

CHAPTER IV

FINDINGS AND DISCUSSIONS

This chapter discusses about the research findings and discussions. It provides the analysis and interpretation of data that had been collected to answer the research question about the consistency analysis of English teachers of Al Amin Islamic Boarding Senior High School Mojokerto in scoring essay test.

As explained in chapter III, grade eleven teachers held a weekly writing examination; an exposition essay test to know students' achievement in writing that kind of text. The students wrote the essay in handwriting to keep the originality of their writing. Even it might increase the subjectivity but the researcher had kept it by omitting the students' name before copying and giving to the examiners or raters; mean the teachers. Each six teachers or raters were asked to grade 10 papers of 40 essays which were chosen randomly by the researcher in pre- and post-scoring. 5 of 10 papers in post-scoring were the same essay that they had actually rated in pre-scoring. After two months interval, some teachers admitted that they did not remember about ever seen those papers before. In addition, the others said that although they remembered having ever seen the papers but they could not remember the grades that they gave.

After getting all teachers' pre-scoring in 22nd of February 2016 and post-scoring in 18th of April 2016, the researcher analyzed the data by using Cronbach

RATER 2	2nd	10.5	10	9.5	10	35	30	11.5	10	4.3	4	70.8	64
	3rd	13	13	14	13	46	41	14	14	4.3	4	91.3	85
	4th	14	13	14	13	43	35	13	10	4.6	3.5	88.6	74.5
	5th	12	13	11	13	38	41	12	12	4.6	4	77.6	83
RATER 3	1st	15	14	15	13	41	46	15	14	4.6	4.6	90.6	91.6
	2nd	14	12.5	14	12.5	38	40	14	14.5	3.8	4	83.8	83.5
	3rd	14	14	12.5	13	40	43	14	13	4.2	4.2	84.7	87.2
	4th	15	14	15	13	40	43	12.5	13	4.2	4.3	86.7	87.3
	5th	14	14	12	14	40	40	15	13	5	5	86	86
RATER 4	1st	12	11.5	12	11	38	35	11	11	4	3.8	77	72.3
	2nd	11	11	11	11	35	35	10	11	4	3.7	71	71.7
	3rd	11	11	11	11	36	35	11	11	3.7	3.7	72.7	71.7
	4th	12	12	12	12	38	40	11.5	12.5	3.8	3.8	77.3	80.3
	5th	11	11	11	11	35	35	10	11	3.7	3.7	70.7	71.7
RATER 5	1st	13	12	13	12	40	40	13	12	4	4	83	80
	2nd	10	14	13	12.5	40	36	13	12	4	4	80	78.5
	3rd	13	12.5	13	12.5	35	41	12	12.5	4	4.2	77	82.7
	4th	13	12.5	12.5	12.5	35	41	13	12.5	4.2	4	77.7	82.5
	5th	10	14	13	12.5	46	36	13	12	4	4	86	78.5
	1st	11	12.5	10.5	12	36	43	11.5	13	3.8	4.2	72.8	84.7

the score was included in the range of $0.800 < \alpha \leq 1.000$. Therefore, the table above presented each result in the word.

As Table above shows, 1st rater got Very High reliability level in Grammar whereas High level in Organization. Besides, it indicated reliable Enough for Content but Very Low in Vocabulary. Mechanic was the worst as it got negative score which means that it was very unreliable. Luckily, his total score was very reliable as it got Very High level.

2nd rater was different from the 1st. Most of categories got High level, such as Organization, Grammar, Vocabulary and it might influence the Total score. The most reliable was in Content because 0.831 means it existed in “Very High” level. In contrary, Mechanic presented unreliable as it got negative score like the 1st rater.

Mechanic and the Total score of 3rd rater were almost perfect as it presented Very High reliability. He got High level in Grammar, reliable Enough in Content and Low level in Vocabulary. Unfortunately, it was the same as two raters before that they have minus value in one of their categories which means unreliable, this rater was in Organization.

This rater was the most stable than the others, 4th rater. It did not have negative value for any categories. Four categories had achieved High level, like Organization, Grammar, Vocabulary and absolutely it were impacted to the Total. The highest level, almost perfect, was Content. Besides, Mechanic

was the lowest reliability level as it got .375 score and it was still in positive value.

