# Honour thesis assessment: the role of guidelines in achieving inter-rater agreement

**Alan Henderson**
School of Engineering, University of Tasmania, Hobart, Australia
alan.henderson@utas.edu.au

**Marcus Guijt**
School of Chemistry, University of Tasmania, Hobart, Australia
marcus.guijt@utas.edu.au

**Michael Breadmore**
School of Chemistry, University of Tasmania, Hobart, Australia
michael.breadmore@utas.edu.au

**Anna Carew**
Australian Maritime College, University of Tasmania, Launceston, Australia
anna.carew@utas.edu.au

**Rosanne Guijt**
School of Chemistry, University of Tasmania, Hobart, Australia
rosanne.guijt@utas.edu.au

*Abstract: Engineering honours projects are often self selected or negotiated by students to be in areas of their own interest. While this encourages motivation and engagement in a self directed research project, it also considerably increases the diversity of honours project types. Such project diversity raises questions about the most suitable form of guidelines to provide good inter-rater agreement. Project diversity can also result in academics assessing projects outside their primary area of research specialisation. This is particularly true of transdisciplinary projects that cross over conventional discipline boundaries. To investigate these issues, a team of academics from two Schools at the University of Tasmania assessed a collection of engineering theses using a variety of different guidelines. All guidelines were found to produce poor inter-rater agreement, however inter-rater agreement was improved when both assessors were of the same discipline. An account of the academics' comments on use of the guidelines reveals conflicting opinions of good and bad features. Guidelines that were viewed as easy to use and less subjective were found not to substantially improve inter-rater agreement. The implications of these findings are discussed and suggestions made in relation to improving assessment guidelines for honours theses.*

## Introduction

Honours projects in the engineering discipline are often regarded as the capstone of the degree program. Students are required to work with a high degree of independence on a major research project which requires them to demonstrate many skills that they have developed throughout the degree. The honours mark is an important result for both graduate and postgraduate opportunities.

Assessment of project based work of this nature is difficult due to the diverse nature of projects undertaken. Honours projects commonly fall into one of three main types (Littlefair and Gossman, 2008): innovative projects that are proposed by a supervisor, often linked to particular interests of a supervisor; original projects proposed by students; and standard problems that have been studied before and have a known result or proven methodology. In addition to these types, projects may also

fall into any one of Boyer's categories of academic scholarship (Boyer, 1997): discovery, integration, application, and teaching.

The process of assessing honours theses can vary greatly across institutions and even among schools of the same institution (Hannan *et al.*, 2009). However, in most cases the final honours mark receives a substantial contribution from assessment of the thesis component alone. This mark is usually derived by combining independent marks from a team of assessors. This team may consist of the primary supervisor, co-supervisors, industry supervisors, and/or other staff members not involved with the project.

Assessment of honours theses requires a framework that does not restrict the type of scholarship, and should be assessed in a fair and transparent manner. Tariq *et al.* (1998) explains that honours marks should reference learning attributes and not external factors beyond students' control. They also note that project work requires a wide range of skills such as problem-solving, time management, project management, information retrieval, that are difficult to isolate. Originality and creativity are highly sought after attributes but are also difficult to assess.

In an attempt to improve learning and teaching, Universities are increasingly embracing criterion based assessment (CRA), where students are able to clearly identify the performance standards by which they will be assessed (Tariq *et al.*, 1998). Marks are assigned based on performance relative to these standards, and not by making reference to the performance of their peers, such as in norm referenced assessment (NRA). Tariq *et al.* (1998) gives an account of the introduction of a CRA system to assess honours theses and found numerous benefits including increased objectivity through reference to a number of clearly defined explicit criteria, improved feedback to students, flexibility with regard to weightings, and ease of use. Tariq *et al.* (1998) also noted that NRA is less transparent, and can be problematic in assessing a diverse range of projects. Criterion referenced assessment is not without criticism and some educators suggest that it can limit student experimentation, creativity and originality (Hay, 1995). Biggs (2003) also notes that NRA has deep roots within the higher education system. For example, the term *high distinction* is a comparative term that refers to the 'few that are highly distinguished'. Tariq *et al.* (1998) commented that devising a perfect assessment strategy is an elusive pursuit, but the introduction of CRA system was viewed as a positive improvement in teaching and learning.

