# Telecommunications use in Education to provide interactive visual feedback on Automatically Graded Essays

Robert Williams
School of Information Systems
Curtin University of Technology
Perth, Australia
bob.williams@cbs.curtin.edu.au

Heinz Dreher
School of Information Systems
Curtin University of Technology
Perth, Australia
h.dreher@curtin.edu.au

## Abstract

This paper discusses the idea of automatic computer marking of assignments and linking learners and teachers to virtually immediately available interactive visualizations of the assessment results. MarkIT is an Automatic Essay Grading System being developed by the authors, and for which they have a provisional patent application filed. It makes use of proprietary algorithms using Natural Language Processing techniques. At its heart is an English language thesaurus used under license from Macquarie Dictionary Pty Ltd. MarkIT provides, in addition to a numerical essay score, comprehensive visual and textual feedback on essay content to enable the teacher to discuss the strengths and weaknesses of the essay, and areas for improvement with the student within a very short time of the essay being submitted for assessment. An example of MarkIT system functionality and performance is presented. Further system development possibilities are presented to permit upload, grading and interactive feedback provision via the Web or wireless networks.

**Key Words:** assignment assessment, Automatic Essay Grading, Natural Language Processing, NLP, telecommunications, feedback interaction, feedback visualization.

## Introduction

Automated Essay Grading (AEG) systems are attracting increased attention from users and developers alike (Williams, 2001; Shermis and Burstein, 2003; Valenti, Neri and Cucchiarelli, 2003; Williams and Dreher, 2004). Articles about new systems are frequently being published on the World Wide Web (WWW), and the systems are now in use by a number of educational organisations (Phelan, 2003).

To achieve a score for an essay, AEG systems employ a variety of scoring algorithms. As an example, the algorithms for the five AEG systems featured below use different conceptual models of essay scoring.

MarkIT (Williams and Dreher, 2004), the system being developed by the authors, is based on content and document linguistic features, derived by implementing modified Natural Language Processing (NLP) techniques.

Project Essay Grade (PEG) (Page, 1994), one of the earliest and longest-lived implementations of automated essay grading, uses various linguistic features of the essay documents to determine a score or grade.

E_rater (Attali and Burstein, 2004) has been developed at the Educational Testing Service in the USA. This system uses a hybrid approach combining linguistic features derived from use of NLP techniques with other document structure and statistical features.

The Intelligent Essay Assessor (Landauer et al., 1998) makes use of Latent Semantic Analysis (LSA) and the "bag of words" approach, and has been developed and evaluated at the University of Colorado at Boulder. It ignores document linguistic and structure features.

The final system, which uses text categorisation techniques (TCT) (Larkey, 1998), has been developed at the University of Massachusetts. It uses a combination of modified key words and linguistic features.

The performance of these systems approaches that of human graders – typically correlation coefficients for human-human and computer-human grades are similar (Shermis and Burstein, 2003). Table 1 shows some of the reported performances of four of the AEG systems discussed above.

**Table 1: Comparative performance of models**

| System | Measure | Values | Reported in |
|--------|---------|--------|-------------|
| PEG | r | 0.389-0.743 | Page, 1994 |
| E_RATER | % | 87-94 | Burstein, et al, 1998 |
| LSA | % | 85-91 | Landauer, 1999 |
| TCT (soc) | r | 0.69-0.78 | Larkey, 1998 |
| TCT (G1) | r | 0.69-0.88 | Larkey, 1998 |

(Source: Williams, 2001, p181)

**MarkIT**

MarkIT is the name given to the AEG system being developed by researchers at Curtin University, and we are able to present an insight into the technology used by MarkIT in some detail. The fundamental question that had to be answered by the developers of MarkIT was: "How can we represent essay semantic content in a robust computationally effective manner, in such a way that unseen text can be processed within seconds?" The following sections explain how we implemented the functionality implied by this question.

**Contemporary Semantic Representations**

Poesio (2000) discusses some current techniques for representing the meaning of sentences, including First-Order Logic and Semantic Networks. First-Order Logic uses mathematical expressions representing set membership relations for objects in the sentence belonging to sets to represent the meaning of a sentence. Computationally, this is very difficult to use for unlimited unseen text, as generally domain specific information needs to be hand coded. Semantic Networks classify objects into a network of nodes and relationships. Arc traversal of the nodes can be used to imply relationships amongst the nodes. Substantial domain specific knowledge needs to be hand coded prior to their use. We decided that these approaches to semantic representation would not be practical for AEG algorithms.

**"Chunking" as an Alternative to Full Parsing**

Many NLP systems use some kind of a parser to initially extract the syntax of sentences in a document as an initial step prior to further processing. Semantic analysis then follows. The use of Context Free Phrase Structure Grammar (CFPSG) parsers is commonly suggested in the literature. However CFPSG parsing cannot be used in all but simple toy domains. For

large volumes of unseen text, a large number of grammar rules are needed for typical parsers, and this generally means that processing time becomes unrealistic for practical applications. So while CFPSG parsing has been tried with the prototype system described in this article, it has been abandoned in favour of using "Chunking" to determine the phrases and clauses used for further processing.

Useful preliminary linguistic computation can be done with less structured parsing. Phrase chunking can be an effective alternative to full parsing at the initial processing stage. Chunking has the advantage that it does not require an extensive set of grammar rules – a few simple rules suffice. Specifically chunking breaks a sentence into syntactically structured components representing noun clauses or phrases and verb clauses or phrases. Often, these structures are sufficient as a preliminary to further processing.

