

A Review on IELTS Writing Test, Its Test Results and Inter Rater Reliability

Veeramuthu a/l Veerappan

Centre For Foundation Studies And Extension Education, Jalan Persiaran Multimedia, 63100, Multimedia University,
Cyberjaya, Malaysia

Email: veeramuthu@mmu.edu.my

Tajularipin Sulaiman

Faculty of Educational Studies, University Putra Malaysia

Abstract—This research paper is based on the Academic Modules of IELTS test. A study is conducted to evaluate the performance of IELTS test takers and to analyze inter rater reliability in evaluating the test. The test is held at a private university and there are 9 participants who are the students of Foundation, Diploma and Degree courses were assessed only on writing component. This test is rated by three different raters to observe inter rater reliability in evaluating the test. The findings show that rating is not reliable and there is less agreement between the raters as different raters awarded different marks for the students. The total accumulated marks are 588 out of maximum 900 marks and the average score of these 9 students are 65.3% and falls under band 6. Since the IELTS test uses holistic marking, the raters find that the scoring rubrics are too vague in defining the qualities of writing. This study proposes analytic scoring system as a better procedure in assessing IELTS test, whereby it uses separate scales to assess different aspects of writing rather than a single score. The inter rater reliability issue can be overcome if test providers are able to furnish grading scales that could clearly equate the total marks achieved by a candidate to an appropriate band. The result and discussion from this study cannot be generalized to all IELTS candidates as there are many limitations in this study such as student preparedness, tester marking experience, time constraint, and reliability of the test and the validity of the marks tabulation method.

Index Terms—IELTS, writing assessment, holistic marking, scoring rubrics, inter-rater reliability

I. INTRODUCTION

The IELTS writing test is used as an example of a large scale high stake test in evaluating the test performance of a group of test takers. The IELTS test is jointly provided by three organizations: the British Council, IDP: IELTS Australia and Cambridge ESOL, with its development and validation unit based in Cambridge, UK. There are two versions of the IELTS test: the Academic Modules and the General Modules. Both versions contain four components: listening, speaking, reading and writing. However, exam providers of any large scale high stake tests need to clearly specify to the purpose of the test. As such according to the IELTS handbook, the Academic Modules of the IELTS are designed to assess a candidate's English language proficiency for academic studies at the undergraduate or post graduate level, whereas the General Modules are developed to assess test candidates who intend to go to English-speaking countries to complete their secondary education or undertake work experience or training programmes at below degree level. People who need to demonstrate their English proficiency in order to immigrate to Australia or New Zealand are also required to sit the General Modules (British Council, IDP: IELTS Australia and University of Cambridge ESOL Examinations, 2005).

However in this study, the evaluation of this test is based on the Academic Modules of IELTS test. The objectives of this action research are to evaluate the performance of its test takers and to analyze the inter rater reliability of different raters in evaluating the test. The test was held at a private university and the participants who are the students of Foundation, Diploma and Degree courses were assessed only on the writing component. The writing component, lasting for 60 minutes, is made up of two academic-oriented tasks; Writing Task 1 and Writing Task 2. It is suggested that about 20 minutes is spent on Task 1. The first task asks the candidate to describe in about 150 words a chart, a diagram, a graphic or a table which they might encounter during their study at university. As such, in this section candidates are assessed on their ability to:

- i. Organize, present and possibly compare the given data
- ii. Describe the stages of a process or procedure
- iii. Describe an object or event or sequence of events
- iv. Explain how something works

On the other hand, the second task usually requires the candidate to write an argumentative essay of about 250 words based on a controversial topic supplied in the question paper. Task 2 should take about 40 minutes. The candidates are assessed on their ability to:

- i. Present the solution to a problem
- ii. Present and justify an opinion
- iii. Compare and contrast evidence
- iv. State opinions and implications of a given issue
- v. Evaluate and challenge ideas by giving a view point
- vi. Argue in support of or against a given statement

However, the issues raised in both the tasks are of general interest which are suitable for and easily understood by candidates entering undergraduate studies or seeking professional registration. The IELTS candidates' writing performance is rated by a single certified rater at local test centers for the purpose of score reporting. Reliability of writing assessment is ensured through a sample monitoring process, where a sub-sample of the candidates' performances are collected and later re-rated by senior examiners for quality check. In 2003, the overall correlation agreement between the local raters and senior examiners were 0.91 for writing scores. Reliabilities based on these correlations were therefore 0.84, using the Spearman-Brown Formula (British Council, IDP: IELTS Australia and University of Cambridge ESOL Examinations, 2004). Finally, the results from the IELTS are reported on a scale from 1 (non-user) to 9 (expert user). However, the writing scores are reported in whole bands only. Table 1 below lists the descriptions of the nine overall bands.

