Improving the standard and consistency of multi-tutor grading in large classes

Keith Willey

School of Computing and Communications, Faculty of Engineering and Information Technology, University of Technology Sydney, keith.willey@uts.edu.au

Anne Gardner

School of Civil and Environmental Engineering, Faculty of Engineering and Information Technology, University of Technology Sydney, anne.gardner@uts.edu.au

For several years the authors have coordinated a large engineering design subject, having a typical cohort of more than 300 students per semester. Lectures are supported by tutorials of approximately 32 students that incorporate a combination of collaborative team and project-based learning activities. Each tutor is responsible for grading the assessment tasks for students in their tutorial. A common issue is how to achieve a consistent standard of marking and student feedback between different tutors. To address this issue the authors have used a number of methods including double-blind marking and/or random re-marking to support consistent grading. However, even when only small variations between the overall grading of different tutors were found, students still complained about a perceived lack of consistency. In this paper we report on an investigation into the use of a collaborative peer learning process among tutors to improve mark standardisation, and marker consistency, and to build tutors' expertise and capacity in the provision of quality feedback. We found that students' perceptions of differences in grading were exacerbated by inconsistencies in the language tutors use when providing feedback, and by differences in tutors' perceptions of how well individual criterion were met.

Keywords: academic standards; marker moderation; self- and peer assessment; SPARK^{PLUS}.

Introduction

As a result of changes in the past two decades, Australian and UK universities have seen a reduction in staff—student ratios that have often resulted in large classes. Furthermore, research funds are often used to buy permanent academic staff out of teaching, resulting in an increasing number of less experienced casual or sessional teaching staff to conduct core teaching activities such as tutorials and marking of student work (Price, 2008; White, 2006).

Grading is an activity that often results in anxiety for both teachers and students, in part due to the difficulties in justifying grading decisions to students. This issue is further complicated in large classes by the fact that often a number of staff are involved in marking an assessment task. It is recognised that even experienced staff differ in their understanding of academic standards. The fact that marking is increasingly being completed by less experienced sessional teachers and tutors compounds this problem. These issues contribute to the fact that some students feel that grades are a function of whose tutorial they find themselves in, described by McCallum, Bondy and Jollands (2008) as "tutorial lotto".

Differences in academic grading may also result from lecturer self-interest given that grade leniency is a significant factor in positive evaluations of teaching (Greenwald & Gilmore, 1997; Marsh & Roche, 2000). Such activities contribute to unrealistic student expectations with respect to the marks they should be awarded. Furthermore, students' self-assessment of their work is often closely related to the time and/or effort required to produce the work rather than its quality.

Nesbit (2006) notes that students who have concerns over the justice of assessment report less satisfaction with the assessment task and the academic. Students with lower marks were found to more often have concerns over the justice of marking (believe the marking was fair) and "appeared to be more negatively affected by these concerns" (Nesbit, 2006, p. 668). Successfully addressing these justice concerns has potential to result in higher student satisfaction and improved performance on future assessment tasks (Nesbit, 2006).

ATN Assessment Conference 2010 University of Technology Sydney

While there are many reasons for variations in academic standards, in this investigation we focus on improving inter-marker standards and consistency, marker understanding of the assessment criteria and the quality of feedback provided to students.

Background

In an effort to achieve consistent grading between multiple markers, double-blind marking and/or re-marking a random selection of assessment tasks is often undertaken. However, with high student numbers and teaching loads these activities are fast becoming unrealistic.

For consistent marking between tutors it is important for all assessors to share a common view of the value of a given mark. Tomkinson and Freeman (2007) suggest that some form of induction, for example, a small number of 'yardstick' assessments be used as a basis for discussion about marking standards.