## Slot Structures contain Semantic Details of Chunks of Text

The technique we have selected makes use of chunking to get the structure of sentences in terms of subject and predicate, as represented by Noun Phrases (NP) and Verb Phrases (VP). Generally the NP nominates the subject of discussion, and the VP the actions being performed on or by the subject. However VPs are notoriously complex to deal with in comparison to NPs, because they typically can have many clusters of a Verb Clause (VC) and a NP together. It is far easier to identify VCs instead of the complex VPs. The basis of the technique used is to represent the meaning of the words making up the NPs and VCs in a sequence of structured slots containing a numerical value representing the thesaurus index number for the root meaning of the word in the slot. A numerical summary of the meaning of the sentences in the document being considered is thus built up.

## NP Structure

Kolln (1994, p433) states a rule for defining an NP under transformational grammar as follows:
 (1) NP = (DET) + (ADJ) + N + (PREP PHR) + (S)

and a PREP PHR (Kolln 1994, p429) as follows:
PREP PHR = PREP + NP

When considering the slots to be provided for a  NP,
(1) above can now be rewritten as
 (2) NP =  DET ADJ N PREP NP S

The basic component of an NP appears to be
 (3) NP = DET ADJ N
and some appended structures. It has been found in practice that
 (4) NP = DET ADJ ADJ ADJ N
is a better structure. If we take this as a basic core structure in a NP, the complete NP structure can be built in terms of this core structure by linking multiple occurrences of this core structure by PREPs. It has been found in practice that we should also allow linking by CONJs. So finally we conclude that the basic component should be

(5) NP = CONJ PREP : DET ADJ ADJ ADJ N

where the two slots before the colon are the linking slots, and those following are the content slots. Figure 1 illustrates this structure. Practice indicates that we should allow about 40

| |
|---|
| ADJ = adjective<br>ADV = adverb<br>AUX = auxiliary<br>COMP = complement<br>CONJ = conjunction<br>DET = determinant<br>N = noun<br>NP = noun phrase<br>PHR = phrase<br>PREP = preposition<br>S = sentence<br>V = verb<br>VC = verb clause<br>VP = verb phrase |

occurrences of this basic component as the NP slot template should handle many practical NPs encountered in general English text. In the current version of the system, we allow for an unlimited number.

| CONJ | PREP | DET | ADJ | ADJ | ADJ | N |
|------|------|-----|-----|-----|-----|---|
|      |      |     |     |     |     |   |

**Figure 1: Noun Phrase Semantic Structure**

The first core component in a sentence generally will have the CONJ and PREP slots set to a zero. Any empty slots will likewise be set to zero.

## VC Structure

Kolln (1994, p428) states a rule for defining a VP under transformational grammar as follows:
 (6) VP = AUX + V + (COMP) + (ADV)

COMP is explained as an NP or ADJ, so by removing this from the VP we end up with a VC as follows:
(7) VC = AUX + V + ADV

It has been found in practice that if we modify this VC definition by the addition of extra AUXs and ADVs we obtain a more useful structure as:
 (8) VC = AUX AUX ADV ADV V AUX AUX ADV ADV

VCs can often be introduced with CONJs, and it has been found in practice that we should also allow PREPs in a VC, so a complete VC definition would be:
 (9) VC = CONJ PREP AUX AUX ADV ADV V AUX AUX ADV ADV

Practice indicates that we should allow about 40 occurrences of this basic component as the NP slot template should handle many practical VCs encountered in general English text. Figure 2 illustrates this structure. In the current version of MarkIT we allow for an unlimited number of occurrences.

| CONJ | PREP | AUX | AUX | ADV | ADV | V | AUX | AUX | ADV | ADV |
|------|------|-----|-----|-----|-----|---|-----|-----|-----|-----|
|      |      |     |     |     |     |   |     |     |     |     |

**Figure 2: Verb Clause Semantic Structure**

If a sentence happens to start with a VC, then the CONJ slot will be set to zero as are any empty slots.

## Sample Processing

Consider the following sentence:
"For example if people working on a group project did work their own way and on their own schedule it would be extremely difficult to coordinate their work and assure the quality and timeliness of the end product."   (Source: Alter, 1992, p168)

The above sentence, taken from an introductory Information Systems textbook is processed as follows (slots containing blanks and zeroes have been eliminated because of space limitations; an 8888 indicates a missing thesaurus entry for the corresponding word).

NOUN PHRASE
For example                    if people
0 3027 0 0 0 0 22 0 4012 0      0 0 0 0 997 0

VERB PHRASE
working
0 0 0 0 0 0 677 0 0 0 0 0

NOUN PHRASE
on a group                    project
0 3034 2000 0 0 0 712 0      0 0 0 0 0 0 8888 0

VERB PHRASE
did work
0 0 5008 0 0 0 677 0 0 0 0 0

NOUN PHRASE
their                  own way                  and on their
0 0 0 0 0 0 2008 0      0 0 8888 0 0 0 26 0      4003 3034 0 0 0 0 2008 0
own schedule              it
0 0 8888 0 0 0 86 0      0 0 0 0 0 0 8888 0