TABLE 1
SHOWS THE BAND, SCALE AND CRITERIA FOR IELTS TEST

Criteria	Band
Has fully operational command of the language: appropriate, accurate and fluent with complete understanding.	Expert user Band 9
Has fully operational command of the language with only occasional unsystematic inaccuracies and inappropriacies. Misunderstandings may occur in unfamiliar situations. Handles complex detailed argumentation well.	Very good user Band 8
Has operational command of the language, though with occasional inaccuracies, inappropriacies and misunderstandings in some situations. Generally handles complex language well and understands detailed reasoning.	Good user Band 7
Has generally effective command of the language despite some inaccuracies, inappropriacies and misunderstandings. Can use and understand fairly complex language, particularly in familiar situations.	Competent user Band 6
Has partial command of the language, coping with overall meaning in most situations, though is likely to make many mistakes. Should be able to handle basic communication in own field.	Modest user Band 5
Basic competence is limited to familiar situations. Has frequent problems in understanding and expression. Is not able to use complex language.	Limited user Band 4
Conveys and understands only general meaning in very familiar situations. Frequent breakdowns in communication occur.	Extremely limited user Band 3
No real communication is possible except for the most basic information using isolated words or short formulae in familiar situations and to meet immediate needs. Has great difficulty understanding spoken and written English.	Intermittent user Band 2
Essentially has no ability to use the language beyond possibly a few isolated words.	Non user Band 1
No assessable information provided.	Did not attempt Band 0

The IELTS test development processes ensure that the test is of a comparable level of difficulty, so that it provides valid and consistent results. The testing system is underpinned by test materials developed using the following stages: commissioning, editing, pre-testing, analysis, banking of material, standards fixing and question paper construction. In addition to this test development, test writers from different English-speaking countries also develop IELTS content, so that it reflects real-life situations around the world. This would be unbiased and fair to all IELTS candidates who come from different cultures and countries.

II. THE IMAGED WRITING TEST SPECIFICATIONS

Test specification is an official statement about WHAT a test tests and HOW it tests it and a test specification is also known as a blueprint of a test. It is used internally among test constructors and test administrators as well as test raters and usually kept confidential. A test specification can also be described as a simplified document which indicates what the test contains. The development of a test specification is crucial part of test construction and evaluation process which serves as a blue print and it provides answer to what is tested, who are the targeted audience or test takers and how does the test look like. In this project, we have designed a test specification based on our imagination for the writing skills. Below is the test specification that is used in conducting this research. First of all the Writing test that we have conducted takes 60 minutes. There are two tasks to complete and in task 1, students are asked to transfer the information from non-linear graphs and diagrams to a linear form of paragraph by using not less than 150 words. We