While many tutors try to provide helpful and detailed feedback, this practice is often inconsistent. Some academics appear to lack the knowledge of how to provide effective feedback; others may simply be overloaded and feel they cannot find the time to provide thoughtful feedback. Finally, there will always be those who remain cynical about the entire process of student learning in general and the purpose of feedback in particular (Weaver, 2006). Students report that unhelpful feedback includes comments which were too general or vague, lacked guidance, focused on the negative, or were unrelated to assessment criteria (Weaver, 2006). Even when quality feedback is provided, students may need advice on understanding and using the feedback before they can engage with it.

Several researchers report combining software tools and a marker meeting to increase student satisfaction with the feedback they receive and to improve the consistency and efficiency of grading, particularly in large classes (Debuse, Lawley, & Shibl, 2007; Moodie, Brammer, & Hessami, 2007; Anglin, Anglin, Schumann, & Kaliski, 2008; Lilje, Breen, Lewis, & Yalcin, 2008). These researchers concluded as did Price that "an assessment standards discourse is needed to support the functioning of assessment communities of practice" (Price, 2008, p. 226). That is, tutors develop their understanding of the assessment criteria and language of feedback by discussing marking with other academics. This aligns with a social constructivist view of learning, that is, learning requires "active engagement and participation" this being true for tutors no less than for students (Rust, O'Donovan, & Price, 2005, p. 237).

With a plethora of purpose-designed software tools available, academic staff run the risk of suffering from software overload. The authors' view is that any innovation is simpler to introduce if it uses software that staff are already familiar with. At the University of Technology, Sydney, the software tool SPARK^{PLUS} (see acknowledgement) is regularly used to facilitate self- and peer assessment tasks. SPARK^{PLUS} has the capacity not only to assess a student's contributions to a team project, but also to allow students to self- and peer assess individual work and improve their judgment through benchmarking exercises (Willey & Gardner 2009a, 2009b, 2009c).

The tutors in the large class chosen for this study were already familiar with this tool (students and subsequently tutors used it four times for various activities within the class) making it ideal to facilitate our tutor development activities. Our aim was to increase markers' confidence and improve the quality and consistency of the feedback markers provide to students. In this paper we investigate the impact of a tutor benchmarking and discussion exercise on developing an agreed marking standard, increasing consistency of marking and feedback to students and tutor engagement.

Research method

The research was conducted in the subject Design Fundamentals in the Spring semester of 2009. Design Fundamentals is a second-year engineering degree course at the University of Technology, Sydney, with typical enrolments of more than 300 students per semester. These students are distributed among tutorials which support the lecture program. Each tutorial has a maximum of 32 students. Individual tutors mark assessment tasks for students in their tutorial. The authors had previously reported on the results of a tutor benchmarking exercise which was undertaken in the Autumn semester of 2009 when one of the authors was course coordinator and lecturer in this subject. The authors are aware that the success reported in this previous research may have been influenced by their involvement with the subject. The authors wanted to further investigate and verify their previously reported

improvements in marking standardisation by observing a third-party implementation of the process. In Spring semester 2009 the authors were not involved in the teaching of the subject Design Fundamentals. The new lecturer (teaching the subject for the first time) implemented the tutor benchmarking exercise described in Willey and Gardner (2010). The authors observed the activity and had access to the research instrument and student results for the semester. Of the nine tutors involved in the subject that semester, six agreed to participate in the exercise. Of these six tutors, two were tutoring the subject for the first time, two for the second time, one for the third time and one for the fifth time