VERB PHRASE
would be extremely
0 5032 5002 8888 0 0 0 0 0 0 0

NOUN PHRASE
to difficult
0 3046 0 868 0 0 0 0

VERB PHRASE
coordinate
0 0 0 0 0 0 60 0 0 0 0 0

NOUN PHRASE
their
0 0 0 0 0 0 2008 0

VERB PHRASE
and      work                    assure
4003 0 0 0 0 0 677 0 0 0 0 0      0 0 0 0 0 0 858 0 0 0 0 0

NOUN PHRASE
  the quality          and timeliness          of the end                  product
0 0 2007 0 0 0 5 0      4003 0 0 0 0 0 8888 0      0 3032 2007 0 0 0 620 0      0 0 0 8888 0 0 0 0

VERB PHRASE
0 0 0 0 0 0 0 0 0 0 0 0

## Scoring Student Essays by Matching a Model Answer against Student Answers
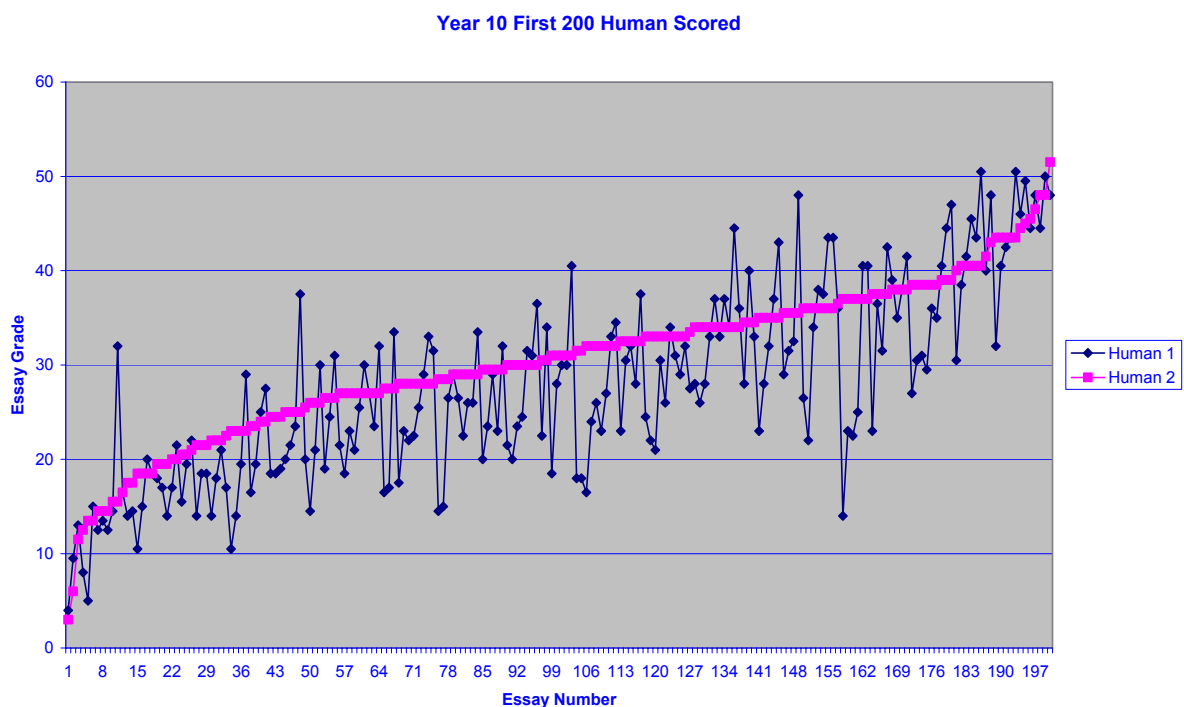
We now give an overview of the scoring algorithm. A model answer is prepared by the instructor that contains the core knowledge required to achieve a 100% score. The system may eventually be able to score a student essay against a number of model answers, in case the instructor wishes to use numerous content models. The instructor can also provide about 100 - 200 human graded essays (ideally each graded by three humans) and their scores, for training purposes. These model and training answers are processed as described above. The

system then performs a content matching task in which the model answer content summary is compared against each of the training essay content summaries. Many aspects of the relationships between the model and training essays are then computed, and a linear regression model computed to derive a scoring equation. Unmarked student essays are then processed to build the content summaries that are to populate the NP and VP content structures. Finally, the scoring equation is used to produce a score for each essay.