The worst unreliable raters are 5th and 6th raters. All of their categories presented negative value. It can be said that they got 0.000 score or were admitted as zero reliability. Simply, it was regarded that they were included in inconsistent or unreliable level.

As the various marks gotten, it was needed to make the average of all grades so that it could conclude the result which represented and covered all raters in all categories. Here is the table of average result. The table shows the average result of pre- and post-scoring of five same essays. For example, the pre-scoring of 1st rater was the average result from all pre-scoring in all categories and so was the post-scoring.

4.4 Table of All Essays' Average Results

No. Essay	PRE	POST
1 st	81.8	77.4
2 nd	75.8	72.9
3 rd	83	80.7
4 th	81.9	81.5
5 th	79.3	77.6

The first table is Paired Sample Statistics that showed the statistic summary of pre- and post-scoring. The table provides that the average score in pre-scoring was 80.360 and in post-scoring was 78.020. It indicated reduction for about 2.340. The standard deviation presented the data variation in each variable, that in pre-scoring was 2.887 and in post-scoring was 3.394. Also N was the number of data which there were five essays graded twice by raters in two-week interval.

Paired Sample Correlation showed the correlation between two variables that produce 0.902 with 0.036 for the significance. It means that the correlation between pre- and post-scoring was so related.

The last is Paired Sample Test. It can be interpreted as:

- Hypothesis

H_0 = the intra-rater reliability of English teachers at Al-Amin Islamic Boarding School Mojokerto in scoring essay test is not consistent.

H_1 = the intra-rater reliability of English teachers at Al-Amin Islamic Boarding School Mojokerto in scoring essay test is consistent.

- Significance level

Sig = 0.05

- Critical area

Based on t -test:

Reject $H_0 = t\text{-test} > t\text{-table} (5\%, N-1)$

Accept $H_0 = t\text{-test} < t\text{-table (5\%, N-1)}$

Based on $p\text{-value (Sig.)}$:

Reject $H_0 = p\text{-value} < 0.05$

Accept $H_0 = p\text{-value} > 0.05$

- Decision

$t\text{-test} = 3.541 > t\text{-table (5\%, N-1)} = 2.776$;

$\text{Sig.} = 0.02 < 0.05$;

means H_0 is rejected.

The intra-rater reliability of English teacher at Al-Amin Islamic Boarding School Mojokerto in scoring essay test was consistent.

B. Discussion

According to the finding, teachers were mostly very consistent in scoring the same paper. Table 4.3 showed that most of raters achieved upper consistency level beside only some raters got enough and lower level. The most consistent rater in all categories was the fourth. Even the fourth rater achieved Low level in Mechanic but no one categories got negative value or zero reliability that means inconsistent. The first rater was very consistent in Grammar and Total, but did not do well in his ratings of Mechanic. The second was almost the same as the fourth yet he was very inconsistent in Mechanic too, like the first. The third was very consistent in Mechanic and the Total score as both got Very High level but very

inconsistent in Organization. Unfortunately, the fifth and sixth teacher or rater seemed to be the least consistent. Their ratings were abysmal in all categories. In fact, they even contradicted in their own ratings in the pre-scoring so that the coefficient is negative.

In order to be easy in taking the conclusion, all various results were taken the average and calculated in Cronbach alpha coefficient. Based on the reliability interpretation, it produced Very High consistency as it got .924 of intraclass correlation in SPSS 23. This value was more than 0.7 as the standard of Cronbach alpha coefficient in deciding the reliability. Simply, it was proved that English teachers of Al-Amin Islamic Boarding School Mojokerto had good reliability. In addition, paired t-test result as the significant calculation of inferential statistic also qualified the rules of rejecting the null hypothesis. The rules are: 1) t-test was more than t-table; $3.541 > 2.776$ and 2) Sig. = 0.02 was less than 0.05 as the level of significant. It means that the result was the real score, not incidentally. Even there were two raters got inconsistent or unreliable in all categories but it did not give any impact to the calculation which proven that the intra-rater reliability of English teachers at Al-Amin Islamic Boarding School Mojokerto was internally consistent. It can be said that the inconsistent scores gotten happened by chance with many exceptions from the raters' self that can be investigated in the next research.