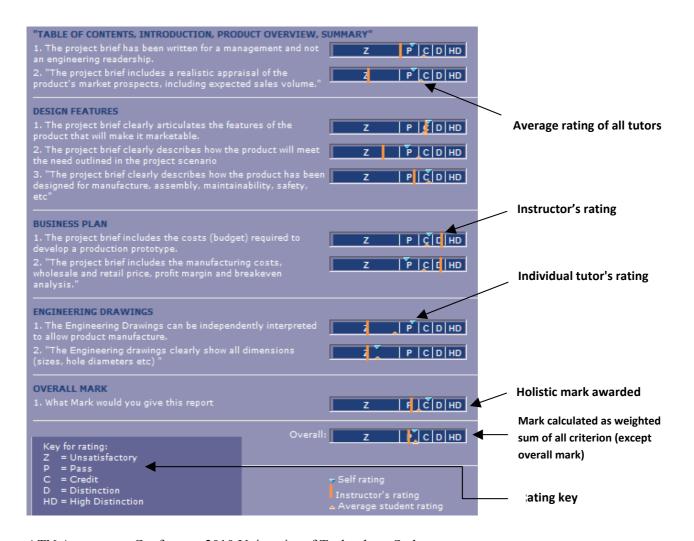
The investigation had a number of stages, as follows.

Stage 1

The average marks for each tutorial for task one (Requirement Specification Report) were calculated and compared. This task was marked by the tutors before they participated in the benchmarking exercise. However, prior to marking, tutors were given two exemplar marked reports with feedback comments by the subject coordinator as examples of the expected marking standard.

Stage 2

Before tutors marked task two (Project Brief Report) they were provided with a copy of two student reports from the current semester (one of satisfactory quality and the other of higher quality). Tutors graded these reports against specified criteria (the same criteria they would be using for marking their tutorials' reports) and then entered their assessment (grading and feedback comments) into a benchmarking task in SPARK (Figure 1). The lecturer also entered his assessment (grading and feedback comments) into SPARK to allow comparison with the tutors' grading and feedback.



ATN Assessment Conference 2010 University of Technology Sydney

Note: Boxes containing written feedback are not shown.

Figure 1: Benchmarking results screen in SPARK PLUS

Stage 3

All tutors gave their permission for their SPARK^{PLUS} ratings to be shared with the other tutors. The authors believe this acknowledgement and respect for the individual helps to create an environment that facilitates peer learning, one where differences in thinking can be safely explored to progress and consolidate the learning of all participants.

Part A

At a tutor meeting, tutors were asked to discuss their individual grading of the two reports (previously recorded in SPARKPLUS) paying particular attention to the reasons why they awarded a particular grade and exploring any differences in grading between them.

Part B

The lecturer explained the thinking behind their marking of the reports as recorded in SPARKPLUS.

Part C

Where significant differences in opinion were evident an individual tutor's grading as recorded in SPARKPLUS was then displayed on the screen as a focus for discussing these differences. SPARKPLUS allows a tutor to compare their holistic mark and rating on individual criterion with both the lecturers and the average tutors' mark/rating.

Part D

Tutors were asked to compare and assess and discuss the difference in the feedback comments provided by individual tutors. The results of this step are not reported in detail in this paper.

Part E

Subsequently, all participants collaboratively re-graded the reports, i.e. they were required to reach a consensus about the appropriate grade for each assessment criterion, the reasons for this grade, the associated feedback comments they would provide students and agree on an overall holistic grade.

Stage 4

Tutors were asked to complete a survey instrument containing a combination of Likert scale and free response questions. The instrument included questions to ascertain tutors' confidence in marking the report and whether they thought they marked harder or easier than the lecturer, and harder or easier than the average of all the tutors.

Stage 5

After the meeting tutors could logon to SPARK^{PLUS} and see how their individual grading of the reports for each criterion and their feedback comments compared with the lecturer's and the average tutor assessment (Figure 1), i.e. they were able to benchmark their judgement against the lecturer and the average of all the tutors. Tutors were encouraged to reflect on the differences before marking their tutorial's reports.

Stage 6

Tutors marked task two and the average marks for each tutorial were compared to the average marks for all tutorials. The marks for task two were also compared to the marks for task one, to determine the impact of the benchmarking exercise.

Throughout the tutors' meeting both authors recorded their observations of the process and tutor discussions.

Results

The results from the collected data are recorded in Tables 1–3.