## Comparison of Human and MarkIT Scores

In a large scale test of MarkIT, 390 essays hand written by year 10 high school students on the topic of "The School Leaving Age" were transcribed to Microsoft Word document format. These essays were graded on a number of categories by three different human graders. The essays and scores were forwarded to the MarkIT development team at the School of Information Systems at Curtin University of Technology for processing. A model answer was chosen from amongst the essays by selecting the essay with the highest average score given by the three human graders – this essay had a score of 48.5 out of a possible 54, representing an overall score of 90%.
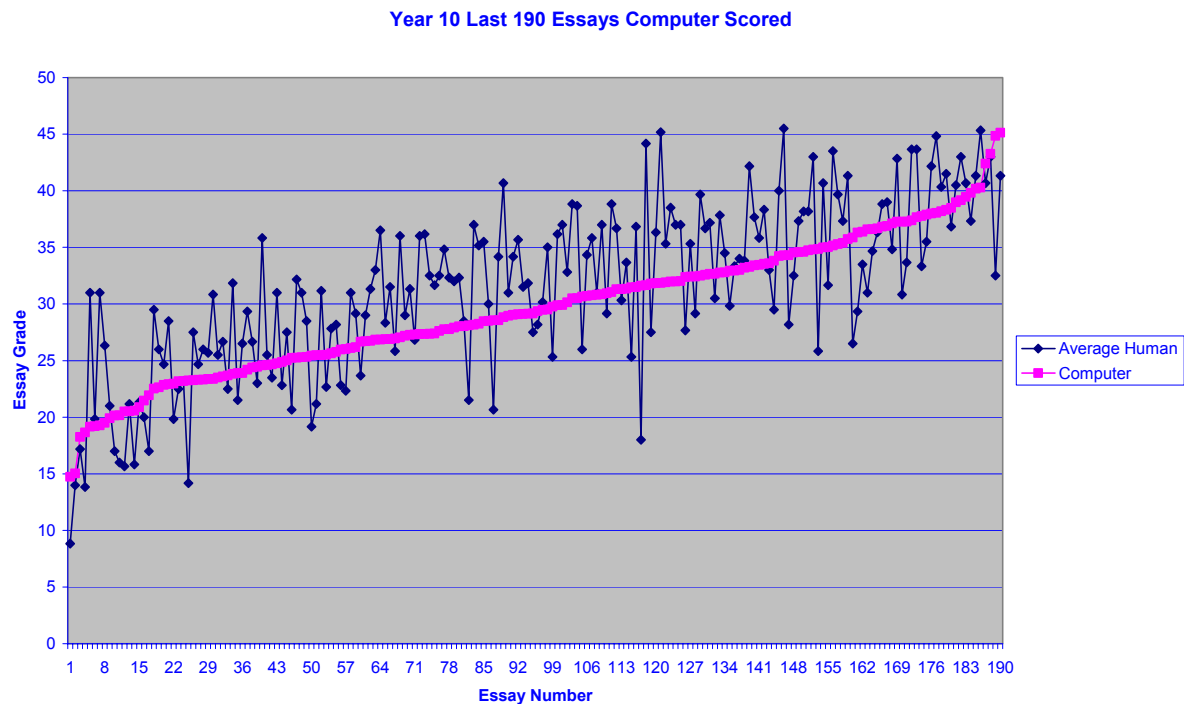
Figure 3 shows the variation amongst the first two graders on 200 essays. The essays scores are arranged in ascending order of one of the human assigned grades. Note the substantial disagreement in the scores for some essays.



**Figure 3: Comparison two Human Grader Scores on 200 Essays**

The mean score for Human1 for these 200 essays was 27.74, while the mean grade given by Human2 was 30.37, a difference of 2.63. The correlation between the two humans was 0.80. The mean absolute difference between the two was 5.22, representing an average error rate of 9.67% when scored out of 54 (the maximum possible human score). After a scoring algorithm was built using the 200 essays above as training data, the remaining 190 essays were scored

by MarkIT. Figure 4 shows the results, arranged in ascending order of the computer assigned score.

**Year 10 Last 190 Essays Computer Scored**



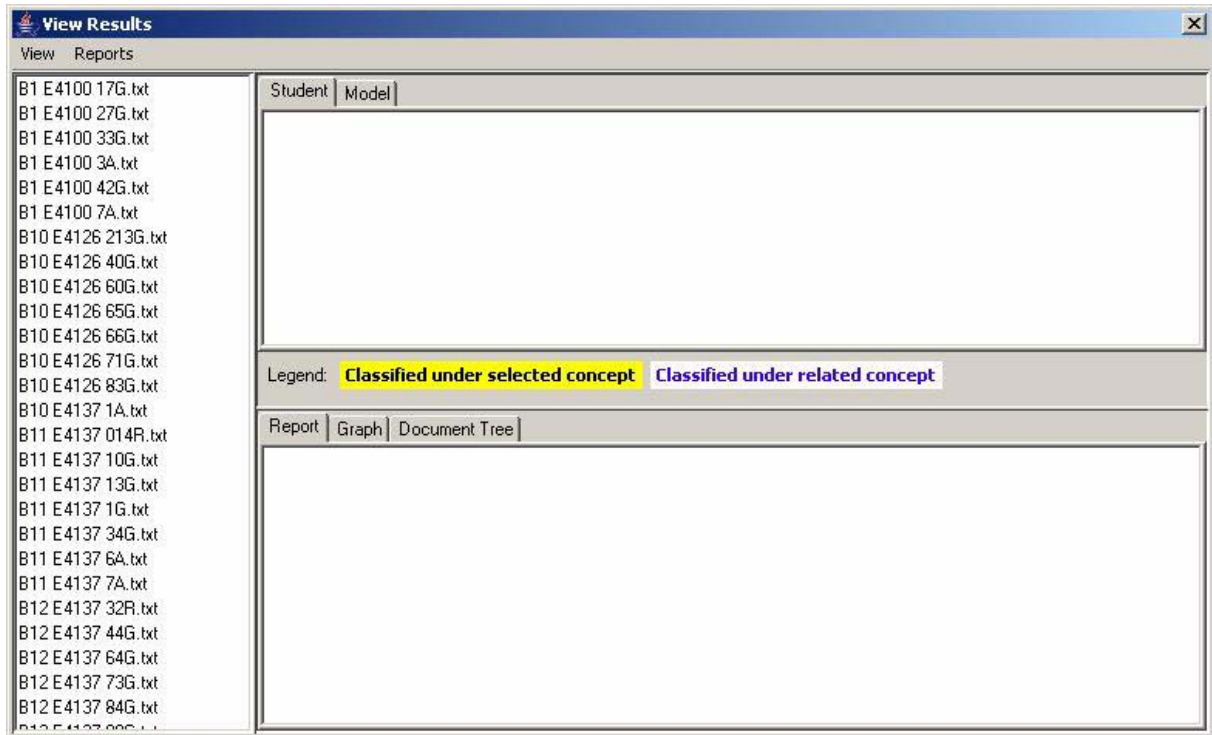**Figure 4: Results of Computer Scoring of 190 Essays vs Average Human**