The motivation for this study stems primarily from an interest in transdisciplinary projects, which are not uncommon in the engineering discipline. Transdisciplinary research projects are increasingly being encouraged by institutions as it is recognised that they present new avenues for innovation (Wall and Shankar, 2008). Assessment is known to be problematic as projects span across discipline boundaries. In particular the paper aims to investigate several research questions. What form of guidelines improve inter-rater reliability? What qualities of assessment guidelines are viewed as useful by assessors? What is the effect of mixed assessor disciplines on the inter-rater reliability? The inter-rater reliability of rubrics itself is rarely assessed (Stellmack *et al.*, 2009). This study investigates these issues through assessment of a collection of engineering theses using a number of different guidelines. An account of the academics' views is presented and then suggestions for developing improved guidelines.

## Research Methodology
In this study a team of academics from the University of Tasmania assessed a set of previously marked engineering honours theses. The team consisted of two academics from the Engineering discipline and three from the chemistry discipline. An attempt was made to select a diverse set of engineering honours theses in terms of Boyer's scholarship type, but interestingly, none could be identified in the pure discovery area. Projects of a technical nature often displayed mixed elements of integration, application and discovery, but not a clear designation into one category. The resulting set of theses may be best described as four science-focused and two teaching-focused.

A diverse set of five assessment guidelines were chosen so that each thesis would be assessed in a range of different ways. The guidelines, designated by letters A to E, are presented in Appendix 1, and have been used with permission of the authors. Authors and institutions have not been identified to maintain anonymomity in this study. Guidelines A to D are presented as rubrics, Guideline E could

be described as a single criterion rubric. Guideline E makes reference to NRA, whereas all others are clearly CRA.

Each thesis was examined by a team of two assessors working independently. Some teams were composed of assessors from the same discipline, termed 'same-field', others were 'cross-field'. Cross-field assessor teams are commonly used when assessing projects of a transdisciplinary nature.

The individual results from each assessor were then processed into grade bands according to Guideline C, and the assessment results from each team examined for a agreeing or disagreeing grade result. Disagreeing results were broken down into sub groups of disagreement by one grade and disagreement by two grades. An additional algorithm was applied to Rubric B to determine a grade result from the set of criteria results.

There were two main reasons for using this approach. First, processing the results into grade bands allowed a direct comparison with Guidelines B and C, as these produced grades and not marks. Second, while the small sample size did not allow for a quantitative statistical analysis, it did provide an indication of inter-rater agreement. It is recognised that a disagreement may stem for a relatively small inter-rater difference near grade boundaries. However, these results are only used to indicate basic trends and are discussed for the purpose of making suggestions on how to improve guidelines.

## Results

The results presented in the left part of Figure 1 show that all assessment guidelines produced a poor level of inter-rater agreement. There is a general trend of more assessor team disagreements than agreements, with exception to Guideline C, where more teams agreed than disagreed on grades. Assessment Guideline E had by far the worst level of grade agreement– none at all.
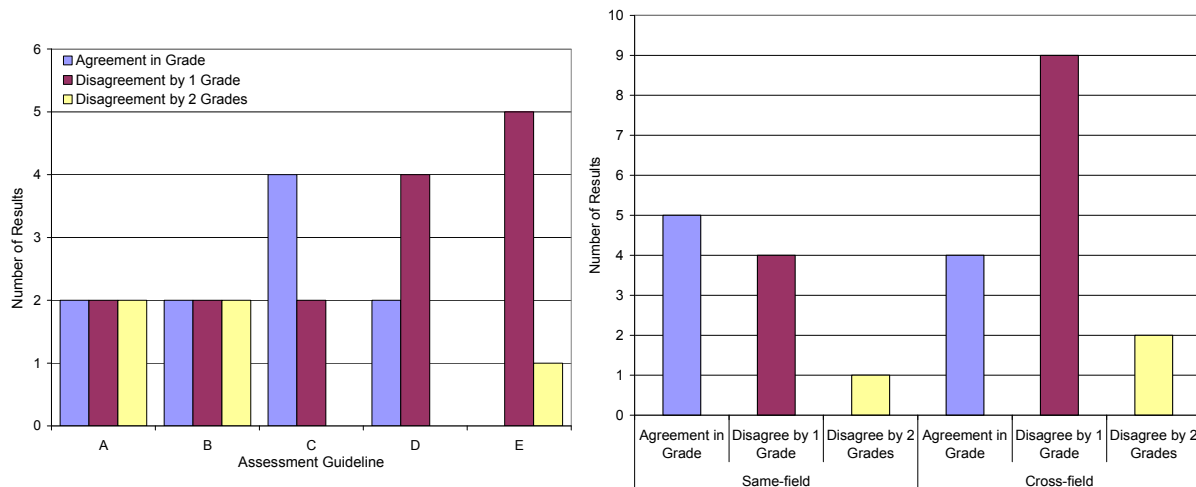