Findings

Table 1 shows that while the average mark for the whole class, i.e. all tutorials, remained the same from task one to task two, the variance between the average marks for each tutorial reduced significantly from 13 to four. This suggests the success of the benchmarking activity in improving mark standardisation in that there was far less variability in the marking between different tutors for task two after the moderation activities compared to task one.

Table 1: Average mark for each tutor's tutorial for each report and standard deviation of the means

	Average task 1 (requirement specification) mark per tutorial (%)	Task 1 (requirement specification) marking within 1 STD of average of all submissions	Average task 2 (design brief) mark per tutorial (%)	Task 2 (design brief) marking within 1 STD of average of all submissions	
Tutor 1	65	Yes	67	Yes	
Tutor 2	79	Yes	74	Yes	
Tutor 3	68	Yes	66	No lower	
Tutor 4	56	No lower	61	No lower	
Tutor 5	72	Yes	78	No higher	
Tutor 6	76	Yes	68	Yes	
	Requireme	nt specification	Design brief		
All individual submissions	Mean 70	STD 13	Mean 70	STD 4	

The average marks for task one (requirements specification) and task two (design brief) reports for each tutorial were compared to the average mark for the whole class. While the authors acknowledge that the average ability of students may vary from one tutorial to another, for the purpose of this study we have used the following definition: if the average mark for a particular tutorial was more than one standard deviation from the class average, that tutor was regarded as marking higher (average mark more than one standard deviation higher than class average) or lower (average mark more than one standard deviation lower than the class average).

Table 2: Response of each tutor's perceptions of their marking collected using the survey instrument

Survey question	Tutor 1	Tutor 2	Tutor 3	Tutor 4	Tutor 5	Lecturer/ Tutor 6
I think I mark harder than the average tutor	Slightly agree	Slightly agree	Slightly agree	Agree	Slightly agree	Slightly agree
Students in my tutorial think I mark harder than the average tutor	Slightly agree	Slightly agree	Slightly agree	Strongly agree	Slightly agree	Slightly agree

Despite, as reported in Table 2, all six tutors thinking they marked harder than the average tutor and their students thinking they marked harder than other tutors, the evidence suggests most tutors had an inaccurate perception of their marking. Tutor four, who was the most experienced, was found to provide lower marks on average for task one and task two (Table 1). Hence, their perception that they marked harder may be correct. All other tutors on task one were found to have marked relatively the same. On task two, tutor three and tutor four were found to mark harder, while tutor five was found to mark easier. The figures for task two, while indicative, need to be treated cautiously as the standard deviation used as a metric was relatively small (four) due to the reduced variation in the marking of different tutors after the benchmarking activity.

Table 3: Tabulation of results collated from individual tutor's grading entered into SPARK^{PLUS}

Report 1	Tutor 1	Tutor 2	Tutor 3	Tutor 4	Tutor 5	Lecturer/ Tutor 6
Criterion rating the same as average of all tutors	3	8	7	5	7	3
Criterion rating different to average of all tutors H (higher), L (lower)	6Н	1L	2Н	4L	1H, 1L	2H, 4L
Mark calculated from independently weighted criterion	64	61	70	51	67	53
Final holistic mark awarded	74	73	72	52	57	60
						Lecturer/
Report 2	Tutor 1	Tutor 2	Tutor 3	Tutor 4	Tutor 5	Tutor 6
Criterion rating the same as average of all tutors	8	6	7	7	6	7
Criterion rating different to average of all tutors H (higher), L (lower)	1L	3L	1H, 1L	1H, 1L	3Н	1H, 1L
Mark calculated from independently weighted criterion	55	50	63	63	71	66
Final holistic mark awarded	54	53	65	76	72	65

There were nine criteria recorded in Figure 1 that were used to assess the two reports in the benchmarking activity. The same standard deviation test was used to compare the differences in tutor marking for each criterion used to mark the reports in the benchmarking activity (Table 3). It is interesting to note that there was more variation (different understanding of what was required) when marking the first report, the poorer of the two, than the second report. In marking report one, three tutors were more than one standard deviation from the average mark on at least four criteria (tutor one was higher on six criteria (6H), tutor four lower on four criteria (4L) and tutor six was high up on two criteria and lower on four (2H, 4L)). While in marking report two (the better report) no tutors differed on more than three of the nine criteria. While there is insufficient evidence to draw solid conclusions, the results indicate that tutors may find it easier to be more consistent in judging the standard of higher quality submissions.