The mean score for the human average grade for these 190 essays was 31.41, while the mean grade given by the computer was 29.62, a difference of 1.79. The correlation between the human and computer grades was 0.75. The mean absolute difference between the two was 4.39, representing an average error rate of 8.14% when scored out of 54 (the maximum possible human score).

The computer assigned scores were close to the agreement between the humans amongst themselves, and the error rates similar. We can conclude that in this particular test, MarkIT performed as well as human graders.

## Visual Feedback

A key strength of MarkIT over its counterparts is the emphasis on providing feedback other than just a grade or number. Numerous aspects of assessed assignments which are useful to students from an improvement point of view include: spelling, grammar, reading ease, and grade level statistics. Such data is derived from existing technology and is incorporated by MarkIT into a comprehensive Report on each assignment. Assignment content relative to model answer content is presented as a graph of concepts juxtaposing student answer concept content with model answer concept content. This graph is interactive in that one can drill down to the thesaurus level and also to the assignment level in order to discover where, and to some degree how, errors and omissions can be rectified.
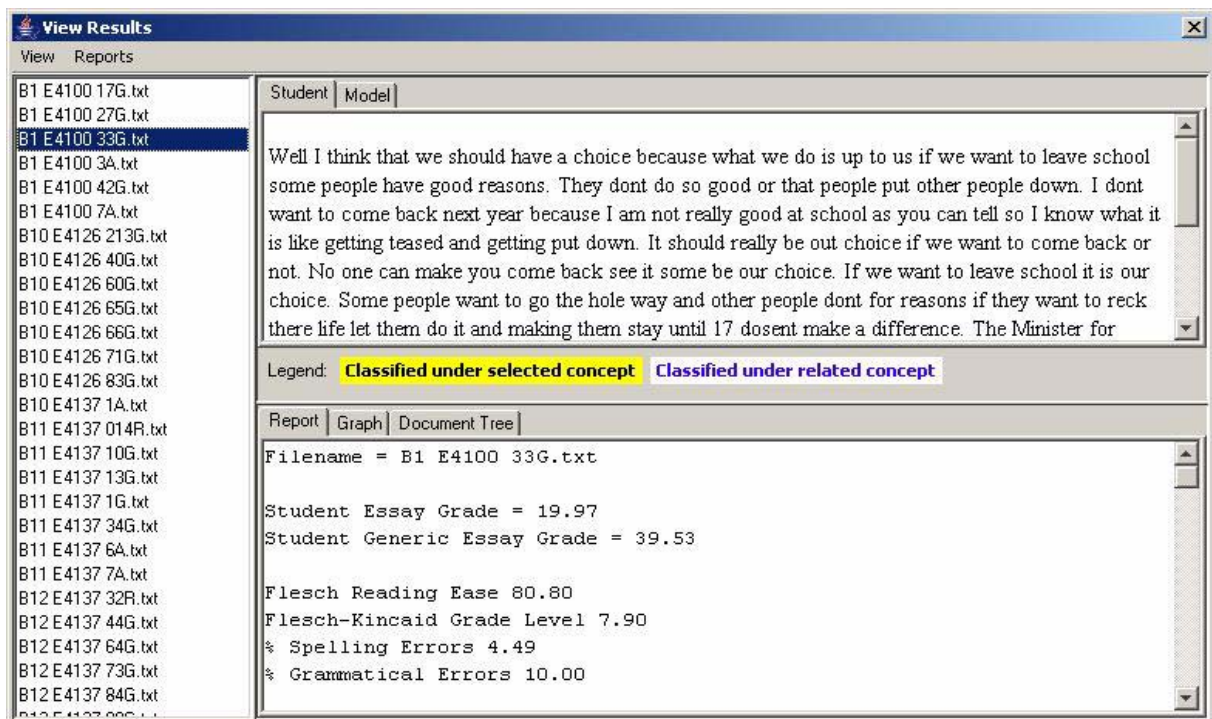
The following figures show MarkIT's visual feedback components and are available to the teacher and student on completion of the grading. Figure 5 shows the essay selection screen, with essay identifiers appearing in the left window.

**Figure 5: Main Control Panel for Visual Feedback**

When an essay identifier is selected, the screen shown in figure 6 results. The upper window can be toggled via the tabs to display the Student essay or the Model essay. The lower window can be toggled via the tabs to show further features:

- the grading Report (Figure 6),
- the Graph of the concept counts for Student and Model essays (Figure 7), and
- a Document Tree (Figure 9) representing the grammatical structure of the essay.