**Figure 1: Assessor grade agreement results for different guidelines (left) and assessor grade agreement results for same-field and cross-field assessor teams (right).**

The results were then examined to reveal whether the inter-rater agreement was improved if both assessors shared a common discipline background. The results shown in the right part of Fig. 1 indicate that there was still substantial disagreement when both assessors were within the same field, but there was clearly more disagreement when the assessor team was cross-field. The following section gives an account of the academics' experiences in using the five guidelines.

## Assessor's Comments on the Guidelines

Guideline A was quick and straight forward to use.  Most markers found that the prescribed weightings for each criterion gave guidance and one less item to consider. However, all markers disagreed with the weightings provided and felt that the descriptions given for levels of achievement were prescriptive and not sufficiently explained. The criteria were based on a standard report / thesis

structure. There was also general feeling that the marking scheme was too narrowly focused and would disadvantage theses not following the prescribed 'recipe'. The mark range assigned to each performance standard was a useful improvement over rubrics where only one standard can be chosen. It was suggested to increase this numeric range as having a range of only 2 was not sufficient. The mathematical recipe for combining the marks from different criteria was straightforward and easy to use.

Guideline B focused on criteria such as excitement, creativity and noting the context of the work. There was also an element of student growth from the experience that could only really be assessed by the supervisor, which proved very hard unless such a commentary was provided in the thesis. The rubric rewarded the student who had taken risks and learned from the experience. It was very interesting to apply this marking scheme to engineering honours theses and some markers saw high value in incorporating some of the criteria from this guideline into honours assessment guidelines.  It was generally concluded that this rubric was not very practical to the engineering situation and most markers found this guideline very subjective.

Most markers found it hard to distinguish between performance levels in Guideline C. The guideline used subjective descriptors such as significant and reasonable, that rely heavily on the marker's own opinion. The rubric had a separate criterion for the supervisor to complete which was acknowledged as a good idea but some markers still found it hard to judge the performance without knowing the original research questions. It was not possible to translate the letter grades into a final mark and it is unclear how grades from several assessors were to be combined into a final grade.

All markers noted that Guideline D suffered from having too many items to consider within each criterion. This made it quite subjective as so many items needed to be weighted up in assigning a mark for each criterion. No performance standard was indicated, assessors were only required to assign a mark to each criterion. Interestingly, some markers in this study developed their own marking scheme for determining a mark for each criterion. The supervisor section was judged a useful addition.

The comments on Guideline E were mixed. Some markers found the guideline easy to use while others found it difficult due to a lack of guidance. One problem was in determining a final mark when some elements of the thesis were determined to be at different grade levels. The guideline was judged as quite subjective, and did not separate out items that could only be known by the supervisor. A big jump was noted in performance standard between credit and distinction. This guideline relied on assessors views of each performance standard relative to potential to enrol in a higher degree. The guideline also made reference to norm-based statistics.

## Discussion

The results from this study are far from conclusive, but provide a basis for an interesting discussion about the relative merits of various approaches to assessing the penultimate undergraduate engineering activity – the honours thesis.

Guidelines that were perceived as easy to use and less subjective did not result in substantially improved inter-rater agreement.  For example, Guideline C which produced the best inter-rater agreement was viewed as subjective relying on hard to distinguish descriptors such as 'significant' and 'reasonable'. This suggests that amongst this group of raters there was some consistent understanding of the intended meaning of these words.  Nevertheless, the considerable variation in mark with this assessment guideline suggests that there will always be some level of disagreement and it may be unrealistic to expect a guideline to give perfect inter-rater agreement. This observation echoes the view of Tariq *et al.* (1998). While this presents some concern in itself, it should be noted that within the School in which the study took place, theses with a substantial inter-rater disagreement are subjected to a moderation procedure. For example, the review by Hannan *et al.* (2009)  showed that schools typically introduce a moderation procedure in cases where the assessors marks disagree by more than a set range, typically 5-15%.