This result was very different from the intra-rater reliability result in the journal from Viphavee Vongpuvimitch entitled *Classroom Writing Teachers'*

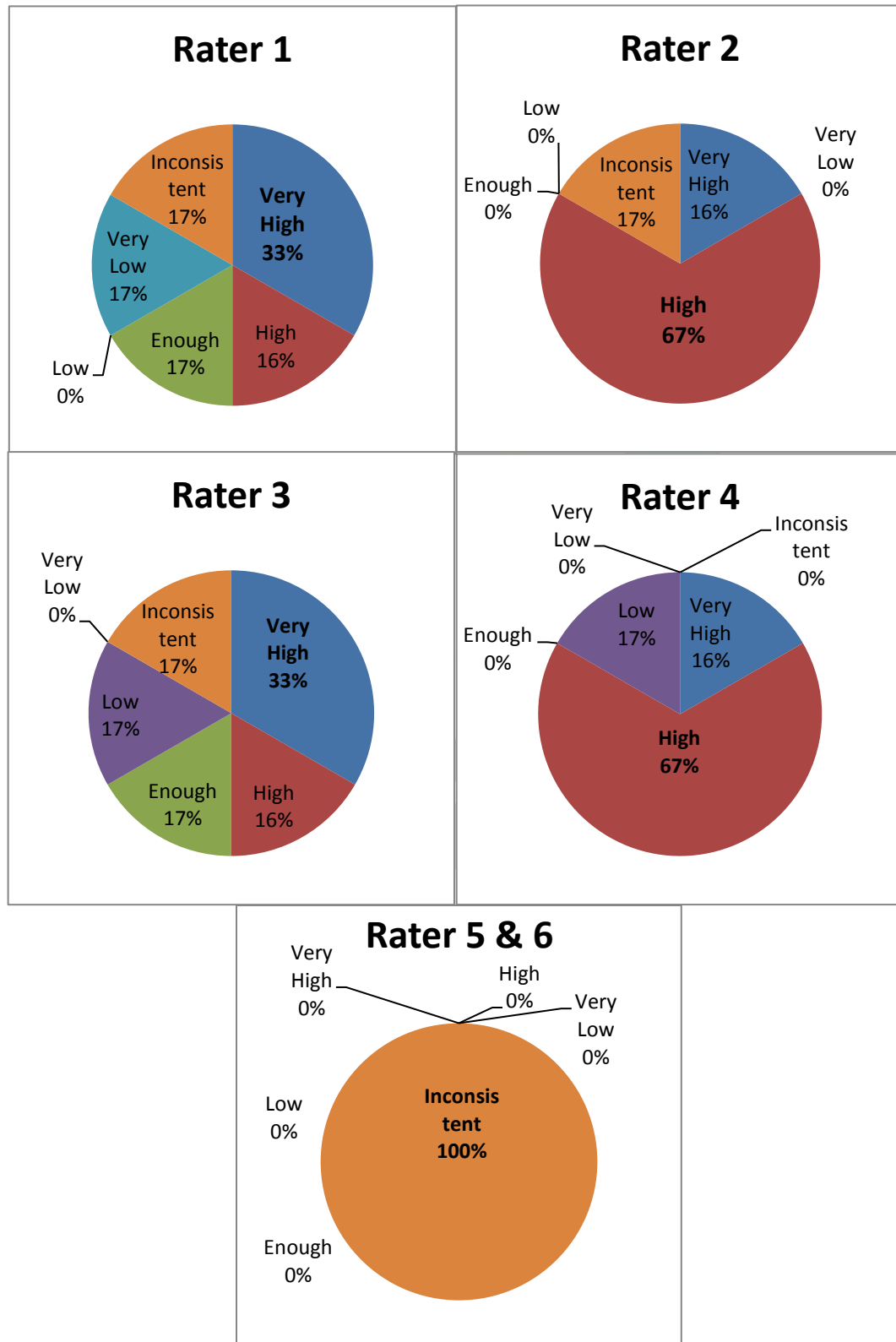
from this study. This journal used holistic whereas the study used analytic. Many factors that make the result of different scoring tool in two studies were different. Therefore, it was needed to investigate deeper in the next research whether the tool used influence the result or not.