**Figure 6: Selected Essay Grading report**

Figure 7 presents a graph of the 'concepts' associated with both the model answer and the student answer. Naturally, the better the correspondence between the 'concepts' in both, the better the score. If we focus on the rightmost bar, labelled "Busyness", we see that the student answer contains a frequency of 2 (vertical axis) where the model answer called for no discussion on this topic or concept. We may say the student has introduced irrelevancies into the answer; or perhaps the student has waffled and provided 'filler'.

The concept labelled "+Being/In…" is a case where the model answer concept is not matched by an equal student contribution - this would correspond to a deficit in knowledge on the part of the student. The concept labelled "+Addition" is not matched by any student contribution – thus we may say the student is ignorant, or unaware of this concept or content. Such visual feedback is rather informative to student and teacher alike. The teacher is able to interactively explain to the student the strengths and weaknesses of the student's answer. If the teacher double clicks on a bar in the graph, the thesaurus text for the category represented by the bar is displayed (Figure 8). The student can then see the type of discussion that should have been devoted to the topic, and also get a good feel, from the many words in that category, how to express that content. The instructor can also switch to the model answer at any time to demonstrate the type of response that was expected.

In the upper window of Figures 7 and 9 one can see some highlighted or underlined words in the selected assignment or model answer. This feature visually marks the words in the essay which are associated with the 'concept' as selected by a user from the graph in the lower window. This matching of 'concept' with the words from the assignment which belong to it, is an excellent learning aid.
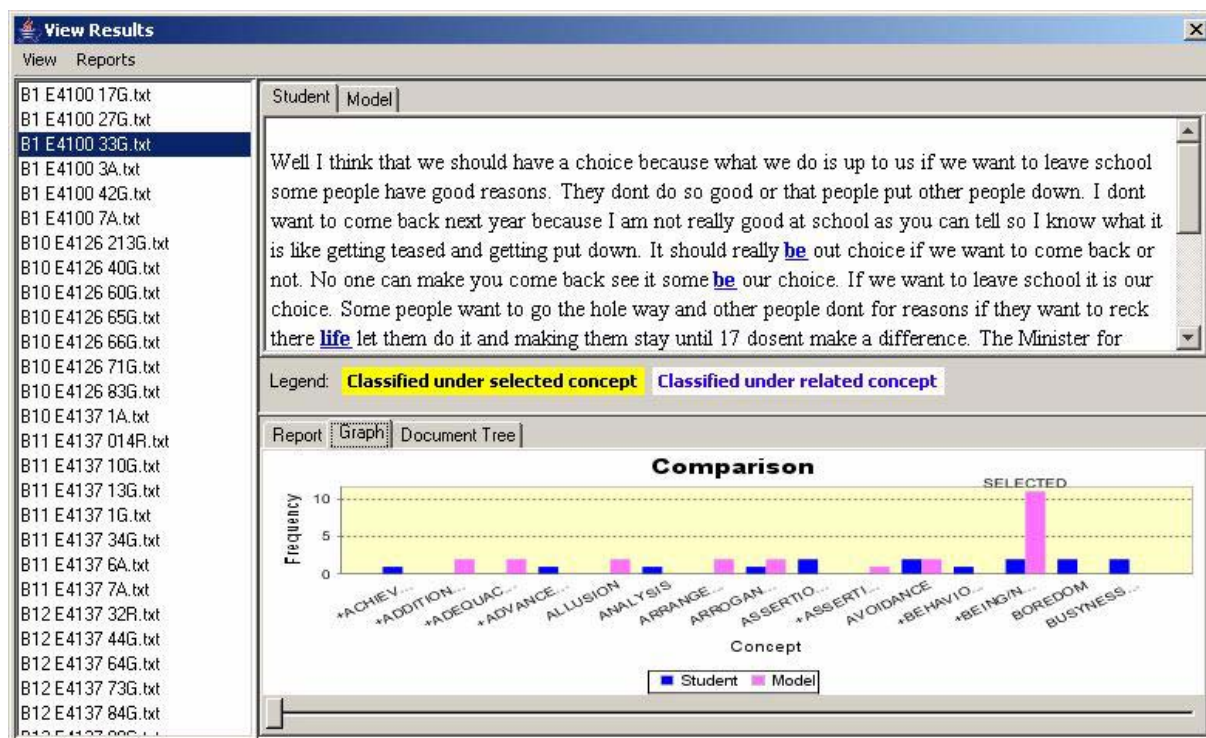
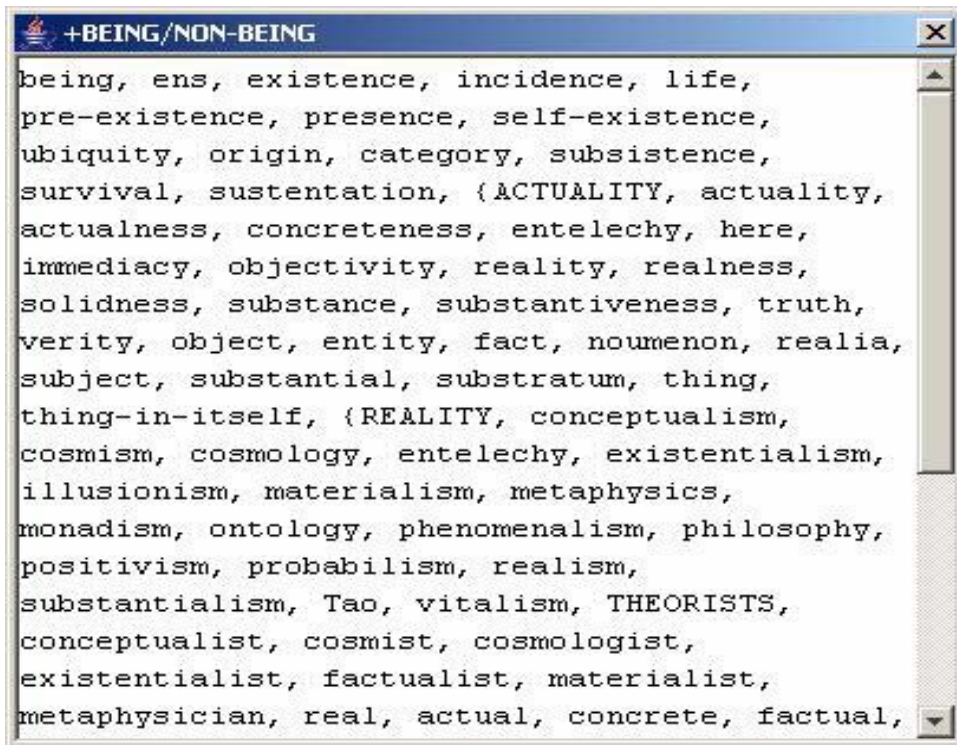