suggested the students to spend 20 minutes to complete task 1 and in task 2 the students are asked to write a continuous essay that requires them to write at least 250 words and to use not more than 40 minutes to attempt this task. Task 1 is an academic writing and the test takers are asked to describe some information based on graph, table, chart and diagram. The test takers are asked to present the description in their own words. The assessment of the task is based on the candidate's ability to organize, present and possibly make comparison based on a non-linear text provided. The candidates are also expected to describe the stages of a process and procedure. To add on, the candidates must be able to describe an object or event or sequence of events and finally they should be able to explain on how something works. In task 2, the test takers are presented with a point of view or argument or problem. Test takers are assessed on their ability to present the solution to a problem, present and justify an opinion, compare and contrast evidence, opinions and complications. The candidates are also required to evaluate and challenge ideas, evidence and argument. The questions are catered to raise general interest among candidates' general interest and it can be easily understood by candidates entering undergraduate or postgraduate studies and also useful for people who are seeking for professional registration. In task 1, test takers are assessed based on the following criteria: 1. Task fulfillment: Have you followed the instructions clearly? Have you given a clear, accurate and relevant description of the information? 2. Coherence and cohesion: Is your writing well organized? Are sentences logically linked? 3. Vocabulary and sentence structure: Have you used a variety of appropriate vocabulary? Are the sentences well constructed? Meanwhile, in task 2, test takers are assessed based on their ability to present a point of view or argument or problem. They are required to: 1. Present the solution to a problem. 2. Present and justify an opinion. 3. Compare and contrast evidence, opinions and implications. 4. Evaluate and challenge ideas, evidence or an argument. In task 2, test takers are assessed based on the following criteria: Good organization-paragraphing, thesis statement, topic sentences and main ideas, clear ideas that directly address topic, coherent argument or point of view and discussions with reasons provided to support, evaluate and challenge ideas, cohesion by the use of transitional phrases, linkers, sentence connectors, synonyms, pronouns and references, accurate and appropriate structures and vocabulary and good punctuation and spelling.

III. DESCRIPTION OF TEST ADMINISTRATION AND DATA COLLECTION

Three groups of students were chosen from a range of available students at tertiary level in a private university. All these students are second language learners. The process of identifying students for the administration of the test was carried out after carefully analyzing our goals and the capability of these learners. The first thought that crossed our mind was that, will these individuals be competent enough to do the writing task. Since IELTS is a comprehensive language testing system, we were quite worried if the test results would give us a clear picture and thereafter proceed to analyze the performance of the test takers. The first group of students was from the foundation program and they are in their final semester. These learners were picked at random from a class of seventeen students. These international students do not have English certificates from neither TOEFL nor IELTS and in order to proceed with the university programme, they have gone through and successfully completed an Intensive English language programme conducted by the university as to meet this university's language requirement. Before being selected, a brief introduction of the study was done and immediately six students volunteered for the exercise. Since all of them are not familiar with the writing task, we decided to hold another briefing session to explain the requirements of the task. Thereafter, the time and venue was decided after consulting the students. The time factor was crucial because we wanted the students to do well, and if they chose the most appropriate time for them to sit for the task, it would be motivating. After this, all six of them sat for the paper and were given three minutes of reading time before starting on the task whereby only twenty minutes was allocated as per the IELTS requirement. After the reading time, one student gave up and mentioned that he could not understand the question and did not want to continue. As for task two, which requires 40 minutes, we decided to give the test takers five minutes of reading time. Since this task is more stringent and requires the student to write 250 words we thought that the reading time should be longer. The next group was from the diploma course whereby three of them volunteered to try the task. After a quick briefing session on the requirements, they started without any reading time and managed to complete both the papers within the total time of an hour. These students have completed two levels of English and are more competent. Their level of confidence was also high as compared to the foundation students because they managed to complete the paper without much difficulty. This was evident while administering the test. The next group was the degree students who had also volunteered but were keen to know if they would be allowed to have more time to complete the task. The only possible answer to this was to explain to the students that this is a study on language assessment and evaluation and no other bearing on what they are required to do. These degree students were asked to sit for this task at 5pm after classes and they took an hour to complete both the papers. This may not be the most suitable time for them to take the test, but due to time constraints the students managed to complete and hand in their tasks. The scores of all the tests taken are shown in table 2. It is the usual practice of the said test for two markers to evaluate each task.

TABLE 2
SHOWS INTER RATER SCORE, AVERAGE MARKS AND BAND

No	Name	Level	Rater 1		Rater 2		Rater 3		Average Marks/Band
			Task 1 30%	Task2 70%	Task 1 30%	Task2 70%	Task 1 30%	Task2 70%	
1	Bharati	Degree	23	42			20	39	124/2= 62 Band 6
2	Reza	Diploma	20	63			18	48	149/2= 74.5 Band 7
3	Ahmed	Foundation	22	56			20	55	153/2= 76.5 Band 7
4	Syaley	Foundation	26	64			21	58	169/2= 84.5 Band 8
5	Yuliana	Degree	27	59	13	25			124/2=62 Band 6
6	Fayzali	Degree	24	62	15	35			136/2= 67.5 Band 6
7	Meng Xiang	Diploma	22	43	11	20			96/2=47.5 Band 4
8	Brahmuda	Diploma			13	28	21	51	113/2= 56.5 Band 5
9	Nikita	Foundation			22	30	19	45	116/2=57 Band 5
	Average Score/ task/rater		23.4	55.5	14.8	27.6	19.8	49.3	Total Average 588/9 = 65.3