The exercise also uncovered a number of other interesting results. Firstly, tutors one, two and three regarded the first report as being the better of the two and hence awarded it the highest holistic grade. This was contrary to tutors four, five and six who felt the second report was of a higher standard. On analysing these results the authors first suspected that tutors one, two and three may have marked the reports in the wrong order, however, on checking, this was not the case. It was during the discussion part of the activity when individual tutor's SPARK grading was displayed to the whole group that it become evident there were vastly different understandings of what was required in respect to each criterion. For example, a significant difference was apparent between the grading of tutor four and the grading of the other tutors for the criterion within the Engineering Drawing category for report one. Tutor four rated the drawings as being extremely poor and a significant contribution to the poorer mark they awarded to this report. It turned out that as an experienced structural engineer with industry experience, tutor four had the expertise to appreciate the deficiencies in the drawings and was able to clearly explain this to all the tutors during the benchmark activity, informing the thinking of the other tutors as well as the lecturer. This resulted in the other five tutors increasing their knowledge of engineering drawings and their understanding of what was ATN Assessment Conference 2010 University of Technology Sydney

important. A similar situation occurred with tutor one, who gave the highest mark for report one (marking significantly higher on six of the nine criteria). Their main discrepancy with the other tutors was in the quality of the business plan and appraisal of the product's market prospects (see Figure 1 for criteria). Tutor four was able to clearly and comprehensively articulate their reasons for marking as they did, after which, all the remaining tutors agreed they had not appreciated certain aspects of the report. Interestingly, tutor four was the only tutor who was not an engineer, but rather had an information technology degree and a good deal of administrative experience. The authors observed many similar conversations where the tutors inclusively and respectfully debated their different understanding of the assessment criteria, what was expected and the standard of work required to achieve different grades. These conversations were an integral part of the collaborative peer learning exercise. The tutors were reassured the intention of the exercise was to improve their knowledge, understanding, skills, confidence and the standard of their marking, and was not to judge their performance. Like students, the tutors needed to approach this exercise as an opportunity to learn and improve skills rather than a summative assessment activity.

The authors observed that tutors actively discussed the assessment of the reports during the meeting with comments such as "I marked too hard" and "I missed that", showing that they felt comfortable enough to acknowledge that they didn't identify something or understand something as well as somebody else. The conversations assisted them to discover and address what they didn't understand about the criteria. The discussion part of the benchmarking exercise provided an opportunity for all tutors to learn from the expertise and experience of each other.

The tutor grading patterns and feedback recorded in SPARK^{PLUS} for the two reports marked as part of the benchmarking exercise were further analysed. The analysis showed there were frequently large differences in the total grade awarded from the individual weighted criterion and the grade awarded holistically. For example, when marking report one, tutor one awarded a holistic final mark of 74 per cent, while the weighted aggregate of the individual criterion equated to a mark of 64 per cent (see Table 3). A similar situation occurred in tutor four's marking of report two, where they holistically awarded a mark of 76 per cent, while the weighted aggregate of the individual criterion was only 63 per cent. Conversely, the holistic mark and weighted aggregate of the individual criterion for both reports marked by tutor three were nearly the same. The fact that students may receive feedback against individual criterion that does not necessarily reflect their awarded grade contributes to student's dissatisfaction with marking consistency (Willey & Gardner, 2010).