**Figure 7: Concept Frequencies**

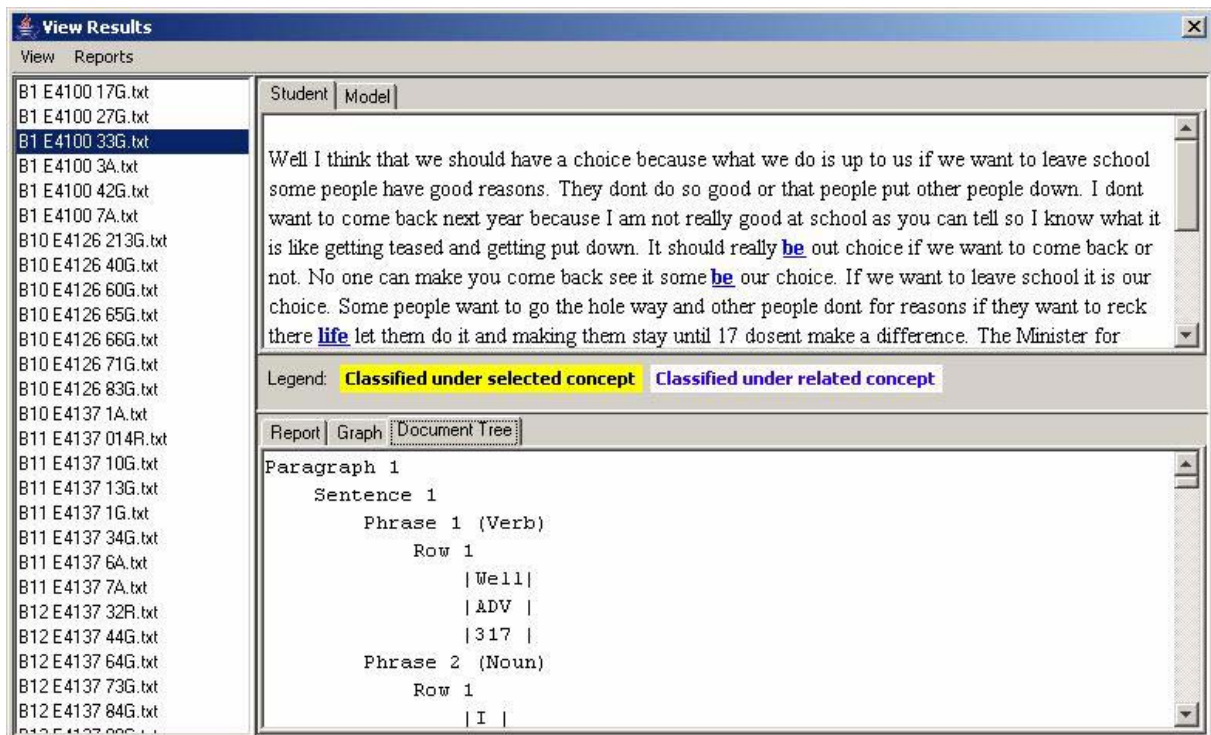**Figure 8: Thesaurus Entry for a Chosen Concept (+BEING/NON_BEING)**



**Figure 9: Document Tree (Semantic Content Structure)**

## MarkIT enabled E-Learning Systems

MarkIT is being developed so that all interaction between the teacher, student and the system can be performed over the Web and other networks. Teachers will be able to upload the model answer, training essays, and the student essays to be graded. The system will then undertake

the necessary processing to grade the submitted essays. Users will be notified within the hour that their results are available for viewing. Typically an essay can be graded in a few seconds on a standard workstation of year 2002. Output from the system, as described above, can then be accessed via the Web.

E-Learning systems in use today (e.g. Blackboard, WebCT, WBT-Master) often have good provision for assignment upload and some related assessment and management functions (for example as in WBT-Master, Dreher, H., Scerbakov, N., and Helic, D. 2004), but interactive feedback for the purpose of improvement as described above has not been implemented to any significant extent. Clearly such feedback is suitably provided in an on-line, real-time mode as and when the need exists. Connectivity to the Internet/WWW or indeed to wireless networks in a classroom setting would be a significant advantage to teacher and student alike. Students and teacher need to be connected, and one can readily imagine a class using wireless network technology to support the interaction with MarkIT in assignment upload and grading results visualization activities.