IV. PRESENTATION AND DISCUSSION OF THE TEST RESULT

Table 2 above shows inter rater score, average marks and band obtained by nine candidates who have sat for IELTS writing test. This test is conducted for the purpose of evaluating the performance of its test takers. This test is rated by three different raters to observe inter rater reliability in evaluating the test. Based on the table 2 above, the findings show that the rating is not reliable and there is no agreement between the raters as different raters awarded different marks for the students. There is not much agreement between the raters on the official rating as the statistic shows that the average score awarded by rater 1 to student 5 is 86 marks for both task 1 and task 2 meanwhile rater 2 awarded 38 marks for the same student. Huge differences can be seen in the distribution of marks and the inference that could be made is that there is no clear rating scale which explicitly states the scoring criterion such as the multiple trait scoring to guide the raters. Since the scores on this test are independent estimates of these judges or raters, the average score from each rater is taken as the final score to be awarded for the students. These would not be fair to the test takers as the standard measurement error of the test is not conducted. The disagreement among raters could also be due to inaccuracy and inconsistencies in marking as all raters are experiencing marking the IELTS writing section for the first time. The error within the raters could also be due to lack of experience in marking IELTS writing examination. The statistic also shows that there is a rater reliability and agreement between rater 1 and rater 3 in awarding marks to student 3 where rater 1 awarded a total of 78 marks for this student while rater 3 awarded 75 marks to the same student. There is only a slight disagreement between these raters as both demonstrate their independency in marking. This shows that there is an inter rater reliability and mutual agreement among these two raters. On the other hand, the statistics show the achievement of all nine test takers and out of these nine students, who have sat for this test, three of the test takers scored band 6, 2 students obtained band 5, the other two students managed to get band 7 and one student achieved band 4. The total accumulated marks are 588 out of maximum 900 marks and the average score of these 9 students is 65.3% and falls under band 6. All the degree students fall under band 6. These are the raw scores of the IELTS test as there are other skills and components that will determine the overall achievement and band of IELTS candidates.

TABLE 3
SHOWS PERFORMANCE ANALYSIS ON EACH COMPONENTS/ELEMENTS

No	Student's Name	Task 1 30%				Task 2 70%					Score
		TF	C&C	V	SS	O	AIE	CQ	VSS	M	
1	Bharati	3.5	3.5	7.5	7	13	10	5.5	5	7	62
2	Reza	4	2.5	7	5.5	16	17	7.5	8	7	74.5
3	Ahmed	3.5	3	7.5	7	17	16	7	7.5	8	76.5
4	Syaley	4.5	4.5	8	7	18	18.5	8.5	8	8	84.5
5	Yuliana	3.5	3.5	6	6.5	11	9.5	6.5	5	5	62
6	Fayzali	3	3.5	7	6	13	15.5	6.5	6	7	67.5
7	Meng Xiang	3	3.5	5	5	8	8.5	5	4.5	5	47.5
8	Brahmuda	3.5	3	6	4.5	12.5	11.5	5.5	5	5	56.5
9	Nikita	3.5	2.5	7.5	7	9.5	10	4.5	6	6.5	57
Achieved Percentage for each element		71%	65.5%	69%	55.5%	65.5%	64.7%	62.7%	61.1%	65%	Total Score 63.5%
											Total Average 64.7%

Task 1 Scoring Criteria

TF= Task Fulfillment	5 marks
C&C= Coherence and cohesion	5 marks
V= Vocabulary	10 marks
SS= Sentence Structure	10 marks
TOTAL	30 marks

Task 2 Scoring Criteria

O= Organization	20 marks
AIE= Arguments, ideas, evidence	20 marks
CQ= Communicative Quality	10 marks
VSS= Vocabulary and sentence structure	10 marks
M= Mechanics	<u>10 marks</u>
TOTAL:	70 marks
GRAND TOTAL:	100 marks