There is some evidence that inter-rater agreement was improved when assessor teams are from the same discipline. Joyner (2003) notes that in assessing PhD theses, the examiner should "be sufficiently aware of the intellectual frontiers of their subject that they can judge whether the thesis makes a

contribution to knowledge or scholarship sufficient to justify the award". Honours theses are not of the same academic 'standard', but the findings of this study do raise some interesting questions about selecting the team for assessing student theses. This is particularly the case for transdisciplinary or non-traditional projects which will inevitably involve some out-of-field assessors. It is difficult for a out-of-field assessor to make an informed judgement of the intellectual frontiers of a different discipline. The exact definition of 'out-of-field' and 'in-field' warrants some consideration given the potential diversity of engineering research projects. Consideration is needed to ensure that students are not disadvantaged by having a different number of in- and out-of-field raters. This has much broader implications for research that moves beyond the traditional disciplines of engineering.

Assessment guidelines must be able to distinguish the role of the assessor in terms of in-field or out-of-field and supervisor or non-supervisor. For example guidelines that require non-supervisors to answer supervision related equations rely on the assessor soliciting and being able to access the supervisor's view of the student's performance, attitude or output relative to opportunity.

Students should be informed as to the general makeup of their assessment panel at the beginning of the research project so that they can present their thesis at an appropriate level for all of the assessors. Further, if assessors are to take account of supervision related factors such as attitude or output relative to opportunity, students need to be encouraged and guided in how to express their 'learning journey'. This could be done as part of the framework of the written thesis, or in a different form such as a reflective journal or e-portfolio.

The references to norm-based statistics made in one of the guidelines were viewed as very subject by the out-of-field assessors, as they found it very difficult to judge the quality of the work without knowledge of similar projects for comparison. This observations suggests that a CRA guidelines with well explained performance measures, might be better suited for assessing transdisciplinary research theses. Irregardless of which type, assessment guidelines must contain criteria that can be adequately considered by the rater. This requires each criterion and performance standard to be clearly explained to remove any ambiguity that may lead to a subjective interpretation.

Guideline B was not particularly well suited for assessment of engineering honours theses, but nonetheless provided an interesting experience for the assessors. In particular, it raised an important question of how well assessment guidelines reward qualities such as creativity and originality, which all assessors found to be lacking in the other guidelines.

As a final note, it is emphasised that the findings of this study should be interpreted with care considering the small sample size used in this study. The results and academics' views on using the guidelines have provided a basis for an interesting discussion and suggestions for improving assessment guidelines

## Conclusion

This paper has provided an interesting study of some of the challenges posed in developing reliable guidelines to assess a diverse range of engineering honours theses. The results suggest that some inter-rater disagreement is likely irregardless of the guidelines, and that inter-rater agreement is improved for assessment teams composed of the same discipline. With this in mind, assessment guidelines must clearly define all assessment criteria to avoid ambiguity and subjectivity. Guidelines must also make careful consideration of roles of each assessor with regard to their ability to accurately assess each criterion. This is particularly important when assessing transdisciplinary research projects that span across traditional discipline boundaries. While honours assessment will remain an inexact science, it is important that we strike a balance: on the one hand, affording the marker independence to apply their expertise and intellectual instinct to assessing student research, and on the other, providing robust, clear, flexible and fair assessment schemas capable of rewarding student efforts across the wide scope of disciplinary and transdisciplinary engineering research.

## References

Biggs, J. (2003). *Teaching for quality learning at University*, Second Ed., Berkshire, United Kingdom: Open University Press.

Boyer, E. L. (1997). *Scholarship Reconsidered: Priorities of the Professoriate*, San Francisco, Carnegie Foundation for the Advancement of Teaching.

Hannan, G., Burke, K., and Donovan, N. (2009). Assessment in Honours Programs in the Faculty of Science, Engineering and Technology at the University of Tasmania, unpublished report.

Hay, I. (1995). Communicating geographies: development and application of a communication instruction manual in the geography discipline of an Australian University, *Journal of Geography in Higher Education,* 19(2), 159-176

Joyner, R. W. (2003). The selection of external examiners for research degrees, *Quality Assurance in Education*. 11(2), 123-127

Littlefair, G. and Gossman, P. (2008). BE (Hons) final year project assessment – leaving out the subjectiveness. In L. Mann, A. Thompson, and P. Howard (Eds.), *Proceedings of the 2008 AaeE Conference, Yeppoon, Queensland, 6 pp.*

