below shows how individual '*Equipment*' concerns (identified by the TA engine, as illustrated in Fig. 4 above) can be compared across different aircraft types, in this case - B737 and MD-80.

The blue line on the graph represent an overall average frequency occurrence across all aircrafts (World data set contains all investigated incident reports). Green and pink lines represent relative frequencies of 'Equipment' concerns for B737 and MD-80 compared to their overall background. The spikes indicate that



**Fig. 6. Compare relative importance of equipment-related issues across different aircraft types.**

the occurrence of these terms is relatively more frequent than average. For example, B737 (green line) has issues related to '*screw', 'controller', 'lever', 'toggle switch', 'wiper'* and '*spool*' while MD-80 (pink line) has issues related to '*flat tire', 'blade', 'receiver', 'faucet'*, etc.

Now the management can see problematic issues and put in place preventive mechanisms to avert future problems. The Snake Chart enables a quick comparison of different tracked parameters across any dimensions of interest.

The user can also click on a chosen concept, say '*lever*' issue associated with Boeing 737 and drill down to view the associated records (see Fig. 7).

**Fig. 7.  Drill-down feature allows viewing records related to the investigated issue.**

### 3) Identify high correlation entities

Calculating and visualizing mutual correlations of attribute values, one gains knowledge of stable patterns of co-occurrences of different values of individual attributes.  Fig. 8 suggests a quick way to view the most important correlations between items of interest and learn if any of the pilot concerns have a high correlation with specific aircraft types.



**Fig. 8.  Correlations between Aircraft type and pilot concerns extracted from free text Narratives.**

The intensity of the line is a measure of the strength of the corresponding correlation.  The user can infer that *B737-300, B727-200 and Regional Jet CL65* aircraft types (on the left side) have high correlations with the term *'MAINT'* mentioned in the narrative (right side of the figure).  Another inference from the above chart could be the high correlation of the concepts *'pax', 'gate'* and *'attendant'* to *MD-80 Super 80* aircraft

type. Note that the user can easily visualize correlations between important items from both structured and unstructured parts of the database.

## 4) Identify Patterns of concerns:

An ability to capture stable patterns of terms, not known in advance but describing main issues of concern without having to read through all the records, can provide valuable insights for quick comprehension of past experience and save time of an FSO for more advanced analysis. PolyAnalyst Link Terms engine can be used to reveal clusters of terms from the narrative portion of pilot reports. The Link Terms diagram (Fig. 9) displays the discovered patterns of terms and relations between them.



**Fig. 9. View clusters of terms that represent stable patterns discovered in textual Pilot narratives.**

Link Terms produced ten clusters, each denoted by a different color. These clusters now instigate the user to find out why

- *'Rudder'* is highly correlated with *'UNCOMMANDED', 'trim', 'trapezoid', 'anomaly'* and *'logbook'* (cluster shown in light green)?
- *'TCASII'* is mentioned together with *'AUTOTHROTTLES'* (cluster shown in light blue)?

Additional insight into individual clusters is provided by drilling down and viewing the corresponding narratives with terms of interest highlighted. Fig. 10 presents the results of drilling down on the *'UNCOMMANDED' <--> 'rudder'* link from the "rudder" (light green) cluster of the above Link Terms diagram, thus giving an analyst the ability to quickly verify significance of patterns of interest.