Table 3 above shows the statistics of candidates in their IELTS writing examination. Student 4 is the highest achiever among all nine candidates where the score obtained is 84.5 marks for both tasks and fall under band 8. This candidate has good command of the language with occasional inaccuracies and inappropriacies. This candidate is also able to provide valid arguments and ideas supported with evidence. Meanwhile the lowest achievement is by student 7 where this student could only manage to obtain a score of 47.5 for both the tasks. This student is unable to handle different situations clearly and has poor communicating quality in conveying his thoughts in writing and there is no cohesion and cohesiveness in his writing. The statistic shows that most of the students are able to fulfil the task and scored an average of 71 % in task fulfilment. The statistic shows that the candidates are able to obtain a score of 65% in their marks given to mechanics as they are able to use correct punctuation, capitalization and spelling with some minor mistakes and grey errors in this element. The lowest average score is given to candidates' sentence structure as most of the students made many mistakes in sentence structure and unable to write grammatically. Most students are unable to make clear inferences from a non-linear text in task 1 to a linear text as they have average vocabulary knowledge to paraphrase the diagrams and graphs. The percentage obtained by all students in organization of paragraphs are 65.5% as many of the candidates are able to write the essay in paragraphs with interesting introduction, correct thesis statement, supporting details in every paragraphs and supported by details and examples. The conclusions are also found to be summarizing the whole essay with some innovative suggestions and opinions. The total average score obtained by all nine candidates is 63.5% and the total percentage achieved by these candidates for all elements is 64.7%. An average of 22 marks are obtained by candidates in accomplishing task 1 and an average of 59 marks are obtained by all candidates in task 2. This infers that candidates are able to organise, present, make comparison and explain a process or procedure between two sets of data given as the stimulus in task 1 better than to provide a solution, justify an opinion, provide evidence, evaluate and challenge ideas required in task 2. Based on the result obtained from the test, it is clear that students are able to perform on a guided writing task better than a continuous writing task. This writing test is conducted as a small scale action research among international students and the raters are language lecturers from 3 different higher learning institutions in Malaysia. The test takers are given an average period of time ranging from week 1 to week 5 to prepare for the test with intensive classroom instructions and resources as the testers are facing time constraint to complete the study. Limited face to face instructions and writing techniques skills are being taught to these candidates to enhance and prepare them to excel in their test. Hence, the result shows that the average band obtained by these students is only band 6. The candidates who are from different levels such as foundation studies, diploma students and undergraduates are anticipated to perform better than the actual achievement in the writing test in this small scale study if enough classroom input, guidance and facilitation be provided before the actual test. The result and discussion from this study cannot be generalized to all IELTS candidates as there are many limitations in this study such as student preparedness, tester marking experience, time constraint, and reliability of the test and the validity of the marks tabulation method.

V. SUGGESTIONS FOR TEST CONSTRUCTORS

The IELTS writing test is assessed based on a holistic marking. In holistic scoring, raters judge texts as a whole and they are not able to separate parts of the essay and identify them. Despite its wide use in writing assessment, holistic marking has recently been criticized as a procedure that fails to provide sufficient information on writing performance (Elbow 1996). Hamp-Lyons (1995) argued that because of the complex and multi-faceted nature of writing, the writing of second language (L2) students may show varied performance on different traits; subsequently, a great deal of information may be lost when assigning a single score to a piece of writing. Since the IELTS test uses the holistic marking, the raters find that the scoring rubrics are too vague in defining the qualities of writing. Hence, analytic scoring system would be a better procedure in assessing IELTS test, whereby it uses separate scales to assess a different aspect of writing rather than a single score. The scripts can be rated on features such as content, organization, cohesion, vocabulary, grammar and mechanics. These aspects of writing are differently weighted to provide more detailed information about a test taker's performance. The second issue with IELTS writing test is that there seems to be inconsistency in assessing the candidates' essay scores. Although measures have been taken to make sure writing is