We believe that automated assignment marking coupled with interactive visual results feedback systems to be an important, effective and interesting use of telecommunications in the educational sector.

## Summary

MarkIT has been developed to provide automated grading of essay-type documents. Along with its peers in the AEG domain MarkIT performs as well as human graders under certain given conditions. Unlike many of its competitors, MarkIT is now endowed with the added feature of providing meaningful, relevant, and detailed interactive feedback to assist learners improve their performance. The system is being developed to provide all this functionality over the WWW and via wireless networks in classroom situations.

## References

Alter, S. (1992) *Information Systems A Management Perspective*, Addison-Wesley, Reading, Massachusetts.

Attali, Y. and Burstein, J. (2004) Automated Essay Scoring with E-rater V.2.0. Paper presented at the *Conference of the International Association for Educational Assessment (IAEA)*, June 13-18, Philadelphia, USA.
http://www.ets.org/research/dload/IAEA.pdf

Burstein, J., Kukich, K., Wolff, S., Lu, C., and Chodorow, M. (1998) Enriching Automated Essay Scoring Using Discourse Marking, *Proceedings of the Workshop on Discourse Relations and Discourse Markers, Annual Meeting of the Association of Computational Linguistics,* August, Montreal, Canada.

Dreher, H., Scerbakov, N., and Helic, D. (2004) Thematic Driven Learning, *Proceedings of E-Learn 2004 Conference*, Washington DC, USA, November 1-5.
http://www.aace.org/conf/ELearn/

Kolln, M. (1994) *Understanding English Grammar*, MacMillan, New York.

Landauer, T. K., Foltz, P. W., and Laham, D. (1998) An Introduction to Latent Semantic Analysis, *Discourse Processes*, 25, 259-284.

Landauer, T. K. (1999) Email communication with author, 8[th] June.

Larkey, L. S. (1998) Automatic Essay Grading Using Text Categorization Techniques, *Proceedings of the Twenty First Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, Melbourne, Australia, 90-95.

Page, E. B. (1994) Computer Grading of Student Prose, Using Modern Concepts and Software, *Journal of Experimental Education*, 62, 127-142.

Phelan, J. (2003) For Student Essayists, an Automated Grader. *The New York Times*, September 4
http://tech2.nytimes.com/mem/technology/techreview.html?res=9D0DE6D71038F937A3575AC0A9659C8B63

Poesio, M (2000) Semantic Analysis, in Dale, R., Moisl, H. and Somers, H. (editors) *Handbook of Natural Language Processing*, Marcel Dekker, New York

Shermis, M. and Burstein, J. (2003) *Automated Essay Scoring: A Cross-Disciplinary Perspective 2003*, Lawrence Erlbaum Associates, New Jersey.

Valenti, S., Neri F., and Cucchiarelli, A. (2003) An Overview of Current Research on Automated Essay Grading, *Journal of Information Technology Education*, Volume 2, pp319 – 330.
http://jite.org/documents/Vol2/v2p319-330-30.pdf

Williams, R. (2001) Automated Essay Grading: An Evaluation of Four Conceptual Models in Kulski, M. and Herrmann, A. (editors), *New Horizons in University Teaching and Learning: Responding to Change*, Curtin University of Technology, Perth, Australia

Williams, R. and Dreher, H. (2004) Automatically Grading Essays with Markit©. *Journal of Issues in Informing Science and Information Technology,* Vol. 1, pp693-700

**Biographies**



**Robert Williams** has over 25 year's experience in the Information Systems industry, as a practitioner, researcher and lecturer. He currently is a lecturer in the School of Information Systems at Curtin University of Technology in Perth, Western Australia. He has extensive experience in systems analysis and design, and programming, on a variety of mainframe, mini and personal computers, and a variety of operating systems and programming languages. Applications he has worked with include mathematical, statistical, bridge and road engineering, financial, corporate resource allocation, business simulation and educational systems. He has published a number of articles on system users' personalities and satisfaction, decision support systems, and automated essay grading systems. In 2001 he led a team of researchers in the School of Information Systems at Curtin University of Technology which conducted what is believed to be the first trial in Australia of an Automated Essay Grading system. Robert holds a Bachelor of Arts degree with double majors in Mathematics and Economics from the University of Western Australia, a Graduate Diploma in Computing from the Western Australian Institute of Technology, and a Master of Information Systems degree from Curtin University of Technology.



**Heinz Dreher** is senior lecturer and research fellow in the School of Information Systems at Curtin University of Technology. He has published in the educational technology and information systems domain through conferences, journals, invited talks and seminars; is currently the holder of Australian National Competitive Grant funding for a 4 year e-Learning project; is participating in a Digital Library project with TU Graz / Austria; is collaborating on Automated Essay Grading technology development, trial usage and evaluation; has received numerous industry grants for investigating hypertext based systems in training and business scenarios; and is an experienced and accomplished teacher, receiving awards for his work in cross-cultural awareness and course design. In 2004 he was appointed Adjunct Professor for Computer Science at TU Graz, and continues to collaborate in teaching and learning, and research projects with local and overseas partners.