Although this study used rubric for the assessment, it was different from journal of *Reliability and Validity of Rubrics for Assessment through Writing* by Ali Reza Rezaei from California State University and Michael Lovorn from The University of Alabama, USA² that showed the rubric impact in improving raters' reliability and validity. As the newest study about this topic in Indonesia, especially in Sunan Ampel State Islamic University, this research did not investigate it too far. The function of rubric was only to help teachers or raters having the same and specific criteria in assessing writing.

To give detail information about reliability level reached by each teacher, the study presented in diagrams what level they were achieved and its percentage.

² Ali Reza Rezaei – Michael Lovorn. “Assessing Writing”. *Reliability and Validity of Rubrics for Assessment Through Writing* Vol. 15, 2010.

4.1 Circle Diagrams of Raters' Reliability Level Percentage



The diagrams above draws clearly that 1st and 3rd teacher got Very High level in 33% of all categories. In addition, 67% High level was achieved by 2nd and 4th rater. Those percents were the highest percentage earned by each rater. It was enough to prove that most of raters in almost category got upper consistency level so that the average result which represents the final conclusion produced very high value. Even if 1st, 2nd and 3rd teacher had reached good value in Very High and High level but they still got negative value in one of categories, especially Mechanic. Fortunately, the negative did not give big impact in the final result and it can be said that the minus score happened incidentally.

In addition, although 4th rater did not get any values in Very High level but she had no minus score in all categories who is admitted as the most stable rater of all. Besides that, 5th and 6th rater had proven well to be the most unreliable rater of all. As it has explained before that it did not influence the average result of this study, both raters might be disturbed by internal or external factors, such as bias, illness, etc., during scoring the essay that can be analyzed deeply in the next research.

As the various result reached by each rater, the researcher tried to map in diagrams what category that have to be given more attention in scoring essay test. It was aimed to know their weakness so that they can improve their reliability.

Diagrams above told that Very High level was reached well in Content and the Total score whereas High level was achieved in Grammar only. It means that Content and Grammar were the most objective categories of all. The other categories got High level as the same as inconsistent level, like Vocabulary. It can be said that each rater had different views in this category that can cause significant different of raters' achievement. Even if half raters got inconsistent level in Organization but another half got High level here, unlike the Vocabulary that produce more varied result. Therefore, the worst categories were Mechanic as 67% of teachers got unreliable level (minus score) here. In other word, Mechanic was not considered well by English teachers in scoring essay test so that more than a half of them got zero reliability.

The teachers did a better job in grading Grammar that was influenced the Total. Almost raters reached Very High and High consistency level as the most objective of the five subscales. In contrast, for about a half raters were very consistent in Organization but the other half were inconsistent. Also even if two teachers reached Very High level in Content but the other two were only stuck in Enough indeed another were inconsistent. The least objective of all were Vocabulary and Mechanic as almost all raters achieved unsatisfactory consistency level.

Even if the diagrams told about the raters' reliability in each category but it did not compare their achievement in each category, means the inter-rater reliability. It just presented the percentage of reliability level in each category to

define what category must be given more attention to improve teachers' reliability so that it would not get lower consistency level.

According to all findings and discussions presented in this chapter, the result of this study showed that the intra-rater reliability of English teachers at Al-Amin Islamic Boarding School Mojokerto was consistent in scoring essay test. The negative result of some raters was admitted as the incident that was happened in chance. From here, it seems that it is very difficult to derive consistent result from raters. Even if they have many experiences in scoring essay test for many years, it is not assure that they have good reliability in scoring subjective test, like essay test.

The result of this study was tentative. The limitation of this study has explained in Chapter I that the result can be changed in another chance as the scoring process was only twice, pre- and post-scoring. However, it is wished that this study allow all language teachers, raters or testers; especially English language, to consider the importance of rater reliability or consistency in writing assessment, like essay evaluation marking. Also, this research is hoped that it can serve as an example for further research in the same topic to eliminate everything that can influence the objectivity of scoring, particularly scoring subjective test.