Stellmack, M. A, Konheim-Kalkstein, Y. L., Manor, J. E., Massey, A. R. and Schmitz, A. P. (2009), An Assessment of Reliability and Validity of a Rubric for Grading APA-Style Introductions. *Teaching of Psychology,* 36(2): 102–107

Tariq, V. N. and Stefani, L. A. J., Butcher, A.C. and Heylings, D.J. A. (1998), Developing a New Approach to the Assessment of Project Work, *Assessment & Evaluation in Higher Education,* 23(3), 221-240

Wall, S. and Shankar, I. (2008), Adventures in transdisciplinary learning, *Studies in Higher Education,* 33(5), 551–565

## Acknowledgements

## APPENDIX 1

### Guideline A

| Standards | Exemplary 5 – 4 | Satisfactory 3 – 2 | Unacceptable 1-0 | Score | Weight | Total |
|---|---|---|---|---|---|---|
| Abstract | Clearly states problem and question to be resolved; clearly summarizes method, results, and conclusions | Summarizes problem, method, results, and conclusions but lacks some details | Is vague about the problem; does not provide a summary of the whole project | | X 2 | |
| Introduction | Provides background research into the topic and summarizes important findings from the review of the literature; describes problem to be solved; justifies the study; explains the significance of the problem to an audience of non-specialists | Provides background research into the topic and describes the problem to be solved | Provides background research into the topic but does not describe the problem to be solved; fails to explain details to non-specialists | | X 1 | |
| Problem | Addresses the problem with clarity | Addresses the problem but may sometimes digress | Does not address the problem | | X1 | |
| Procedure | Presents easy-to-follow steps that are logical and adequately detailed; sampling is appropriate to the problem | Presents most of the steps but lacks some details | Has steps but some are missing or not sequential | | X1 | |
| Data + results | Explains data and results in a complete fashion | Explains data and results with some features lacking detail | Lacks description of data and results | | X2 | |
| Conclusion | Presents a logical explanation for findings; addresses recommendations for further research or use/application | Presents a logical explanation for findings | Does not adequately explain findings | | X2 | |
| Mechanics + documentation | Is free or almost free of errors of grammar, spelling, and writing mechanics; documents sources | Has errors but they don't represent a major distraction; documents sources | Has many errors that obscure meaning or add confusion; neglects important sources or uses too few | | X1 | |
| | | | | | Grand Total | |

## Guideline B

| Category | Description | Check & comment here | | |
|---|---|---|---|---|
| | | Good | Average | Needs work |
| Growth How does student now compare with earlier lessons? | In knowledge and vocabulary? | . | . | . |
| | In awareness and perception? Noticing? | . | . | . |
| | In interest, involvement, and attitude? | . | . | . |
| | In spontaneous comments about art topics? | . | . | . |
| Creativity | Speculates about meaning and feeling of work. Takes risks. entions alternatives. | . | . | . |
| | How original and innovative are comments, questions, and answers? | . | . | . |
| Fulfills Assignment | How well does written and spoken work solve the problems outlined in this assignment? | . | . | . |
| | Are variations from the assignments made for valid & creative reasons? | . | . | . |
| Knowledge | Gaining knowledge & awareness of art terminology & art purposes & concepts. | . | . | . |
| | Gaining knowledge and awareness of artists and their styles and work. | . | . | . |
| Helpful | Is the student positive, cooperative, tactful, & considerate in discussions? | . | . | . |
| | A thoughtful listener. Asked good questions? | . | . | . |
| Work Habits | Attentive and participatory? | . | . | . |
| | Do conversations with classmates stick to art topics and other appropriate related topics? | . | . | . |
| Composition And Design | Did the student see and mention the principles of design and composition and explain how things worked visually? | . | . | . |
| | Explains visual causes and effects in art? | . | . | . |