THE PROB AROSE WHEN I CALLED OUR MAINT COORDINATOR IN ATLANTA TO RPT INFO ON AN ACFT THAT HAD BEEN GATHERED ON 2 LEGS. I USED A POOR CHOICE OF WORDS TO DESCRIBE WHAT WE HAD EXPERIENCED. THE MAINT COORDINATOR MISUNDERSTOOD WHAT WAS SAID AND DECLARED WE HAD HAD AN UNCOMMANDED RUDDER EVENT. THE MISUNDERSTANDING WAS FURTHER EXACERBATED BY MY SELECTION OF WORDS IN THE LOGBOOK WRITE-UP. THIS IS UNFORTUNATE AS MY INTENTION WAS SIMPLY TO HAVE MAINT CHK THE ACFT AFTER EXPERIENCING AN UNDETERMINED ANOMALY DURING THE LAST TKOF. WHEN I ARRIVED AT THE ACFT IN SLC, I NOTICED THE RUDDER TRIM WAS SET APPROX 2 UNITS OF L TURN. I ZEROED OUT THE TRIM PER PROC AND INFORMED THE CAPT. THE ACFT WAS HVY AND FLAPS WERE AT 1 DEG FOR THE TKOF. ROTATION WAS AT APPROX 150 KTS. I HAD ANTICIPATED THAT THE ACFT MIGHT ROLL RIGHT AFTER LIFTOFF, AND IT DID. I CORRECTED, BUT IT TOOK MORE AILERON THAN I FIRST EXPECTED TO HOLD HDG. I TRIMMED THE RUDDER TO APPROX 2 UNITS L RUDDER AND THE ACFT FLEW FINE. WHEN WE LEFT IAH, THE CAPT DECIDED TO LEAVE THE RUDDER SET WITH APPROX 2 UNITS OF L TRIM FOR HIS TKOF. THE ACFT WAS AGAIN HVY AND ROTATION WAS AT APPROX THE SAME SPD AS IN SLC. FLAPS WERE SET AT 1 DEG. THE RWY WAS VERY ROUGH AND THE CAPT LIFTED THE NOSE SLIGHTLY TO KEEP THE GEAR FROM BOUNCING DOWN THE RWY. I FELT A SLIGHT BUMP IN MY SEAT AT ROTATION. FROM MY PERSPECTIVE, I DIDN'T PERCEIVE IT AS A MOVEMENT OF OR R. IN RETROSPECT, IT COULD HAVE BEEN WHEEL CASTER, TURB, JETBLAST OR THE ROUGH RWY. ASKED THE CAPT ABOUT THE RUDDER PEDALS AND HE SAID THEY DID NOT MOVE. I PERCEIVED NO ROLL, YAW, OR 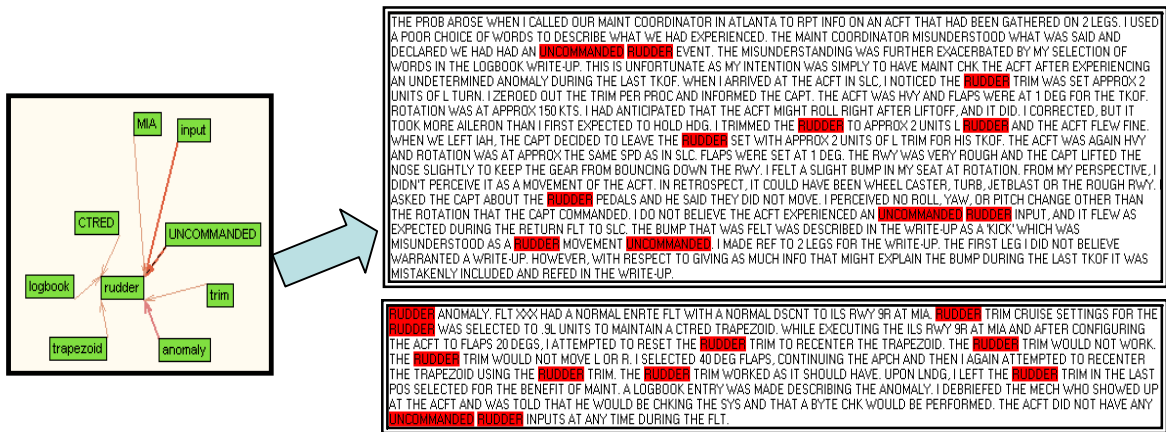PITCH CHANGE OTHER THAN THE ROTATION THAT THE CAPT COMMANDED. I DO NOT BELIEVE THE ACFT EXPERIENCED AN UNCOMMANDED RUDDER INPUT, AND IT FLEW AS EXPECTED DURING THE RETURN FLT TO SLC. THE BUMP THAT WAS FELT WAS DESCRIBED IN THE WRITE-UP AS A 'KICK' WHICH WAS MISUNDERSTOOD AS A RUDDER MOVEMENT UNCOMMANDED. I MADE REF TO 2 LEGS FOR THE WRITE-UP. THE FIRST LEG I DID NOT BELIEVE WARRANTED A WRITE-UP. HOWEVER, WITH RESPECT TO GIVING AS MUCH INFO THAT MIGHT EXPLAIN THE BUMP DURING THE LAST TKOF IT WAS MISTAKENLY INCLUDED AND REFED IN THE WRITE-UP.