being assessed 'correctly', raters cannot consistently agree with each other when assessing the same writing samples or even with judgments about the same samples made on different occasions (Hamp-Lyons, 1992:80). This is so true in the case of IELTS writing test where despite the inter-rater method, the test scores given by the two different raters show a large variance. In order to achieve consistency in assessing writing test, raters must ensure that they apply the standardize features of writing in the same way over time, despite their different background knowledge and experience that might affect their perceptions of the written product. In addition, if test constructors could equip the raters with sufficient information on how much weight they have to give to various writing components, the issue of inconsistency in assessing could further be reduced. Lastly, the IELTS examination board needs to demonstrate and share how to operationalize criteria distinctions between levels in the writing tests. This is in evident as the *IELTS Writing Assessment Guidelines* (WAG) are not always made available to practicing examiners, thus making them not fully conversant with the WAG. As such, focus should be given on providing useful information to help raters interpret test scores. However, the issue can be overcome if test providers could furnish grading scales that can clearly equate the total marks achieved by a candidate to an appropriate band. As such, a range of test scores in term of percentage should be assigned to a particular band to make marking more comprehensive. It could also serve as a reliable benchmark to raters in awarding a fair and accurate grade to the candidates.

VI. CONCLUSION

In conclusion, evaluating writing is a complex process that requires accounting for multiple factors to ensure a fair and accurate judgment of the writer's abilities. Testers must be sure to set purposes for assessment, choose appropriate scoring criteria, and be aware of how multiple variables in the testing situation can influence outcomes. By being aware of these issues, test constructors involved in the testing process have a greater chance of providing a more meaningful and justifiable assessments to the candidates.

REFERENCES

- [1] British Council. (2005). IDP: IELTS Australia and University of Cambridge ESOL Examinations 2004 Test Performance 2003. Retrieved February 26, 2011 from (<http://www.ielts.org/teachersandresearchers/analysisoftestdata/article173.aspx>)
- [2] British Council. (2006). IDP: IELTS Australia and University of Cambridge ESOL Examinations 2005 IELTS Handbook.
- [3] Elbow, P. (1996). Writing Assessment in the 21st century: A Utopian View. In: L. Z. Bloom, D. A. Daiker & E. M. White (Eds.), *Composition in the 21st Century: Crisis and Change*, (p. 83-100). Carbondale: Southern Illinois University Press.
- [4] Hamp-Lyons, L. (1989). *Raters Respond to Rhetoric in Writing*. In H. Dechert & G. Raupach (Eds.), *Interlingual Processes*, (p. 229-244). Tübingen: Gunter Narr Verlag
- [5] Gunter Narr Verlag. Hamp-Lyons, L. (1992). Holistic Writing Assessment for LEP Students. C Proceedings of the Second National Research Symposium on Limited English Proficient Student Issues: Focus on Evaluation and Measurement, OBEMLA.
- [6] Hamp-Lyons, L. (1995). Rating Non-native Writing: The Trouble with Holistic Scoring. *TESOL Quarterly*, 29(4): 758-762.



Veeramuthu a/l Veerappan was born in Johor Bahru, Malaysia on 17th February 1975. He completed his Diploma in Management from Institute of Tun Abdul Razak in 1998, obtained a Bachelor of Education (TESL) 2006 and Master in Applied Linguistics (English) 2010 from University Putra Malaysia. He is currently pursuing PhD in English studies. His major field of study is teaching English as a second language, applied linguistics and ESP. He started teaching career in SMK (LKTP) Tengaroh A in the year 2000. He is currently working as lecturer at FOSEE, Malaysian Multimedia University. He has experience teaching English as Second Language for 12 years for both local and International students. His research area of interest is in ESL writing, sociolinguistics and ESP.

Tajularipin Sulaiman Ph. D is currently a lecturer in Faculty of Educational Studies, University Putra Malaysia. He received his early education in Muar, Johor and continued his secondary education at the Sekolah Teknik Malacca. He continued his studies at the Centre for Foundation Studies in Science, Universiti Malaya in 1998. He obtained a degree in Bachelor of Science with Education in 1994 and Master of Education in 1998 from University Malaya. He holds a doctoral degree in education from University Putra Malaysia. His area of specialization is in pedagogy, science education, and cognitive development. He has involved teaching and researching in science education, thinking skills and primary education. His research interest is in the field of pedagogy and primary education especially in primary science. He has also presented papers in national and international conferences.