## Guideline C

### Honours research thesis assessment

| | HF | HU | HL | HT | HN |
|---|---|---|---|---|---|
| Criterion 1 Ability to plan and manage a scientific investigation, incorporating a substantial amount of original work, with clearly defined constraints of time, finance and technical resources | Played a major role in project development and demonstrated a sophisticated understanding of research methods, with evidence of careful attention to critical design issues in the execution of the project. Demonstrated a reasonable degree of autonomy, while still seeking advice when appropriate. | Good use of advice and resources from supervisory team, indicating a well-designed and competently conducted program and evidence of a solid understanding of research methods. Adequate design of the research project, although possibly containing minor but retrievable errors. | Relied on close control of the project by the supervisory team. Project well planned but limited evidence of creative input with basic but somewhat limited understanding of research methods. Generally adequate design of the research project but is marred by some errors and oversights | An unusual amount of help required by the supervisory team in planning and managing the project. Knowledge of research methods is deficient and serious flaws exist in the design of the research project making it difficult for the research to reach its aims | Evidence of poorly or unplanned research project. Knowledge of research methods is lacking and fatal flaws exist in the design of the research project making it impossible for the research to reach its aims. Limited or no evidence of seeking advice from the supervisory team. |
| Criterion 2 Ability to analyse and interpret scientific results | Evidence of significant insight and original thought in dealing with the critical issues. Clear and coherent interpretation of the thesis data and/or the results of other studies. | Evidence of reasonable insight and some evidence of original thought in dealing with the critical issues. Reasonable interpretation of the thesis data and/or results of other studies. | Occasional evidence of insight into the issues underlying the thesis, but little evidence of original thinking. Interpretation of results or other studies is adequate but limited. | Little evidence of insight and ideas tend to be highly derivative. Interpretations of results are superficial. | Serious misunderstanding of key concepts and issues and misinterpretation of results. |
| Criterion 3 Ability to report scientific data, using an appropriate range of techniques | Thoughtful and appropriate choice of data analysis (where appropriate) and outstanding presentation and reporting of results | Appropriate choice of data analysis for the design, although may not be well justified. Clear presentation of results. | Acceptable choice of data analysis, although other approaches may have been more appropriate. The presentation of results lacks clarity | Data analysis techniques are arbitrary in inappropriate. The results are poorly presented. | Data analysis techniques are inappropriate and the results are presented inadequately. |
| Criterion 4 Ability to discuss results and draw conclusions from a scientific investigation in relation to relevant chemical principles and published work | Superior evaluation and integration of existing literature and comprehensive understanding of the results in the context of the theoretical framework. | Good evaluation and integration of existing literature. Generally sound interpretation of results and their importance to the theoretical context. | Provides an adequate coverage of the literature, although it tends to be more descriptive than evaluative, and arguments are often disjointed. Limited interpretation of results in the context of the theoretical framework. | Coverage of the necessary literature is weak, with insufficient information provided to support the arguments made, or conclusions drawn within the thesis. Poor interpretations of results and their relevance to the theoretical framework. | Coverage of the necessary literature is inadequate, with little information provided relevant to the claims made, or conclusions drawn within the thesis. Inability to show how the results of the research project relate to the theoretical framework. |
| Criterion 5 Evidence of scientific communication skills, using effective written English | Outstanding command of expression and logical argument in a skillfully structured thesis. Correct use of relevant scientific terminology. Virtually free from typographic, grammatical and punctuation errors. Consistent referencing style used throughout. | The thesis is well written logically argued and generally well structured. Minor errors only in the correct use of relevant scientific terminology. Minor errors only in typographical, grammatical and punctuation formats. Consistent reference style used throughout with only minor errors. | The thesis is generally competently written, although some problems exist in the logical organization of the text and the way it is expressed. There are some errors in the correct use of scientific terminology. There are some typographical, grammatical and punctuation errors. Referencing style is consistent but contains errors. | The thesis is not well written and shows flaws in the structuring of logical arguments. Scientific terminology is used inconsistently and/or incorrectly. There are a number of typographical, grammatical and punctuation errors. Referencing style is inconsistent and/or contains many errors. | The thesis is very poorly written and shows a serious inability to structure and present a logical argument. Scientific terminology is used incorrectly or not used at all. Academic written presentation conventions are not adhered to. |

Legend: HF – honours first class (80-100%), HU – honours second class upper division (70-79%), HL – honours second class lower division (60-69%), HT – honours third class (50-59%), HN – honours – failure (0-49%). Criterion 1 is assessed by the supervisor alone, Criteria 2-5 are assessed by all examiners.