Discussion

In response to the marking of individual criterion not always reflecting the overall grade awarded, the authors have changed their approach to marking. Previously, when coordinating subjects with multiple tutors a detailed marking rubric was provided. Tutors awarded marks in accordance to the rubric which was handed back to students. While this provided comfort to the tutors in explaining their marks to students it also resulted in students' arguing about small differences in the marks awarded on individual criterion between different reports. Students became focused on their marks, using any perceived differences to argue for a mark increase rather than reflecting on the overall quality of their submission and understanding its strengths and weaknesses.

The authors now use such rubrics to provide feedback on the standard of a student's work (was the work against each criterion unsatisfactory, satisfactory, credit worthy, distinctive or highly distinctive, etc.) rather than being used to calculate their final grade. Final grades are awarded on a holistic basis using academic judgement. Students who wish to argue their grade are required to justify their arguments holistically rather than focusing on one mark here or one mark there. This has been effective in moving students to focus on learning and the strengths and weaknesses of their contribution rather than marks. For example, where students who were awarded a mark in the range of 46 to 49 per cent would once have focused on arguing for any marks to achieve a pass, now they have to focus on explaining why their work is of satisfactory quality and worthy of a pass.

In addition, as previously discussed, even when tutors awarded the same grade, there were often substantial differences in the language used by different tutors to provide feedback. A major benefit of the described benchmarking and discussion process is that active participation in the associated discussions promotes the use of more consistent assessment language by tutors when providing feedback to students to explain their grades. In the reported exercise all but one tutor (the most experienced having tutored the subject five times) reported that the activity enabled them to give more consistent and higher quality feedback to their students. Not only had the tutors developed a better understanding of the associated learning outcomes and what was required to meet them, they developed a clear language to explain both the strengths and weaknesses of individual submissions to students.

A notable advantage of this benchmarking and discussion process over other reported marker meetings, as described in (Lilje et al., 2008), is that the tutors had to exercise and record their judgement (assessment) of the reports in SPARK PLUS before they came to the meeting. In previous semesters, when a multiple-blind marking process was used, tutors would bring their marked reports to the tutors' meeting to discuss. However, the authors observed that when a dominant tutor, or a tutor who was seen to have more authority, stated a certain position, less experienced or confident staff were reluctant to express a contrary opinion. This stifled the discussions, reducing the learning of the exercise participants. Having the results recorded in SPARK prior to the exercise means that all differences in opinion are available for comparison and discussion. Furthermore, tutors could continue to log onto SPARK after the exercise to compare their marking and feedback to the other tutors' and lecturer's. The authors cannot stress enough that for this exercise to be successful care must be taken to ensure the environment is safe and respectful of all tutors in that differences in opinions are not seen as being correct or incorrect but rather as opportunities to learn.

Conclusions

The reported benchmarking activity was effective in reducing the variability in marking between different tutors in a large class. The process promoted inclusiveness of less experienced and less confident tutors by using a software tool to record tutor assessments and feedback before exploring their understanding in a subsequent discussion activity. This tool also encouraged ongoing reflection by allowing individual tutors to compare their grading for each criterion and their feedback comments to those of the course coordinator and the average tutor assessment.

The tutors used the benchmarking exercise as a way to calibrate their marking. Furthermore, it was particularly effective in helping to develop feedback skills and shared understandings of marking standards in less experienced tutorial staff. We also found that while students' perceptions of difference in grading between tutorials were not unfounded, the problem was exacerbated by inconsistencies in the language tutors use when providing feedback. Hence, to improve student satisfaction with marking justice it is necessary to improve the consistency of both the applied marking standard and the feedback language used to explain student grades. Furthermore, the benchmarking activity proved to be robust in that it provided measurable improvements when it was run by a relatively inexperienced third-party – a lecturer teaching the subject for the first time.