RUDDER ANOMALY. FLT XXX HAD A NORMAL ENRTE FLT WITH A NORMAL DSCNT TO ILS RWY 9R AT MIA. RUDDER TRIM CRUISE SETTINGS FOR THE RUDDER WAS SELECTED TO .9L UNITS TO MAINTAIN A CTRED TRAPEZOID. WHILE EXECUTING THE ILS RWY 9R AT MIA AND AFTER CONFIGURING THE ACFT TO FLAPS 20 DEGS, I ATTEMPTED TO RESET THE RUDDER TRIM TO RECENTER THE TRAPEZOID. THE RUDDER TRIM WOULD NOT WORK. THE RUDDER TRIM WOULD NOT MOVE L OR R. I SELECTED 40 DEG FLAPS, CONTINUING THE APCH AND THEN I AGAIN ATTEMPTED TO RECENTER THE TRAPEZOID USING THE RUDDER TRIM. THE RUDDER TRIM WORKED AS IT SHOULD HAVE. UPON LNDG, I LEFT THE RUDDER TRIM IN THE LAST POS SELECTED FOR THE BENEFIT OF MAINT. A LOGBOOK ENTRY WAS MADE DESCRIBING THE ANOMALY. I DEBRIEFED THE MECH WHO SHOWED UP AT THE ACFT AND WAS TOLD THAT HE WOULD BE CHKING THE SYS AND THAT A BYTE CHK WOULD BE PERFORMED. THE ACFT DID NOT HAVE ANY UNCOMMANDED RUDDER INPUTS AT ANY TIME DURING THE FLT.

**Fig. 10.  Some pilot narratives that support the Link Terms cluster centered on '*rudder*'.**

The number of reports submitted by pilots grows over time causing the relevance of concepts to change too. Thus the above links between terms may assume more or less important significance as time progresses and can serve as a valuable tool for knowing whether there are changing patterns.

## *Application of Results of Analysis*

The outputs of the link analysis and snake charts deliver explicit and actionable results that can be used by the business manager to rectify the observed anomalies.  For example the safety manager can use this type of information to assess Pilot concerns across age groups.  Another example would be for flight operations and maintenance analysts to identify issues related to aircraft makes, types, and specific aircraft.  The scope of the analysis can also be extended to identifying problem parts and supplier quality issues.

Overall, the results obtained by PolyAnalyst can be further investigated and manipulated within the system and exported in a report, while the discovered predictive models can be scheduled for online execution or applied to data in the original database to store the predicted outcome of future situations.  The Megaputer integration platform lets power users of the system to record reusable analytical scripts implementing typical data exploration scenarios, and business users – access over the Internet the resulting periodically updates reports in the format of preset templates.

### Conclusion

This paper outlined just a few standard scenarios for safety data analysis that can be performed with the help of PolyAnalyst.  The described project demonstrated that a synergetic combination of automated text analysis and visual presentation of discovered clusters and correlations can significantly reduce the latency and bias of the analysis, automate the most time-intensive operations and increase the thoroughness and quality of the obtained results.

PolyAnalyst generates significant new business value through:
- Capturing previously unanticipated knowledge from raw data
- Efficient use of analyst's time
- Automation of repetitive processes, which removes latency of manual processing
- Reduced cost and increased accuracy of incident data analysis
- Quick intelligent analysis of textual data
- Consistent and comprehensive utilization of *all* available data (structured and unstructured)

Summarizing, PolyAnalyst helps improving aviation safety and preventing potential accidents through early and accurate detection of problem areas from the analysis of incident reports data.  An evaluation copy of PolyAnalyst can be downloaded from www.megaputer.com.