Grade allocation rules:

Supervisor

| | |
|---|---|
| HF | HF standard in 4 criteria with at least HU standard in the 5th. |
| HU | HU standard in 4 criteria, at least HL standard in the 5th |
| HL | HL standard in 4 criteria, at least HT standard in the 5th |
| HT | HT standard in 4 criteria, HN standard in no more than 1 of criteria 2-5 |
| HN | HN standard in criteria 1, or NH standard in any 2 other criteria |

Examiners (in-field and out-of-field)

| | |
|---|---|
| HF | HF standard in 3 criteria with at least HU standard in the 4th. |
| HU | HU standard in 3 criteria, at least HL standard in the 4th |
| HL | HL standard in 3 criteria, at least HT standard in the 4th |
| HT | HT standard in 3 criteria, HN standard in no more than 1 criteria |
| HN | HN standard in any 2 criteria |

## Guideline D

| Name of student: | Date submitted: | | Thesis title: |
|---|---|---|---|
| Name of marker: | Total no. of pages: | No. of pages of main text: | Title page OK? Abstract OK? Table of contents OK? |
| A) Problem definition<br>• Justification of research<br>• Clear statement of objectives<br>• Definition of research scope<br>• Literature review: relevance, diversity, depth<br>• Project plan (eg. experimental design) | Comments | | Mark<br><br><br><br>/20 |
| B) Technical content<br>• Quantity and quality of data collected<br>• Sophistication of data analysis<br>• Interpretation of results<br>• Logical argument<br>• Achievement of aims<br>• Conclusions supported by data and analysis | Comments | | Mark<br><br><br><br><br>/50 |
| C) Presentation<br>• Grammar, syntax and "visual appeal"<br>• Compliance with thesis guidelines (eg. word limit)<br>• Proper referencing<br>• Adequate use of appendices | Comments | | Mark<br><br><br>/30 |
| D) Sub-total (all markers) | A) +B)+C) | | /100 |
| E) Student Effort (first supervisor only)<br>• Level of understanding<br>• Appreciation of engineering context<br>• Leadership in project management | Comments (first supervisor only) | | Mark<br><br><br>/20 |
| | Total Mark* | D)+E) | /120 |
| | Final Mark | | /70<br>/65 |

\* The sub-total (D) mark from each marker is averaged (a third marker will be used if the marks differ by more than 10%). This sub-total component is marked out of 100. The Student Effort (E) is markes out of 20 by the first supervisor only. The Total Mark is the sum of the average (D) mark and the (E) mark and is out of 120. The final thesis mark is adjusted to a mark out of 70 or to a mark out of 65.

## Guideline E

Guidelines for marking honours theses

<50 Fail. The project either has significant flaws, or the quantity or degree of challenge of work (quality) undertaken falls well short of the supervisor's expectations. For example the thesis may fail to adequately cover the following major components: 1) identification of the rationale for selection of the topic and the basic objectives to be achieved in the project; 2) identification of the basic approaches to be undertaken and have completed the major steps toward achieving the outcomes of the work; 3) demonstrate adequate technical competence and 4) indentify the main applications of the outcomes of the work.

>50 to 59 Pass. There is basic but minimal level of achievement and some minor flaws can be identified by the marker(s). The thesis would include the four basic components mentioned above.

>60 to 69 Credit. The thesis must show clearly demonstrated technical and planning competence, analysis and outcomes, although the work may be routine or not requiring particular extension, independence or initiative. Both technical and research skills and writing skills must be above average. Only very few minor flaws can be identified. It also demonstrates the logical and analytic thinking in problem solving process. It has demonstrated some comparison study to other approaches in the literature. It elaborates the advantage of the method used and results obtained in the thesis.

>70 to 79 Distinction. This level implies a significant degree of initiative, independence, or originality, for example in extending the scope of the project beyond that initially defined by the supervisor. Both technical and research skills and writing skills must be well above average. The thesis has demonstrated thorough understanding of the relevant knowledge to the project, clearly identifies the scope, logically and completely describes the approaches and works, includes the necessary technical details, analysis, clearly summarises the outcomes and theoretical results, and illustrates the advantages of the approach employed.

>80 to 89 High Distinction. Such theses are at the first class honours standard, i.e. the student(s) may possess the capability to enrol in and complete a PhD. Both technical and research skills and writing skills must be excellent. The thesis must stand out when compared with other honours theses and demonstrate clear evidence of at least two of the following beyond that normally expected of a very competent student.
• Originality, innovation
• Initiative, independence
• Extension of the scope of the project beyond that defined by the supervisor

The quality of work would be deserving of a postgraduate scholarship and should be able to be extended to a refereed publication. The level is normally attained by less than around 10% of projects/theses.

>90 High Distinction. This level is for exceptional projects and is normally attained by less than around 5% of projects/theses. The quality of work would be deserving of a postgraduate scholarship and start to approach the level of a MEngSci thesis and should be suitable for a refereed publication.