Finally, our findings support the conclusions of other researchers that found that conversations about assessment standards and marking is an effective method of developing a shared understanding of assessment criteria and improving the standard, consistency and student satisfaction with marking by multiple academics on an assessment task.

Acknowledgement

SPARK^{PLUS} is a joint research project between the University of Technology, Sydney and the University of Sydney. This research was supported through a UTS Learning and Teaching Performance Fund Grant.

References

- Anglin, L., Anglin, K., Schumann, P., & Kaliski, J. (2008). Improving the efficiency and effectiveness of grading through the use of computer-assisted grading rubrics. *Decision Sciences Journal of Innovative Education*, 6(1), 51-73.
- Debuse, J., Lawley, M., & Shibl, R. (2007). The implementation of an automated assessment feedback and quality assurance system for ICT courses. *Journal of Information Systems Education*, 18(4), 491-502.
- Greenwald, A. G., & Gilmore, G. M. (1997). 'No pain no gain': the importance of measuring course workload in student ratings of instruction. *Journal of Educational Psychology 89*(4), 743-751.
- Lilje, O., Breen, V., Lewis, A., & Yalcin, A. (2008). The structure, use and impact of the staff version of ORWET. In A. Hugman & K. Placing (Eds.), *Symposium Proceedings: Visualisation and Concept Development* (pp. 188-192). UniServe Science, The University of Sydney.
- Marsh, H. W., & Roche, L. A. (2000). Effects of grade leniency and low workload on students' evaluations of teaching: popular myth, bias, validity or innocent bystanders? *Journal of Educational Psychology* 92(1), 202-228.
- McCallum N., Bondy J., & Jollands M. (2008). Hearing each other how can we give feedback that students really value. *Proceedings of the Nineteenth Annual Conference of the Australasian Association for Engineering Education*. Yeppoon: Faculty of Sciences, Engineering & Health, CQU University.
- Moodie J., Brammer, N., & Hessami M. (2007). Improving written communication skills of students by providing effective feedback on laboratory reports. *Proceedings of the 2007 AaeE Conference*. Melbourne.
- Nesbit P., & Burton S. (2006). Student justice perceptions following assignment feedback. *Assessment & Evaluation in Higher Education*, 31(60), 655-670.
- Price M. (2005). Assessment standards: the role of communities of practice and the scholarship of assessment. *Assessment & Evaluation in Higher Education*, 30(3), 215-230.
- Rust, C., O'Donovan B., & Price M. (2005). A social constructivist assessment process model: how the research literature shows us this could be best practice. *Assessment & Evaluation in Higher Education*, 30(3), 231-240.
- SPARK PLUS, Retrieved August 6, 2010, from http://spark.uts.edu.au.
- Tomkinson, B., & Freeman, J. (2007). Using portfolios for assessment: problems of reliability or standardisation? *Enhancing Higher Education, Theory and Scholarship, Proceedings of the 30th HERDSA Annual Conference*. Adelaide.
- Weaver, M. R. (2006). Do students value feedback? Student perceptions of tutors' written responses. *Assessment & Evaluation in Higher Education*, 31(3), 379-394.

- White, N. R. (2006). Tertiary education in the noughties: the student perspective. *Higher Education Research & Development*, 25(3), 231-246.
- Willey, K., & Gardner, A. (2009a). Improving self- and peer assessment processes with technology. *Campus-Wide Information Systems*, 26(5), 379-399.
- Willey, K., & Gardner, A. (2009b). Developing team skills with self- and peer assessment: Are benefits inversely related to team function? *Campus-Wide Information Systems*, 26(5), 365-378.
- Willey, K., & Gardner, A. (2009c). Changing student's perceptions of self and peer assessment. *Proceedings of the Research in Engineering Education Symposium*. Queensland.
- Willey, K., & Gardner, A. (2010). Perceived differences in tutor grading in large classes: fact or fiction? *Proceedings of the 40th ASEE/IEEE Frontiers in Education Conference*. Washington, DC.