# RSTTool 2.4 – A Markup Tool for Rhetorical Structure Theory

Michael O'Donnell (micko@dai.ed.ac.uk)
Division of Informatics, University of Edinburgh.

## Abstract

RSTTool is a graphical tool for annotating a text in terms of its rhetorical structure. The demonstration will show the various interfaces of the tool, focusing on its ease of use.

## 1 Introduction

This paper describes the RSTTool, a graphical interface for marking up the structure of text. While primarily intended to be used for marking up Rhetorical Structure (cf. Rhetorical Structure Theory (RST): Mann and Thompson (1988)), the tool also allows the mark-up of constituency-style analysis, as in Hasan's Generic Structure Potential (GSP - cf. Hasan (1996)).

The tool is written in the platform-independent scripting language, Tcl/Tk, and thus works under Windows, Macintosh, UNIX and LINUX operating systems.

RSTTool is easy to use, one creates an RST diagram from a text by dragging from segment to segment, indicating rhetorical dependency. There is a separate interface for text segmentation. The tool can automatically segment at sentence boundaries (with reasonable accuracy), and the user clicks on the text to add boundaries missed by the automatic segmenter (or click on superfluous boundaries to remove them).

The tool was previously described in O'Donnell (1997). However, since then the tool has been substantially revised and extended,(the current version being 2.4). This version is also far more robust due to extensive debugging by one of RST's inventor's, Bill Mann. Particular improvements in the tool include:

1. *GUI for defining relation:* ability to add, rename, delete, etc. the relations used using a graphical user interface.

2. *Statistical Analysis:* a new interface was added, which allows users to be presented with statistics regarding the proportional use of relations in a text.

3. *Output Options:* the new tool allows saving of RST analyses in postscript (for inclusion in Latex documents), or sending diagrams directly to the printer. Files are now saved in an XML format, to facilitate importation in other systems.

4. *Improved Structuring:* the possibilities for structuring have improved, allowing the insertion of spans, multinuclear elements and schemas *within* existing structure.

The Tool consists of four interfaces, which will be described in following sections:

1. *Text Segmentation:* for marking the boundaries between text segments;

2. *Text Structuring:* for marking the structural relations between these segments;

3. *Relation Editor:* for maintaining the set of discourse relations, and schemas;

4. *Statistics:* for deriving simple descriptive statistics based on the analysis.

## 2 What is RSTTool For?

The RSTTool is an analysis tool, but most users of the tool are researchers in the text generation field. For this reason, we present the tool at this conference.

Several reasons for using the tool are:

- *Corpus Studies:* before one can generate text, one must understand the rhetorical patterns of language. By performing analyses of texts similar to which one wishes to generate, one can identify the recurrent structures in the text-type and work towards understanding their context of use.

- Results Verification: often, a particular study may be challenged by other researchers. If the study was performed using RSTTool, the corpus supporting the study can be released for analysis by others. Previously, most RST analysis was done by hand, making distribution of
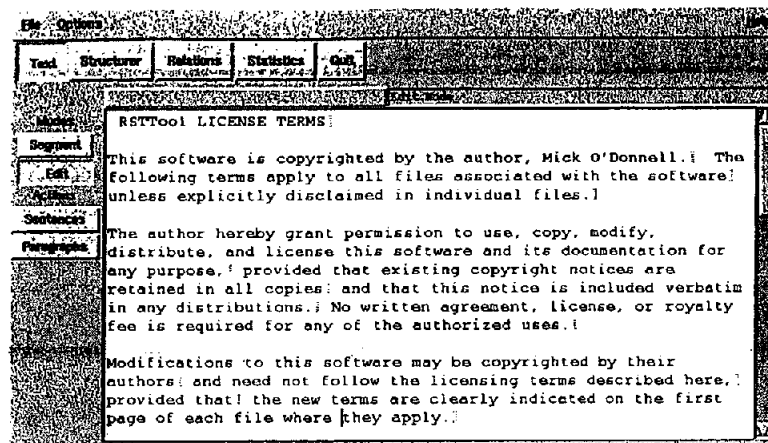
Figure 1: The Segmentation Interface

corpora difficult. RSTTool thus not only simplifies the production of the corpus, but also allows ease of distribution and verification.

● *Diagram Preparation*: the RSTTool can also be used for diagram preparation, for inclusion in papers. The tool allows diagrams to be exported as EPS files, ready for inclusion in LaTeX documents (as demonstrated in this paper). For PCs and Mac, screen-dumps of diagrams are possible (Tcl/Tk does not yet fully support export of GIF or JPG formats, and conversion from EPS to other formats is primitive). Some versions of MS Word allow the inclusion of EPS diagrams.

● *A Teaching Tool*: by getting students to analyse texts with the RSTTool, teachers of discourse theory can increase the student's understanding of the theory.

To allow RSTTool analyses to be more generally usable, the tool now saves its analyses in an XML format, making loading into other systems for processing much simpler.

## 3 Segmentation Interface

The first step in RST analysis is to determine segment boundaries. RSTTool provides an interface to facilitate this task. The user starts by "importing" a plain text file. The user can then automatically segment at sentence boundaries by pressing the "Sentences" button. This segmentation is not 100% reliable, but is reasonably intelligent. The user can then correct any mistakes made by the automatic segmentation, and also add in segment boundaries within sentences.

To add a segment boundary, the user simply clicks at the point of the text where the boundary is desired. A boundary marker is inserted. To remove a boundary, the user simply clicks on the boundary marker. Figure 1 shows the Segmentation interface after clausal segmentation.

The user can also edit the text, correcting mistakes, etc., by switching to Edit mode.

The user then moves to the Structuring interface by clicking on the "Structurer" button at the top of the window. Note that the user can return at any point to the Segmentation interface, to change segment boundaries, or edit text. These changes are automatically accounted for in the structuring component.

## 4 Structuring Interface

The next step involves structuring the text. The second interface of the RSTTool allows the user to connect the segments into a rhetorical structure tree, as shown in figure 2. We have followed the graphical style presented in Mann and Thompson (1988).

The tool supports not only RST structuring, but also constituency structuring. I believe that texts cannot always be analysed totally in terms of rhetorical relations, and that some level of schematic analysis complements the rhetorical analysis. For instance, a typical conference paper (such as this one) can be assigned a top level schematic structure of

```
Title ^ Author ^ Institution ^ Abstract
    ^ Section* ^ Bibliography
```

The RSTTool allows intermixing of such schema with RST analysis.

Initially, all segments are unconnected, ordered at the top of the window. The user can then drag the mouse from one segment (the satelite) to another (the nucleus) to link them.

The system allows both plain RST relations and also multi-nuclear relations (e.g., Joint, Sequence,
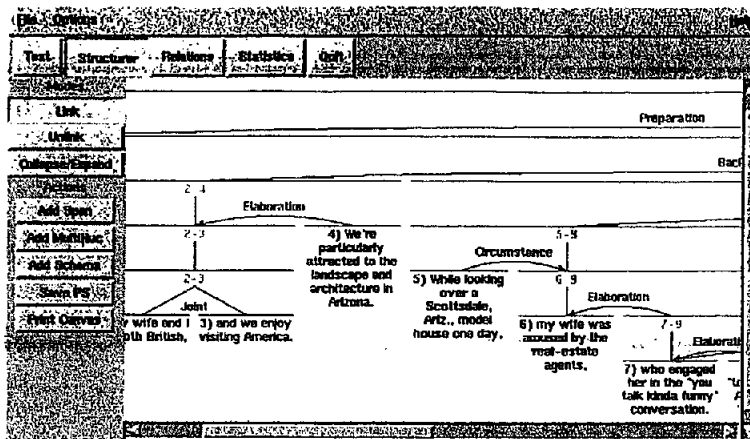
254

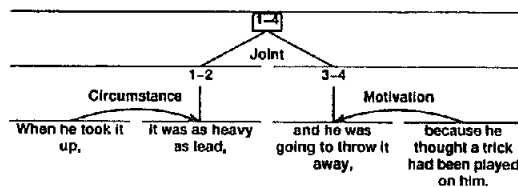Figure 2: The Structuring Interface
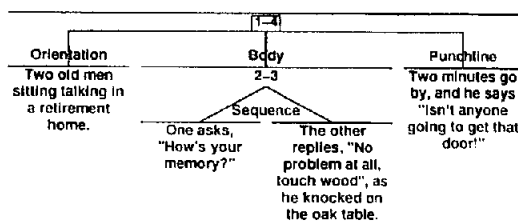


Figure 3: RST Structuring



Figure 4: Schema-based Structuring

etc.). Scoping is also possible, whereby the user indicates that the nucleus of a relation is not a segment itself, but rather a segment and all of its satellites. See figure 3 for an example combining normal RST relations (Circumstance, Motivation); multinuclear structure (Conjunction), and scoping (the nodes marked 1-2 and 3-4). In addition, schemas can be used to represent constituency-type structures. See figure 4.

Because RST-structures can become very elaborate, the RSTTool allows the user to collapse subtrees – hiding the substructure under a node. This makes it easier, for instance, to connect two nodes

which normally would not appear on the same page of the editor.

## 5   Editing Relations

The tool provides an interface for editing relation sets. The user can add, delete or rename relations. If the relation is in use in the current analysis, the changes are propagated throughout the analysis.

## 6   Statistical Analysis

Discussions on the RST mail list have demonstrated that there is a community concern with frequency of different relations in specific text-types. The RSTTool, by providing counts of relations within a text, supports this research goal. See figure 5.

The interface shows not only the frequency of relations, but also the ratio of Nuc Sat orderings to Sat Nuc orderings for the relation (valuable data for both generation and automatic discourse structure recognition).

## 7   Summary

RSTTool is a robust tool which facilitates manual analysis of a text's rhetorical structure. These analyses can be used for a number of purposes, including i) to improve understanding of discourse structure, to aid in either text generation or analysis; ii) diagram preparation, and iii) as a teaching tool.

The main improvement in the latest version of the tool is the statistical analysis interface. Later versions of the tool will extend on this aspect, increasing the range of analyses which can be performed on each text, or collection of texts.

Future versions will also add code for automatic structure recognition, using such work as Marcu's RST recognition tool (Marcu, 1997). While the author believes that automatic recognition is not yet reliable, integrating such a tool into an RST Markup
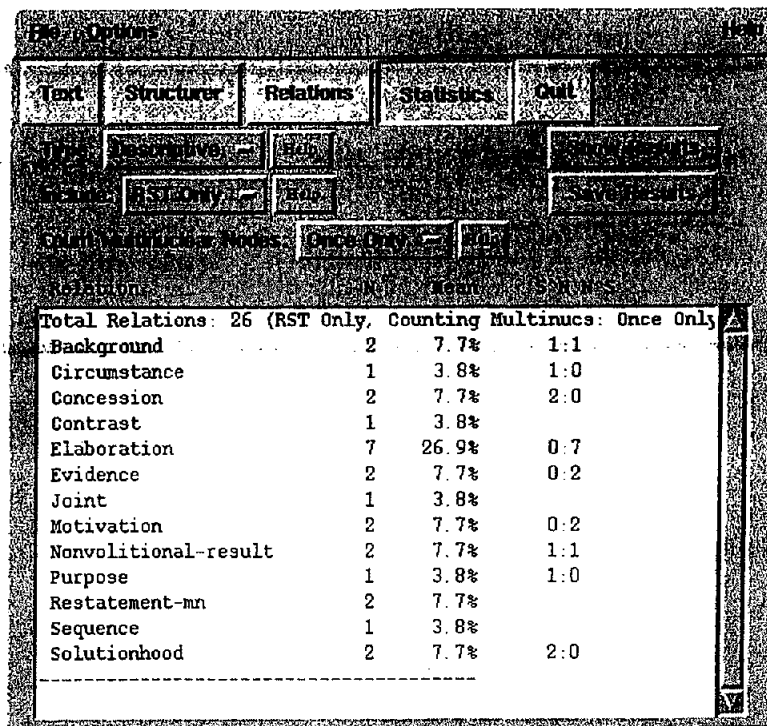
Figure 5: The Statistics Interface

tool allows the recognition software to provide a *first draft*, which the human editor can correct to their liking. At present, such a mixture of automatic and human-directed mark-up is the best way of achieving accurate mark-up of text structure.

# References

Ruqaiya Hasan. 1996. The nursery tale as a genre. In Carmel Cloran, David Butt, and Geoff Williams, editors, *Ways of Saying: Ways of Meaning.* Cassell, London. Previously published in Nottingham Linguistics Circular 13, 1984.

W.C. Mann and S. Thompson. 1988. Rhetorical structure theory: Toward a functional theory of text organization. *Text*, 8(3):243–281.

Daniel Marcu. 1997. The rhetorical parsing of natural language texts. In *The Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics, (ACL'97/EACL'97)*, pages 96–103, Madrid, Spain, July 7-10.

Michael O'Donnell. 1997. Rst-tool: An rst analysis tool. In *Proceedings of the 6th European Workshop on Natural Language Generation*, pages 92 – 96, Gerhard-Mercator University, Duisburg, Germany. March 24 - 26.

256